

PERFORMANCE & LATENCY ISSUES

Real Interview Scenarios & How to Handle Them

A practical guide for Data Engineers to answer
real-world performance and latency questions with confidence

By Ankita Gulati

Senior Data Engineer | Interview Mentor

Interview Edition • Practical • Real Scenarios



Table of Content

| | |
|------------------------------------------------------------|-----------|
| Scenario 1 | 9 |
| Spark Job Runtime Gradually Increases Over Time | 9 |
| Problem Statement | 9 |
| Why This Problem Is Tricky | 9 |
| Clarifying Questions | 10 |
| Before proposing fixes, reduce uncertainty: | 10 |
| Confirmed Facts & Assumptions | 10 |
| Root Cause Analysis | 11 |
| Final Resolution | 13 |
| Key Learnings | 13 |
| Core Principle Reinforced | 13 |
| Scenario 2 | 14 |
| SQL Queries Running Slower Than Expected (Redshift) | 14 |
| Problem Statement | 14 |
| Expected vs Actual Behavior | 14 |
| This gap indicates a planning problem, not a logic problem | 14 |
| Confirmed Facts & Assumptions | 15 |
| What the System Thinks vs Reality | 16 |
| Root Cause Analysis | 16 |
| Final Resolution | 18 |
| Key Learnings | 18 |
| Core Principle Reinforced | 18 |
| Scenario 3 | 18 |
| Spark Shuffle Becomes a Bottleneck | 19 |
| Problem Statement | 19 |
| Expected vs Actual Behavior | 19 |
| Why This Problem Is Tricky | 20 |
| Clarifying Questions | 20 |
| Confirmed Facts & Assumptions | 20 |
| Root Cause Analysis | 21 |
| Final Resolution | 23 |
| Key Learnings | 23 |
| Core Principle Reinforced | 23 |

| | |
|-------------------------------------------------------|-----------|
| Scenario 4..... | 23 |
| Streaming Job Slows Down During Peak Traffic..... | 24 |
| Problem Statement..... | 24 |
| Expected vs Actual Behavior..... | 24 |
| Why This Problem Is Misleading..... | 25 |
| Clarifying Questions..... | 25 |
| Confirmed Facts & Assumptions..... | 25 |
| What the System Assumes vs Reality..... | 26 |
| Root Cause Analysis..... | 26 |
| Final Resolution..... | 28 |
| Key Learnings..... | 28 |
| Core Principle Reinforced..... | 28 |
| Scenario 5..... | 29 |
| High Garbage Collection (GC) Time in Spark..... | 29 |
| Problem Statement..... | 29 |
| Expected vs Actual Behavior..... | 29 |
| Why This Problem Is Misleading..... | 30 |
| Clarifying Questions..... | 30 |
| Confirmed Facts & Assumptions..... | 30 |
| What Spark Expects vs Reality..... | 31 |
| Root Cause Analysis..... | 31 |
| Final Resolution..... | 33 |
| Key Learnings..... | 33 |
| Core Principle Reinforced..... | 34 |
| Scenario 6..... | 34 |
| Query Performance Degrades After Column Addition..... | 34 |
| Problem Statement..... | 34 |
| Expected vs Actual Behavior..... | 34 |
| Clarifying Questions..... | 35 |
| Confirmed Facts & Assumptions..... | 35 |
| Root Cause Analysis..... | 36 |
| Final Resolution..... | 38 |
| Key Learnings..... | 38 |
| Core Principle Reinforced..... | 39 |

| | |
|--------------------------------------------|-----------|
| Scenario 7..... | 39 |
| Network Latency Causes Job Delay..... | 39 |
| Problem Statement..... | 39 |
| Expected vs Actual Behavior..... | 39 |
| Why This Problem Is Misleading..... | 40 |
| Clarifying Questions..... | 40 |
| Confirmed Facts & Assumptions..... | 40 |
| What the System Assumes vs Reality..... | 41 |
| Root Cause Analysis..... | 41 |
| Final Resolution..... | 43 |
| Key Learnings..... | 44 |
| Core Principle Reinforced..... | 44 |
| Scenario 8..... | 44 |
| Skewed Join Causes Spark Job Slowness..... | 44 |
| Problem Statement..... | 44 |
| Expected vs Actual Behavior..... | 44 |
| Why This Problem Is Deceptive..... | 45 |
| Clarifying Questions..... | 45 |
| Confirmed Facts & Assumptions..... | 46 |
| What Spark Assumes vs Reality..... | 46 |
| Root Cause Analysis..... | 46 |
| Final Resolution..... | 48 |
| Key Learnings..... | 49 |
| Core Principle Reinforced..... | 49 |
| Scenario 9..... | 49 |
| Slow Downstream Writes to Redshift..... | 49 |
| Problem Statement..... | 49 |
| Expected vs Actual Behavior..... | 50 |
| Why This Problem Is Misleading..... | 50 |
| Clarifying Questions..... | 50 |
| Confirmed Facts & Assumptions..... | 51 |
| What the System Expects vs Reality..... | 51 |
| Root Cause Analysis..... | 51 |
| Final Resolution..... | 54 |
| Key Learnings..... | 54 |
| Core Principle Reinforced..... | 54 |

