# CLOUD COST &
# RESOURCE EXPLOSIONS

## Real Interview Scenarios
## & How to Handle Them

A practical guide for Data Engineers to answer

real-world cloud cost and resource management questions with confidence

## By Ankita Gulati

Senior Data Engineer | Interview Mentor

Interview Edition • Practical • Real Scenarios

# Table Of Content

## Scenario 1

# Sudden Spike in AWS EMR Costs

―――

## Problem Statement

An AWS EMR cluster's **monthly cost spikes by 3× overnight**. Investigation shows the spike coincided with an **auto-scaling misconfiguration**. Jobs must continue to meet SLAs, the **budget is constrained**, and the **cluster is shared** across teams.

**Key Details**

- EMR cost increased ~3× overnight
- Auto-scaling recently changed/misconfigured
- Jobs must still complete within SLA
- Shared cluster with multiple workloads
- Budget constraints apply

## Expected vs Actual Behavior

| Expected | Actual |
|---|---|
| Auto-scaling adjusts moderately | Aggressive scale-out |
| Costs align with workload | Costs spike unexpectedly |
| Jobs efficient and stable | Inefficient jobs trigger scaling |
| Budget predictable | Budget exceeded |

This is a **cost efficiency failure**, not a functional outage.

# Why This Problem Is Costly

Because:

- Auto-scaling masks inefficiencies
- Jobs still "succeed," hiding waste
- Shared clusters amplify cost impact

Teams often react by:

- Disabling auto-scaling
- Shrinking the cluster immediately

But **these actions can jeopardize SLAs without fixing the root cause**.

# Clarifying Questions

Before acting, a senior engineer asks:

- Which steps triggered scale-out?
- Did job runtimes or shuffles increase?
- Are there straggler stages or skew?
- Which jobs consumed the added capacity?
- Are scaling policies aligned with workload patterns?

These questions focus on **why scaling happened**, not just stopping it.

# Confirmed Facts & Assumptions

After investigation:

- Auto-scaling reacted to prolonged executor demand
- One or more jobs became inefficient
- Scaling followed policy, not a spike in business demand
- Disabling scaling would risk under-provisioning
- Optimizing jobs would reduce scale-out pressure

**Interpretation:**
The spike is driven by **inefficient resource usage**, not true demand.

# What the System Assumes vs Reality

| Assumption | Reality |
|---|---|
| Scaling means more work | Scaling masks inefficiency |
| Bigger cluster fixes slowness | Bigger cluster increases cost |
| Costs rise with demand | Costs rose due to waste |
| Scaling policy is enough | Jobs still need optimization |

Cloud cost control requires **efficient workloads**, not just controls.

# Root Cause Analysis

## Step 1: Identify Cost Drivers

Observed:

- Specific jobs triggered extended scale-out
- Executors stayed busy due to long stages

**Conclusion:**
Auto-scaling responded correctly—to inefficient workloads.

## Step 2: Inspect Job Behavior

Observed:

- Data skew or excessive shuffles
- Suboptimal partitioning or joins
- Long GC or I/O waits

These inefficiencies **forced the cluster to grow.**

## Step 3: Conceptual Root Cause

The root cause is **job inefficiency amplified by auto-scaling:**

- Scaling responds to demand signals
- Inefficient jobs create false demand
- Costs balloon without improving outcomes

This is a **workload optimization gap**, not a scaling bug.

## Step 4 : Wrong Approach vs Right Approach

**Wrong Approach**

- Disable auto-scaling blindly
- Downsize the cluster immediately
- Ignore the spike temporarily

**Right Approach**

- Analyze job logs and Spark UI
- Optimize inefficient stages
- Tune scaling policies after optimization

Senior engineers **optimize first, restrict later.**

## Step 5 : Validation of Root Cause

To confirm:

- Optimize identified jobs
- Re-run with auto-scaling enabled
- Observe reduced scale-out and lower costs

**Outcome:**
 Jobs meet SLAs with significantly lower spend.

## Step 6 : Corrective Actions

- Profile jobs triggering scale-out
- Fix skew, shuffles, and partitioning
- Tune executor sizing and parallelism
- Set sensible auto-scaling thresholds
- Monitor cost per job, not just cluster cost

These steps align **performance with budget.**

## Step 7 : Result After Fix

| Before Fix | After Fix |
|---|---|
| 3× cost spike | Predictable spend |
| Scaling hides inefficiency | Scaling reflects real demand |
| Budget risk | Budget control |
| Reactive mitigation | Proactive optimization |

## Final Resolution

- **Root Cause:** Inefficient jobs triggering excessive auto-scaling
- **Fix Applied:** Job optimization and scaling policy tuning

## Key Learnings

- Auto-scaling amplifies inefficiencies
- Cost issues often start in job design
- Disabling scaling is not optimization
- Monitor cost at job and stage level

## Core Principle Reinforced

**Cloud cost problems are usually workload problems in disguise.**

■ ■ ■

## Scenario 2

# Long-Running ETL Jobs Cause Cloud Resource Overuse

## Problem Statement

Critical ETL jobs begin running **significantly longer than expected,** consuming **nearly double the cluster resources** and pushing **cloud costs beyond budget.** While jobs still complete within the **2-hour SLA,** the **cost impact is unsustainable.**

**Key Details**

- ETL runtimes increased unexpectedly
- Cluster resource usage doubled
- Cloud budget exceeded
- Jobs are business-critical
- SLA: 2 hours

## Expected vs Actual Behavior

| Expected | Actual |
|---|---|
| Jobs complete within planned runtime | Jobs run much longer |
| Predictable resource usage | Excessive cluster consumption |
| Costs aligned with workload | Cloud spend spikes |
| SLA met efficiently | SLA met but at high cost |

This is a **performance inefficiency and cost optimization issue,** not a functional failure.

# Why This Problem Is Dangerous

Because:

- Jobs still "succeed" technically
- SLA compliance hides inefficiency
- Costs accumulate silently over time

Teams often react by:

- Killing long-running jobs
- Increasing cluster size

But **both actions either break data pipelines or increase costs further.**

# Clarifying Questions

Before acting, a senior engineer asks:

- Which stages or transformations are slow?
- Has data volume or distribution changed?
- Are there skewed joins or shuffles?
- Is GC time or I/O wait unusually high?
- Which jobs contribute most to cost?

These questions focus on **why jobs are slow**, not just stopping them.

# Confirmed Facts & Assumptions

After investigation:

- Job logic has inefficiencies (skew, shuffles, poor partitioning)
- No business requirement changed
- Scaling the cluster would only mask inefficiency
- Killing jobs would break downstream reporting
- Profiling can pinpoint expensive stages

**Interpretation:**
 This is a **job optimization problem**, not a capacity problem.

# What the System Assumes vs Reality

| Assumption | Reality |
|---|---|
| Jobs scale linearly | Inefficiencies grow non-linearly |
| SLA compliance means healthy jobs | Cost efficiency is missing |
| Bigger clusters fix slowness | Bigger clusters increase cost |
| Failures drive cost | Inefficiency drives cost |

Cloud bills grow fastest when **inefficient jobs run longer,** not when they fail.

# Root Cause Analysis

## Step 1: Profile Job Execution

Observed:

- Long-running stages dominate runtime
- Executors busy with skewed partitions or heavy shuffles

**Conclusion:**
Specific transformations are inefficient.

## Step 2: Analyze Resource Utilization

Observed:

- CPU and memory underutilized in parts
- Excessive data movement and serialization

This confirms **logic-level inefficiency,** not lack of resources.

## Step 3: Conceptual Root Cause

The root cause is **unoptimized ETL logic:**

- Inefficient joins or aggregations
- Poor partitioning strategy
- Lack of profiling over time

This is a **performance debt issue.**

## Step 4 : Wrong Approach vs Right Approach

**Wrong Approach**

- Kill long-running jobs
- Increase cluster size blindly
- Ignore cost overruns

**Right Approach**

- Profile job execution (Spark UI / query plans)
- Optimize joins, partitions, and transformations
- Reduce runtime before scaling

Senior engineers **optimize before they scale.**

## Step 5 : Validation of Root Cause

To confirm:

- Optimize identified slow stages
- Re-run job with same cluster
- Compare runtime and cost

**Outcome:**
Jobs complete faster with significantly lower resource usage.

## Step 6 : Corrective Actions

- Profile jobs regularly
- Fix skewed joins and repartition data
- Optimize transformations and serialization
- Monitor runtime trends per job
- Track cost per job, not just SLA

These steps control both **performance and spend.**

## Step 7 : Result After Fix

| Before Fix | After Fix |
|---|---|
| Long runtimes | Optimized runtimes |
| High cloud cost | Controlled spend |
| Reactive scaling | Proactive optimization |
| Budget risk | Budget stability |

## Final Resolution

- **Root Cause:** Inefficient ETL logic causing extended runtimes
- **Fix Applied:** Job profiling and targeted optimization

## Key Learnings

- Long-running jobs silently drain cloud budgets
- SLA success does not imply efficiency
- Profiling reveals cost drivers
- Scaling is a last resort, not a first fix

## Core Principle Reinforced

**If your jobs run longer than necessary, the cloud will charge you for every wasted minute.**

■ ■ ■

# Scenario 3

# Storage Costs Skyrocket Due to Mismanaged Retention Policy

## Problem Statement

Cloud storage costs (Amazon S3) **increase sharply** because **raw data is retained indefinitely** with no cleanup or tiering strategy. While pipelines continue to meet the **1-hour SLA,** the **cloud budget is under serious pressure,** and **data compliance requirements must still be respected**.

**Key Details**

- Raw data retained without expiry
- S3 storage cost increasing rapidly
- Compliance requires controlled retention
- Budget constraints exist
- SLA: 1 hour

## Expected vs Actual Behavior

| Expected | Actual |
|---|---|
| Old data archived or expired | Data retained indefinitely |
| Storage cost grows predictably | Storage cost explodes |
| Compliance respected | Compliance unclear |
| Budget under control | Budget exceeded |

This is a **data governance and cost management failure**, not a pipeline performance issue.

# Why This Problem Is Dangerous

Because:

- Pipelines keep working normally
- Costs grow silently month over month
- Manual cleanup is error-prone and risky

Teams often react by:

- Deleting old data manually
- Ignoring storage growth

But **manual deletion risks compliance violations,** and ignoring the issue guarantees runaway costs.

# Clarifying Questions

Before acting, a senior engineer asks:

- What is the legal/compliance retention period?
- Which data must remain accessible vs archived?
- How frequently is raw data accessed?
- Can storage be tiered automatically?
- Are lifecycle rules already partially defined?

These questions focus on **policy-driven automation,** not ad-hoc cleanup.