| | |
|-------------------------------------------------|-----------|
| Scenario 10..... | 54 |
| Spark Job Slows Down After Dataset Growth..... | 54 |
| Problem Statement..... | 54 |
| Expected vs Actual Behavior..... | 55 |
| Why This Problem Is Misleading..... | 55 |
| Clarifying Questions..... | 55 |
| Confirmed Facts & Assumptions..... | 56 |
| What Spark Assumes vs Reality..... | 56 |
| Root Cause Analysis..... | 56 |
| Final Resolution..... | 59 |
| Key Learnings..... | 59 |
| Core Principle Reinforced..... | 59 |
| Scenario 11..... | 59 |
| DAG Task Queueing Causes Latency (Airflow)..... | 59 |
| Problem Statement..... | 59 |
| Expected vs Actual Behavior..... | 60 |
| Why This Problem Is Misleading..... | 60 |
| Clarifying Questions..... | 60 |
| Confirmed Facts & Assumptions..... | 61 |
| What the Scheduler Assumes vs Reality..... | 61 |
| Root Cause Analysis..... | 61 |
| Final Resolution..... | 64 |
| Key Learnings..... | 64 |
| Core Principle Reinforced..... | 64 |
| Scenario 12..... | 64 |
| Slow Join on Large Tables in Production..... | 64 |
| Problem Statement..... | 64 |
| Expected vs Actual Behavior..... | 65 |
| Why This Problem Is Misleading..... | 65 |
| Clarifying Questions..... | 65 |
| Confirmed Facts & Assumptions..... | 66 |
| What Spark Assumes vs Reality..... | 66 |
| Root Cause Analysis..... | 66 |
| Final Resolution..... | 69 |
| Key Learnings..... | 69 |
| Core Principle Reinforced..... | 69 |

| | |
|---------------------------------------------|-----------|
| Scenario 13..... | 69 |
| Slow Job After Code Refactor..... | 69 |
| Problem Statement..... | 69 |
| Expected vs Actual Behavior..... | 70 |
| Why This Problem Is Misleading..... | 70 |
| Clarifying Questions..... | 70 |
| Confirmed Facts & Assumptions..... | 71 |
| What the Developer Intended vs Reality..... | 71 |
| Root Cause Analysis..... | 71 |
| Final Resolution..... | 74 |
| Key Learnings..... | 74 |
| Core Principle Reinforced..... | 74 |
| Scenario 14..... | 74 |
| Slow ETL Due to Excessive Logging..... | 74 |
| Problem Statement..... | 74 |
| Expected vs Actual Behavior..... | 75 |
| Why This Problem Is Misleading..... | 75 |
| Clarifying Questions..... | 76 |
| Confirmed Facts & Assumptions..... | 76 |
| What the System Assumes vs Reality..... | 76 |
| Root Cause Analysis..... | 76 |
| Final Resolution..... | 79 |
| Key Learnings..... | 79 |
| Core Principle Reinforced..... | 79 |
| Scenario 15..... | 79 |
| Slow Job Due to Skewed Aggregations..... | 79 |
| Problem Statement..... | 79 |
| Expected vs Actual Behavior..... | 80 |
| Why This Problem Is Misleading..... | 80 |
| Clarifying Questions..... | 81 |
| Confirmed Facts & Assumptions..... | 81 |
| What Spark Assumes vs Reality..... | 81 |
| Root Cause Analysis..... | 81 |
| Final Resolution..... | 84 |
| Key Learnings..... | 84 |
| Core Principle Reinforced..... | 84 |

Scenario 1

Spark Job Runtime Gradually Increases Over Time

Problem Statement

A daily Spark ETL job that consistently completed in ~25 minutes now takes 50–60 minutes over the past week, despite no code changes. The SLA is 30 minutes, cloud costs matter, and the slowdown is gradual rather than sudden.

Key Details

- No recent code or configuration changes
- Runtime degradation observed over multiple days
- Data volume increased by ~10% only
- SLA breach risk
- Cost-sensitive environment

Why This Problem Is Tricky

The job still **completes successfully**.
There are **no failures, no alerts, no obvious errors**.

This makes it easy to:

- Ignore early warning signs
- Overreact by scaling infrastructure
- Miss the real cause until costs or SLAs spiral

This is a **classic slow-burn production issue**.

Clarifying Questions

Before proposing fixes, reduce uncertainty:

- Is the slowdown uniform across stages or localized?
- Do some tasks take significantly longer than others?
- Has key distribution changed over time?
- Are executors idle while one task runs long?
- Is the increase linear or exponential?

Confirmed Facts & Assumptions

After investigation:

- Runtime increase is gradual, not sudden
- Data volume growth (~10%) does not justify 2× runtime
- Cluster configuration is unchanged
- No retries or failures observed
- Job stages show long tail execution

Interpretation:

This is not normal scale-related behavior.

Key Observation

A small data increase causing a large runtime increase almost always indicates:

- Data skew
- Hot keys
- Uneven partition sizes

This shifts focus from **compute capacity** to **data distribution**.

Root Cause Analysis

Step 1: Analyze Spark Execution (Spark UI)

Key signals observed:

- Large gap between average and max task duration
- One or two tasks running far longer than others
- Many executors idle while a few remain busy

Conclusion:

Straggler tasks dominate stage completion time.

Step 2: Understand Why This Happens

Spark processes data in partitions:

- Each task handles one partition
- A stage completes only when the slowest task finishes

If:

- Certain keys grow disproportionately
- Aggregations or joins concentrate data

Then:

- One partition becomes a bottleneck
- Parallelism collapses silently

Step 3: Conceptual Root Cause

The root cause is **gradual data skew accumulation**:

- Certain keys grow over time
- No alerts fire because the job still “works”
- Performance degrades invisibly until SLA breaks

This is a **data evolution problem**, not an infrastructure one.

Step 4 : Why Obvious Fixes Are Suboptimal

- **Increasing cluster size:**
Temporarily reduces runtime but increases cost and leaves skew untouched.
- **Increasing executor memory:**
Helps marginally, but does not redistribute skewed data.
- **Restarting the job:**
Wastes time; the same data produces the same behavior.

Senior engineers fix **data imbalance**, not symptoms.

Step 5 : Validation of Root Cause

To confirm:

- Analyze key frequency distributions
- Inspect partition sizes before heavy joins/aggregations
- Identify outlier keys with extreme record counts

Outcome:

A small set of keys accounts for a large portion of data.

Step 6 : Corrective Actions

- Apply key salting for skewed dimensions
- Repartition data before expensive operations
- Use broadcast joins where applicable
- Add skew detection metrics to monitoring

These changes restore parallelism and control costs.

Step 7 : Result After Fix

- Tasks complete in similar time
- Executors are evenly utilized
- Runtime returns below 30 minutes
- Cloud cost stabilizes
- Performance degradation is prevented proactively

Final Resolution

- **Root Cause:** Gradual data skew leading to straggler tasks
- **Fix Applied:** Targeted skew mitigation and partition rebalancing