# Confirmed Facts & Assumptions

After investigation:

- No lifecycle policies are configured
- Raw data is rarely accessed after a few weeks
- Compliance allows retention with tiering
- Manual deletion is unsafe
- Automated lifecycle rules are supported

**Interpretation:**
This is a **missing lifecycle and retention strategy**, not a storage service limitation.

# What the System Assumes vs Reality

| Assumption | Reality |
|---|---|
| Storage is cheap | Storage grows endlessly |
| Data might be useful someday | Most raw data is cold |
| Manual cleanup works | Manual cleanup doesn't scale |
| Compliance means "keep forever" | Compliance means "retain correctly" |

Cloud storage must be **actively governed**, not passively accumulated.

# Root Cause Analysis

## Step 1: Identify Cost Drivers

Observed:

- Old raw data accumulating daily
- No expiration or transition rules

**Conclusion:**
Unlimited retention is the primary cost driver.

## Step 2: Evaluate Access Patterns

Observed:

- Recent data accessed frequently
- Older data rarely accessed

This supports **tiered storage** instead of deletion.

## Step 3: Conceptual Root Cause

The root cause is **absence of automated retention policies:**

- No lifecycle rules
- No transition to cheaper storage
- Cost grows linearly with time

This is a **governance gap,** not a technical limitation.

## Step 4 : Wrong Approach vs Right Approach

**Wrong Approach**

- Delete old data manually
- Ignore storage growth
- Retain everything indefinitely

**Right Approach**

- Implement S3 lifecycle policies
- Transition data to cheaper tiers (IA, Glacier)
- Expire data based on compliance rules

Senior engineers **automate governance,** not cleanup.

## Step 5 : Validation of Root Cause

To confirm:

- Enable lifecycle rules
- Monitor storage growth and cost
- Validate data availability and compliance

**Outcome:**
Storage growth stabilizes and costs drop predictably.

## Step 6 : Corrective Actions

- Define retention periods per dataset
- Configure lifecycle policies
- Transition cold data to cheaper storage
- Expire data only when compliance allows
- Monitor storage growth trends

These steps enforce **cost control with compliance safety.**

## Step 7 : Result After Fix

| Before Fix | After Fix |
|---|---|
| Exploding storage costs | Controlled storage growth |
| Manual cleanup risk | Automated governance |
| Compliance uncertainty | Compliance assured |
| Budget instability | Budget predictability |

# Final Resolution

- **Root Cause:** Indefinite data retention without lifecycle policies
- **Fix Applied:** Automated lifecycle and tiered storage policies

# Key Learnings

- Storage costs grow silently
- Compliance ≠ infinite retention
- Automation is safer than manual deletion
- Tiered storage is a cost optimization tool

# Core Principle Reinforced

**If data retention is unmanaged, storage costs will eventually manage your budget—for you.**

▪ ▪ ▪

# Scenario 4

# Cloud Functions Trigger Excessively, Driving Up Costs

## Problem Statement

Serverless cloud functions (AWS Lambda) begin triggering **far more frequently than expected,** causing a **sharp increase in Lambda costs.** While the functions still meet the **sub-second SLA**, a **budget alert is triggered,** and **multiple event sources** are configured.

**Key Details**

- Lambda invocation count spikes unexpectedly
- Multiple triggers configured (events, schedules, streams)
- SLA: sub-second response time
- Budget alerts triggered
- Functions still working correctly

## Expected vs Actual Behavior

| Expected | Actual |
|----------|--------|
| Functions trigger only when needed | Functions trigger excessively |
| Predictable invocation count | Invocation spikes |
| Costs proportional to usage | Costs escalate rapidly |
| SLA met within budget | SLA met but budget breached |

This is a **cost efficiency and event design issue**, not a performance failure.

# Why This Problem Is Dangerous

Because:

- Serverless hides infrastructure complexity
- Small per-invocation costs add up fast
- Functions "working fine" masks over-invocation

Teams often react by:

- Increasing function memory
- Ignoring small invocation spikes

But **scaling resources does not reduce invocation count—it often increases cost further.**

# Clarifying Questions

Before acting, a senior engineer asks:

- Which triggers are firing the function?
- Are there duplicate or overlapping event sources?
- Are retries or error loops causing re-invocations?
- Is the function idempotent?
- Can triggers be filtered or consolidated?

These questions focus on **event design,** not compute sizing.

# Confirmed Facts & Assumptions

After investigation:

- Multiple triggers fire the same function
- Some triggers are redundant or too broad
- Function logic itself is efficient
- Increasing memory would raise per-invocation cost
- Reducing triggers lowers invocation volume

**Interpretation:**
 This is a **trigger configuration problem**, not a Lambda performance issue.

# What the System Assumes vs Reality

| Assumption | Reality |
|---|---|
| Serverless scales cheaply | Invocations scale costs |
| Memory tuning fixes cost | Invocation count dominates |
| More triggers = better coverage | Redundant triggers waste money |
| Budget alerts catch issues early | Costs already incurred |

Serverless cost control depends on **event discipline,** not compute tuning.

# Root Cause Analysis

## Step 1: Identify Invocation Sources

Observed:

- Multiple event rules invoking the same function
- Some triggers firing more often than intended

**Conclusion:**
Invocation explosion originates from trigger design.

## Step 2: Evaluate Trigger Necessity

Observed:

- Some triggers overlap in purpose
- Others lack proper filtering

This confirms **uncontrolled event fan-out.**

## Step 3: Conceptual Root Cause

The root cause is **poor trigger governance:**

- No ownership of event sources
- Redundant or overly broad triggers
- Invocation count not monitored

This is a **serverless cost governance gap.**

## Step 4 : Wrong Approach vs Right Approach

**Wrong Approach**

- Increase function memory
- Ignore invocation spikes
- Delete the function

**Right Approach**

- Reduce and consolidate triggers
- Add filtering at the event source
- Monitor invocation metrics

Senior engineers **control events before scaling functions.**

## Step 5 : Validation of Root Cause

To confirm:

- Disable redundant triggers
- Observe invocation count drop
- Monitor cost trend

**Outcome:**
Invocation volume and Lambda cost reduce immediately.

## Step 6 : Corrective Actions

- Audit all triggers per function
- Remove redundant event sources
- Add filters to narrow triggers
- Monitor invocation rate, not just duration
- Set budgets and alerts per function

These steps keep **serverless costs predictable.**

## Step 7 : Result After Fix

| Before Fix | After Fix |
|---|---|
| Excessive invocations | Controlled invocation rate |
| Cost spike | Stable cost |
| Hidden inefficiency | Transparent event flow |
| Budget alerts | Budget compliance |

# Final Resolution

- **Root Cause:** Redundant and overly broad triggers causing excessive invocations
- **Fix Applied:** Reduced and consolidated event triggers

# Key Learnings

- Serverless cost is invocation-driven
- Memory tuning doesn't fix over-invocation
- Trigger design matters more than function code
- Monitor events, not just execution time

# Core Principle Reinforced

**In serverless systems, uncontrolled triggers—not slow code—are the fastest way to burn budget.**

■ ■ ■

# Scenario 5

# Under-Provisioned EMR Cluster Triggers Costly Over-Scaling

## Problem Statement

An EMR cluster is **initially under-provisioned** for its workload. As jobs start running, **auto-scaling rapidly spins up many additional nodes**, causing a **sharp increase in cloud costs.** Jobs must still meet the **1-hour SLA,** and a **budget alert has already fired.**

**Key Details**

- EMR base cluster size too small
- Auto-scaling aggressively adds nodes
- Jobs are business-critical
- Budget alerts triggered
- SLA: 1 hour

## Expected vs Actual Behavior

| Expected | Actual |
|----------|--------|
| Baseline cluster handles workload | Baseline cluster overwhelmed |
| Auto-scaling adds minimal capacity | Excessive nodes spun up |
| Costs scale predictably | Costs spike unexpectedly |
| SLA met efficiently | SLA met at high cost |

This is a **capacity planning and cost efficiency issue,** not a job failure.

# Why This Problem Is Common

Because:

- Teams try to save costs by starting small
- Auto-scaling is assumed to "handle everything"
- Baseline sizing is often underestimated

But **auto-scaling reacts to pressure—it doesn't fix poor sizing.**

# Clarifying Questions

Before acting, a senior engineer asks:

- What is the typical vs peak workload?
- Which stages trigger scale-out?
- Is scale-out happening immediately or gradually?
- Are executors starved early in the job?
- What is the minimum stable cluster size?

These questions focus on **right-sizing**, not disabling safeguards.

# Confirmed Facts & Assumptions

After investigation:

- Initial node count is too low for job startup
- Auto-scaling reacts correctly to executor starvation
- Jobs stabilize only after large scale-out
- Disabling auto-scaling risks SLA breaches
- Proper baseline sizing would reduce scaling events

**Interpretation:**
This is an **under-sized baseline cluster problem**.

# What the System Assumes vs Reality

| Assumption | Reality |
|---|---|
| Smaller base cluster saves money | It triggers costly scaling |
| Auto-scaling fixes sizing | Auto-scaling amplifies mistakes |
| Scaling means higher demand | Scaling hides poor planning |
| Budget alerts catch issues early | Costs already incurred |

Cloud cost control starts with **correct baseline sizing.**

# Root Cause Analysis

## Step 1: Observe Scaling Pattern

Observed:

- Immediate scale-out at job start
- Many nodes added quickly

**Conclusion:**
 The cluster cannot handle normal workload at baseline.

## Step 2: Evaluate Baseline Capacity

Observed:

- Executors starved early
- Shuffle and GC pressure high initially

This confirms **baseline under-provisioning**.

## Step 3: Conceptual Root Cause

The root cause is **incorrect cluster sizing:**

- Base capacity too small
- Auto-scaling compensates aggressively
- Costs rise without improving efficiency

This is a **capacity planning gap**, not an auto-scaling bug.

## Step 4 : Wrong Approach vs Right Approach

**Wrong Approach**

- Disable auto-scaling
- Retry jobs repeatedly
- Ignore budget alerts

**Right Approach**

- Resize baseline cluster appropriately
- Let auto-scaling handle only true spikes
- Balance baseline + scaling strategy

Senior engineers **size for normal load, scale for peaks.**

## Step 5 : Validation of Root Cause

To confirm:

- Increase baseline node count
- Re-run jobs with auto-scaling enabled
- Observe reduced scale-out and stable costs

**Outcome:**
Jobs meet SLA with far fewer additional nodes.

## Step 6 : Corrective Actions

- Right-size baseline EMR cluster
- Use historical job metrics for sizing
- Tune auto-scaling thresholds
- Monitor scale-out frequency
- Track cost per job run