Key Learnings

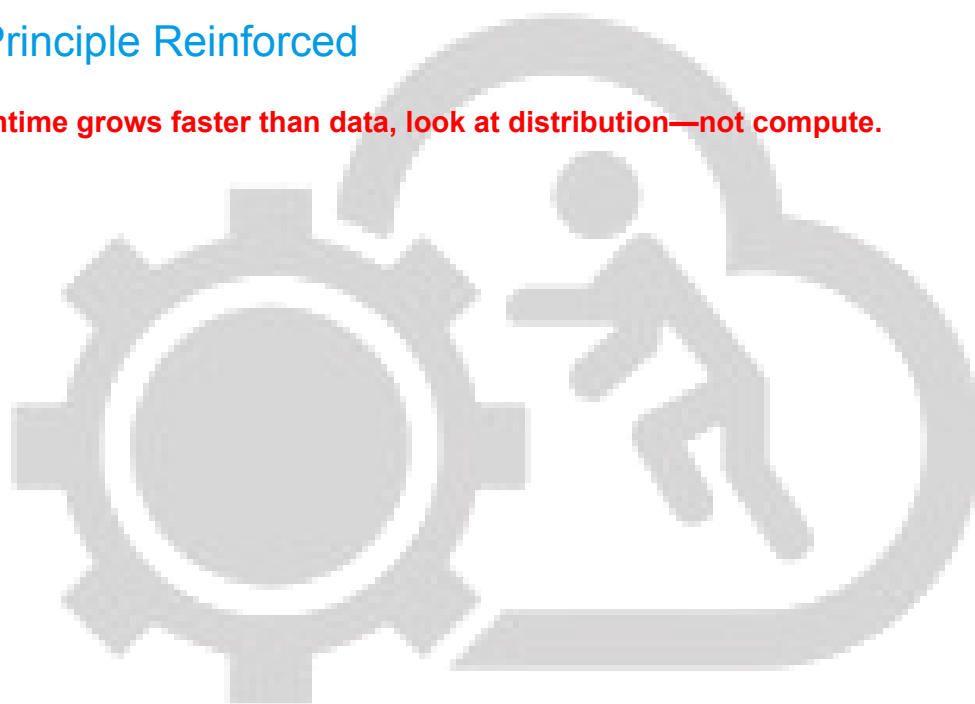
Gradual slowdowns are more dangerous than sudden failures

- Data distribution changes silently over time
- Scaling infrastructure is not a diagnostic tool
- Senior engineers read execution behavior, not just metrics

Core Principle Reinforced

When runtime grows faster than data, look at distribution—not compute.

■ ■ ■



Scenario 2

SQL Queries Running Slower Than Expected (Redshift)

Problem Statement

A critical daily reporting query on Amazon Redshift that previously completed within minutes now takes **10× longer**, even though there have been no schema or query changes. The SLA is 30 minutes, and other workloads must not be impacted.

Key Details

- No schema or SQL changes
- Sudden performance regression
- SLA: 30 minutes
- Query is business-critical
- Other queries must remain unaffected

Expected vs Actual Behavior

| Expected (Earlier) | Actual (Now) |
|-----------------------|----------------------------|
| Runtime: ~5 min | Runtime: ~50 min |
| Stable execution plan | Inefficient execution plan |
| SLA met | SLA breached |
| Low resource usage | High scan and join cost |

This gap indicates a **planning problem**, not a logic problem.

Why This Problem Is Deceptive

Nothing is obviously broken:

- The query still runs
- No deployment occurred
- No errors are reported

This often leads teams to:

- Scale the cluster prematurely
- Rewrite SQL without evidence
- Move workloads unnecessarily

In most cases, the real issue lies in **metadata, not SQL**.

Clarifying Questions

Before making changes, a senior engineer asks:

- When were table statistics last updated?
- Has data volume or distribution changed?
- Did the execution plan change?
- Are joins and filters still selective?
- Is table fragmentation present?

These questions focus on **how the database plans the query**.

Confirmed Facts & Assumptions

After investigation:

- Data volume increased gradually
- **ANALYZE** was not run recently
- Tables were not vacuumed
- Execution plan shows inefficient joins
- Infrastructure is healthy

Interpretation:

The optimizer is making decisions using outdated information.

What the System Thinks vs Reality

| Optimizer Assumption | Reality |
|----------------------------|-------------------------------|
| Table is small | Table has grown significantly |
| Low cardinality | High cardinality |
| Cheap join | Expensive join |
| Sequential scan acceptable | Causes major slowdown |

When statistics are stale, the optimizer guesses—and guesses wrong.

Root Cause Analysis

Step 1: Inspect Execution Plan

Observed issues:

- Sequential scans instead of optimized joins
- Large intermediate result sets
- Poor join order selection

Conclusion:

The optimizer lacks accurate statistics.

Step 2: Understand Redshift Maintenance

Redshift relies on:

- **ANALYZE** → to refresh table statistics
- **VACUUM** → to reduce fragmentation and maintain sort order

Without these:

- Cost estimates drift
- Query plans degrade silently
- Performance worsens over time

Step 3 : Wrong Approach vs Right Approach

Wrong Approach

- Increase node size
- Rewrite SQL blindly
- Run query on another cluster

Right Approach

- Inspect execution plan
- Refresh table statistics
- Vacuum fragmented tables

Senior engineers fix **planning first**, not hardware.