These steps align **performance, stability, and cost.**

## Step 7 : Result After Fix

| **Before Fix** | **After Fix** |
|---|---|
| Excessive scale-out | Controlled scaling |
| Cost spikes | Predictable spend |
| Reactive mitigation | Proactive sizing |
| Budget risk | Budget stability |

## Final Resolution

- **Root Cause:** Under-provisioned baseline cluster
- **Fix Applied:** Resized cluster to match expected workload

## Key Learnings

- Auto-scaling is not a sizing substitute
- Start with the right baseline capacity
- Under-sizing can cost more than over-sizing
- Cost optimization begins with planning

## Core Principle Reinforced

**If your base cluster is too small, auto-scaling will make sure your bill isn't.**

# Scenario 6

# Inefficient Spark Jobs Consume Excess Cloud Resources

## Problem Statement

Critical Spark jobs are configured with **excessive executors and memory,** leading to **unnecessarily high cloud costs.** Although jobs still complete within the **2-hour SLA**, the **shared cluster** experiences resource pressure and the cloud bill continues to rise.

**Key Details**

- Spark jobs over-allocate executors and memory
- Cluster is shared across teams
- Jobs are business-critical
- SLA: 2 hours

## Expected vs Actual Behavior

| Expected | Actual |
|---|---|
| Executors sized to workload | Too many executors allocated |
| Memory used efficiently | Large memory underutilized |
| Costs proportional to workload | Cloud bills inflated |
| SLA met efficiently | SLA met at high cost |

This is a **configuration and efficiency issue,** not a functional or SLA failure.

# Why This Problem Is Costly

Because:

- Spark will happily use all allocated resources
- Over-provisioning looks "safe" under SLA pressure
- Shared clusters amplify waste across teams

Teams often:

- Increase cluster size further
- Retry jobs without changes

But **scaling inefficient jobs only increases waste.**

# Clarifying Questions

Before acting, a senior engineer asks:

- Are executors fully utilized?
- Is memory usage close to allocation?
- Are there many idle executors?
- Is GC time high or low?
- Does the job scale linearly with executors?

These questions focus on **right-sizing,** not brute-force scaling.

# Confirmed Facts & Assumptions

After investigation:

- Executors are frequently idle
- Memory usage is far below allocation
- Job runtime does not improve with more executors
- Cluster contention increases for other jobs
- Proper tuning can reduce resource usage safely

**Interpretation:**
 This is a **Spark configuration inefficiency**, not a workload growth issue.

# What the System Assumes vs Reality

| Assumption | Reality |
|---|---|
| More executors = faster jobs | Diminishing returns |
| More memory = safer execution | Memory sits unused |
| Over-provisioning is harmless | Over-provisioning is costly |
| SLA success = optimal setup | SLA hides inefficiency |

Spark performance depends on **efficient configuration,** not maximum allocation.

# Root Cause Analysis

## Step 1: Analyze Executor Utilization

Observed:

- Many executors idle or lightly loaded
- CPU utilization low

**Conclusion:**
 Executor count exceeds what the job can use.

## Step 2: Analyze Memory Usage

Observed:

- Allocated memory far exceeds actual usage
- Minimal GC pressure

This confirms **memory over-allocation.**

## Step 3: Conceptual Root Cause

The root cause is **poor Spark executor configuration:**

- Executor count and memory not aligned with workload
- Configuration copied without tuning
- Cost grows without performance benefit

This is a **performance tuning gap,** not an infrastructure problem.

## Step 4 : Wrong Approach vs Right Approach

**Wrong Approach**

- Increase cluster size
- Retry jobs repeatedly
- Ignore resource waste

**Right Approach**

- Tune executor count and memory
- Match configuration to workload
- Optimize before scaling

Senior engineers **right-size first, then scale if needed.**

## Step 5 : Validation of Root Cause

To confirm:

- Reduce executor count and memory
- Re-run job
- Compare runtime and cost

**Outcome:**
 Job meets SLA with significantly lower resource usage.

## Step 6 : Corrective Actions

- Tune **executor.instances, executor.memory, executor.cores**
- Avoid copy-paste configurations
- Monitor executor utilization metrics
- Set sensible defaults for shared clusters
- Track cost per job execution

These steps control **both performance and spend.**

## Step 7 : Result After Fix

| Before Fix | After Fix |
|---|---|
| Excessive resource usage | Efficient resource usage |
| High cloud cost | Reduced cloud cost |
| Cluster contention | Stable shared cluster |
| Hidden waste | Measured efficiency |

## Final Resolution

- **Root Cause:** Over-provisioned Spark executors and memory
- **Fix Applied:** Optimized executor configuration

## Key Learnings

- Spark uses what you give it—even if it doesn't need it
- Over-provisioning is a hidden cost multiplier
- SLA success does not equal efficiency
- Executor tuning is a core data engineering skill

## Core Principle Reinforced

**In Spark, unused executors don't fail jobs—but they silently drain your cloud budget.**

▪ ▪ ▪

# Scenario 7

# Idle Cloud Resources Left Running Inflate Costs

## Problem Statement

Non-production resources such as **staging clusters, development environments, and test databases** are left running over weekends and off-hours. Although these resources are **non-critical for SLAs,** they continue to consume cloud capacity, **breaching budget limits** across shared accounts.

**Key Details**

- Staging, dev, and test resources run continuously
- Workloads are non-critical
- Multiple teams share cloud resources
- Budget limits breached

## Expected vs Actual Behavior

| Expected | Actual |
|---|---|
| Non-prod resources shut down when idle | Resources run 24/7 |
| Costs reflect active usage | Costs accrue during inactivity |
| Budgets predictable | Budgets exceeded |
| Teams manage usage responsibly | Idle resources forgotten |

This is a **governance and automation failure,** not a performance issue.

# Why This Problem Is Common

Because:

- Engineers forget to stop resources
- Ownership across teams is unclear
- Manual processes don't scale

Teams often rely on:

- Slack reminders
- Periodic cleanup

But **humans are unreliable cost-control mechanisms.**

# Clarifying Questions

Before acting, a senior engineer asks:

- Which resources are non-production?
- When are they actually needed?
- Can they be safely stopped without impact?
- Who owns each resource?
- Can shutdowns be automated?

These questions focus on **policy-driven automation,** not reminders.

# Confirmed Facts & Assumptions

After investigation:

- Many resources idle for long periods
- No automated shutdown in place
- Manual reminders fail repeatedly
- Auto-shutdown is technically feasible
- Deleting resources is too destructive

**Interpretation:**
This is a **lack of automated resource governance**.

## What the System Assumes vs Reality

| Assumption | Reality |
|---|---|
| Teams will stop unused resources | They often forget |
| Non-prod cost is negligible | Non-prod costs add up |
| Manual controls are enough | Automation is required |
| Budget alerts prevent waste | Alerts arrive too late |

Cost control requires **defaults that save money,** not rely on behavior.

## Root Cause Analysis

### Step 1: Identify Idle Resources

Observed:

- Clusters and databases idle overnight and on weekends
- No activity during off-hours

**Conclusion:**
Resources run without business need.

### Step 2: Evaluate Shutdown Risk

Observed:

- No SLA impact for dev/staging
- Resources can be restarted safely

This confirms **auto-shutdown is low-risk.**

## Step 3: Conceptual Root Cause

The root cause is **missing lifecycle automation:**

- No start/stop schedules
- No ownership enforcement
- Costs accumulate silently

This is a **cloud governance gap.**

## Step 4 : Wrong Approach vs Right Approach

**Wrong Approach**

- Ask teams to stop resources manually
- Ignore idle usage
- Delete resources outright

**Right Approach**

- Implement auto-shutdown policies
- Use schedules and tags
- Restart resources on demand

Senior engineers **design systems that fail safe for cost.**

## Step 5 : Validation of Root Cause

To confirm:

- Enable auto-shutdown for non-prod resources
- Monitor usage over weekends
- Compare cloud spend before vs after

**Outcome:**
Costs drop immediately without impacting productivity.

## Step 6 : Corrective Actions

- Tag non-production resources
- Implement scheduled shutdown/startup
- Enforce policies via IaC
- Alert on resources running outside schedules
- Track non-prod spend separately

These steps automate **cost hygiene.**

## Step 7 : Result After Fix

| Before Fix | After Fix |
|---|---|
| Idle resources running | Resources shut down automatically |
| Budget breaches | Budget compliance |
| Manual reminders | Automated enforcement |
| Hidden waste | Visible savings |

## Final Resolution

- **Root Cause:** Idle non-production resources left running
- **Fix Applied:** Automated shutdown policies

## Key Learnings

- Non-prod environments drive real costs
- Manual cost controls don't scale
- Automation beats reminders
- Idle time is paid time in the cloud

## Core Principle Reinforced

**If a resource doesn't shut itself down, it will eventually shut down your budget.**

■ ■ ■

# Scenario 8

# Over-Provisioned Storage for Frequently Accessed Data Drives High Costs

## Problem Statement

Frequently queried tables are stored entirely in **high-performance (hot) storage**, even though a large portion of the data is **rarely accessed.** While **sub-second query SLAs** are met and dashboards perform well, a **budget alert has been triggered** due to rising storage costs.

**Key Details**

- Hot storage used for all data
- Large portion of data rarely accessed
- Downstream dashboards are business-critical
- SLA: sub-second query performance
- Budget alert triggered

## Expected vs Actual Behavior

| Expected | Actual |
|---|---|
| Hot storage used only for active data | All data stored in hot tier |
| Cold data moved to cheaper tiers | Cold data remains expensive |
| Performance and cost balanced | Performance achieved at high cost |
| Budget predictable | Budget exceeded |

This is a **storage tiering and cost efficiency issue**, not a performance failure.

# Why This Problem Is Common

Because:

- Teams default to "keep everything hot" for safety
- Access patterns are rarely reviewed
- Performance SLAs overshadow cost considerations

But **hot storage pricing assumes frequent access—not indefinite retention.**

# Clarifying Questions

Before acting, a senior engineer asks:

- What percentage of data is queried frequently
- Can historical data tolerate higher latency?
- Are access patterns seasonal or consistent?
- Does the storage engine support tiering?
- Can dashboards remain fast with partial hot data?

These questions focus on **data temperature,** not deletion.

# Confirmed Facts & Assumptions

After investigation:

- Majority of queries hit recent data
- Historical partitions rarely accessed
- Cold data can tolerate slower reads
- Hot storage cost dominates monthly spend
- Tiered storage is supported

**Interpretation:**
This is a **misalignment between access patterns and storage tier**.

# What the System Assumes vs Reality

| Assumption | Reality |
|---|---|
| Fast queries require all data hot | Only recent data needs hot storage |
| Storage cost is fixed | Cost scales with retained volume |
| Deleting data is the only savings | Tiering preserves data + saves cost |
| Performance and cost conflict | They can be balanced |

Cloud storage is cheapest when **data temperature matches access frequency.**

# Root Cause Analysis

## Step 1: Analyze Access Patterns

Observed:

- Recent partitions queried frequently
- Older partitions almost never touched

**Conclusion:**
Not all data requires hot storage.

## Step 2: Evaluate Storage Configuration

Observed:

- No tiering or lifecycle rules
- All partitions treated equally

This confirms **over-provisioned hot storage.**

## Step 3: Conceptual Root Cause

The root cause is **lack of tiered storage strategy:**

- Access patterns ignored
- Hot storage used by default
- Costs grow unnecessarily

This is a **storage governance gap.**

## Step 4 : Wrong Approach vs Right Approach

**Wrong Approach**

- Keep all data in hot storage
- Ignore budget alerts
- Delete historical data

**Right Approach**

- Move rarely accessed data to cheaper tiers
- Keep hot storage for active datasets
- Balance performance and cost

Senior engineers **optimize placement, not availability.**

## Step 5 : Validation of Root Cause

To confirm:

- Move cold partitions to cheaper storage
- Measure query latency for dashboards
- Monitor storage cost trend

**Outcome:**
Performance remains within SLA, costs drop significantly.

## Step 6 : Corrective Actions

- Analyze query access patterns regularly
- Define hot vs warm vs cold data
- Implement lifecycle or tiering policies
- Monitor cost per storage tier
- Review tiering as data grows

These steps align **performance requirements with cost control.**

## Step 7 : Result After Fix

| Before Fix | After Fix |
|---|---|
| All data in hot storage | Tiered storage applied |
| High storage cost | Reduced monthly spend |
| Budget alerts | Budget stability |
| Inefficient usage | Optimized storage layout |

## Final Resolution

- **Root Cause:** Hot storage used for rarely accessed data
- **Fix Applied:** Moved cold data to cheaper storage tiers

## Key Learnings

- Storage costs grow silently
- Not all data deserves hot storage
- Access patterns should drive storage decisions
- Tiered storage preserves both data and budget

## Core Principle Reinforced

**If cold data lives in hot storage, your cloud bill will always feel hot.**

■ ■ ■

# Scenario 9

# Data Transfer Costs Balloon Due to Cross-Region ETL

## Problem Statement

ETL jobs begin **moving large datasets across regions or cloud providers,** causing **unexpected spikes in data transfer costs.** While pipelines still meet the **1-hour SLA,** the workload is **budget-sensitive,** and data replication remains a business requirement.

**Key Details**

- Large datasets transferred cross-region / cross-cloud
- Transfer costs spike unexpectedly
- Data replication still required
- Budget sensitivity high
- SLA: 1 hour

## Expected vs Actual Behavior

| Expected | Actual |
|---|---|
| Data transfer minimized | Large volumes moved unnecessarily |
| Transfer cost predictable | Transfer costs balloon |
| SLA met within budget | SLA met but budget exceeded |
| Replication efficient | Replication inefficient |

This is a **data placement and movement efficiency issue,** not a pipeline failure.

## Why This Problem Is Often Missed

Because:

- Transfer costs are less visible than compute
- Pipelines continue to function normally
- Costs show up only in billing reports

Teams often:

- Retry transfers
- Ignore small spikes

But **data movement is one of the most expensive cloud operations at scale.**

## Clarifying Questions

Before acting, a senior engineer asks:

- Why is data moving across regions?
- Can workloads be co-located with data?
- Is full dataset transfer required every run?
- Can data be compressed before transfer?
- Are transfers incremental or full refreshes?

These questions focus on **reducing bytes moved,** not speeding up transfers.

## Confirmed Facts & Assumptions

After investigation:

- ETL transfers full datasets repeatedly
- Jobs run in a different region than storage
- Compression is not enabled
- Retry does not reduce transfer volume
- Co-location or compression reduces cost

**Interpretation:**
 This is a **data movement design issue**, not a network failure.

## What the System Assumes vs Reality

| Assumption | Reality |
|---|---|
| Transfer cost is minor | Transfer cost is significant |
| Faster transfers solve cost | Volume matters more than speed |
| Replication requires full copy | Incremental copies often sufficient |
| SLA focus is enough | Cost must be designed in |

Cloud cost control starts with **minimizing data movement.**

## Root Cause Analysis

### Step 1: Identify Transfer Pattern

Observed:

- Large volumes moved each run
- Same data transferred repeatedly

**Conclusion:**
Transfer volume, not frequency, drives cost.

### Step 2: Evaluate Placement Strategy

Observed:

- Compute and storage in different regions
- No co-location strategy

This confirms **inefficient workload placement.**

## Step 3: Conceptual Root Cause

The root cause is **poor data locality and transfer optimization:**

- No compression
- No co-location
- No incremental transfer strategy

This is a **cloud architecture design gap.**

## Step 4 : Wrong Approach vs Right Approach

**Wrong Approach**

- Ignore transfer costs
- Retry transfers
- Accept cost as unavoidable

**Right Approach**

- Compress data before transfer
- Co-locate compute with storage
- Transfer only incremental changes

Senior engineers **design for locality, not bandwidth.**

## Step 5 : Validation of Root Cause

To confirm:

- Enable compression
- Move ETL jobs closer to data
- Monitor transfer volume and cost

**Outcome:**
Transfer costs drop significantly without impacting SLA.

## Step 6 : Corrective Actions

- Enable compression for transfers
- Co-locate ETL jobs with storage
- Use incremental or delta-based replication
- Monitor cross-region traffic explicitly
- Track cost per GB transferred

These steps keep **data movement costs predictable.**

## Step 7 : Result After Fix

| Before Fix | After Fix |
|---|---|
| High transfer volume | Reduced transfer volume |
| Ballooning costs | Controlled costs |
| Hidden billing impact | Visible, managed spend |
| Reactive fixes | Proactive design |

# Final Resolution

- **Root Cause:** Large, unnecessary cross-region data transfers
- **Fix Applied:** Compression and improved data locality

# Key Learnings

- Data transfer costs add up quickly
- Bytes moved matter more than speed
- Co-location reduces both cost and complexity
- Replication doesn't mean full copy

# Core Principle Reinforced

**The cheapest data transfer is the one you don't make.**

■ ■ ■

# Scenario 10

# **Spot Instances Terminated Mid-Job, Driving Hidden Costs**

## Problem Statement

Critical ETL jobs are executed on **spot instances** to reduce compute costs. Due to **spot price volatility,** instances are frequently terminated mid-execution, triggering **job retries**. Although jobs eventually complete within the **2-hour SLA**, repeated retries **increase overall cloud costs and reduce reliability.**

**Key Details**

- ETL jobs running on spot instances
- Frequent spot interruptions
- Jobs are business-critical
- SLA: 2 hours
- Spot pricing volatile

## Expected vs Actual Behavior

| Expected | Actual |
|---|---|
| Spot instances reduce compute cost | Retries inflate total cost |
| Jobs run uninterrupted | Jobs terminated mid-run |
| SLA met efficiently | SLA met after retries |
| Cost savings realized | Savings offset by failures |

This is a **cost–reliability trade-off failure,** not a capacity issue.

# Why This Problem Is Misleading

Because:

- Spot instances appear cheaper on paper
- Jobs eventually succeed
- SLA compliance hides instability

Teams often assume:

- "Retries are acceptable"
- "Spot is always cheaper"

But **for critical jobs, retries erase cost benefits quickly.**

# Clarifying Questions

Before acting, a senior engineer asks:

- How often are spot instances interrupted?
- What is the retry cost per job?
- Are jobs checkpointed or restart-safe?
- Which jobs are truly SLA-critical?
- Is predictable execution more valuable than marginal savings?

These questions balance **cost vs reliability,** not just pricing.

# Confirmed Facts & Assumptions

After investigation:

- Spot interruptions are frequent during peak hours
- Jobs restart from scratch on termination
- Retry compute cost exceeds on-demand pricing
- SLA pressure increases operational risk
- On-demand instances provide stability

**Interpretation:**
This is **misuse of spot instances for critical workloads**.

# What the System Assumes vs Reality

| Assumption | Reality |
|---|---|
| Spot is always cheaper | Retries add hidden cost |
| SLA success means good choice | Reliability matters |
| Volatility is manageable | Interruptions are frequent |
| Cost optimization is pricing only | Cost includes retries |

Cloud optimization must consider **failure cost,** not just hourly rates.

# Root Cause Analysis

## Step 1: Analyze Failure Pattern

Observed:

- Jobs terminate mid-execution
- Retries restart full workload

**Conclusion:**
Spot interruptions directly cause waste.

## Step 2: Compare Cost Models

Observed:

- Multiple retries on spot > single on-demand run
- Engineering time spent managing failures

This confirms **false economy of spot for critical jobs.**

## Step 3: Conceptual Root Cause

The root cause is **incorrect workload-to-pricing alignment:**

- Spot instances used for SLA-critical jobs
- No tolerance for interruption
- Retry cost ignored in planning

This is a **resource strategy gap.**

## Step 4 : Wrong Approach vs Right Approach

**Wrong Approach**

- Keep retrying on spot
- Ignore interruption frequency
- Split jobs without stability

**Right Approach**

- Move critical jobs to on-demand or reserved instances
- Use spot only for fault-tolerant workloads
- Separate critical and non-critical pipelines

Senior engineers choose **predictability over theoretical savings.**

## Step 5 : Validation of Root Cause

To confirm:

- Run job on on-demand instances
- Compare completion time and cost
- Observe elimination of retries

**Outcome:**
Jobs complete once, costs stabilize, SLA risk removed.

## Step 6 : Corrective Actions

- Classify jobs by criticality
- Use on-demand/reserved for critical ETL
- Use spot for retry-safe or batch workloads
- Add interruption metrics to cost analysis
- Design checkpoints if spot must be used

These steps align **cost strategy with workload reliability.**

## Step 7 : Result After Fix

| Before Fix | After Fix |
|---|---|
| Frequent retries | Single successful run |
| Unpredictable cost | Predictable spend |
| SLA risk | SLA confidence |
| Operational noise | Stable execution |

## Final Resolution

- **Root Cause:** Using spot instances for interruption-intolerant jobs
- **Fix Applied:** Moved critical ETL jobs to on-demand resources

## Key Learnings

- Spot pricing ≠ guaranteed savings
- Retries are a real cost
- Critical jobs need stable infrastructure
- Cost optimization includes reliability

## Core Principle Reinforced

**Cheap compute isn't cheap if you have to run the job twice.**

■ ■ ■