Step 4 : Validation of Root Cause

To confirm:

- Run **ANALYZE** on affected tables
- Perform **VACUUM** where required
- Compare execution plans before and after

Outcome:

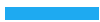
The optimizer selects an efficient plan and runtime drops significantly.

Step 5 : Corrective Actions

- Run **ANALYZE** regularly
- Schedule **VACUUM** jobs
- Review sort and distribution keys
- Monitor execution plan changes
- Automate maintenance tasks

These actions restore performance without increasing cost.

Step 6 : Result After Fix



| Before Fix | After Fix |
|-------------------|-----------------|
| Runtime: ~50 min | Runtime: ~5 min |
| SLA missed | SLA met |
| High scan cost | Optimized joins |
| Rising cloud cost | Stable cost |

Final Resolution

- **Root Cause:** Stale statistics and table fragmentation
- **Fix Applied:** ANALYZE and VACUUM to restore optimizer efficiency

Key Learnings

- SQL performance depends on metadata, not just queries
- Databases degrade gradually, not suddenly
- Scaling is not diagnosis
- Execution plans should always be checked first

Core Principle Reinforced

If the optimizer doesn't understand your data, even correct SQL will perform poorly.



Scenario 3

Spark Shuffle Becomes a Bottleneck

Problem Statement

A Spark batch job processing ~500 GB of data runs efficiently until it reaches the shuffle stage, where task execution time increases disproportionately. The SLA is one hour, the cluster is shared, and the job cannot be fully rewritten.

Key Details

- Data volume: ~500 GB
- Performance degradation starts at shuffle stage
- SLA: 1 hour
- Shared cluster environment
- Full rewrite not feasible

Expected vs Actual Behavior

| Expected | Actual |
|---------------------------|-----------------------------|
| Even task execution | A few tasks run much longer |
| Balanced executor usage | Some executors idle |
| Shuffle completes quickly | Shuffle dominates runtime |
| SLA achievable | SLA at risk |

This pattern indicates a **shuffle imbalance**, not a general compute shortage.

Why This Problem Is Tricky

- The job does not fail
- Earlier stages complete normally
- The slowdown appears only at shuffle

This often leads teams to:

- Add more memory
- Scale executors
- Blame cluster load

However, shuffle issues are usually **data distribution problems**, not hardware problems.

Clarifying Questions

Before acting, a senior engineer asks:

- Are all shuffle tasks slow or only a few?
- Do task runtimes show a long tail?
- Is key distribution uneven?
- Are shuffle partitions too large?
- Are executors waiting on a small number of tasks?

These questions focus on **execution imbalance**, not capacity.

Confirmed Facts & Assumptions

After investigation:

- Only a subset of shuffle tasks run much longer
- Large variance between average and max task time
- Executors are underutilized during shuffle
- No recent code or cluster changes
- Data volume alone does not explain the delay

Interpretation:

This is a **shuffle skew / partition sizing issue**.

What Spark Expects vs Reality

| Spark Assumption | Reality |
|-------------------------|-------------------------------|
| Even key distribution | Some keys are heavy |
| Similar partition sizes | Few partitions are very large |
| Parallel execution | Parallelism collapses |
| Fast shuffle completion | Straggler tasks dominate |

When a shuffle is unbalanced, Spark waits for the slowest task.

Root Cause Analysis

Step 1: Inspect Shuffle Metrics (Spark UI)

Observed:

- High max task duration vs average
- Few tasks handling most of the shuffle data
- Idle executors waiting for stragglers

Conclusion:

Shuffle parallelism is insufficient.

Step 2: Understand Shuffle Mechanics

During shuffle:

- Data is repartitioned by key
- Each task processes one partition
- Stage completion depends on the slowest task

If partitions are too large or skewed:

- One task becomes a bottleneck
- The entire stage slows down

Step 3: Conceptual Root Cause

The root cause is **insufficient shuffle parallelism**, often caused by:

- Skewed keys
- Too few shuffle partitions
- Uneven data distribution

This is a tuning and data-distribution problem.

Step 4: Wrong Approach vs Right Approach

Wrong Approach

- Increase executor memory
- Ignore shuffle imbalance
- Wait for job completion

Right Approach

- Increase shuffle partitions
- Rebalance data during shuffle
- Improve parallelism

Senior engineers fix **work distribution**, not just resource size.

Step 5 : Validation of Root Cause

To confirm:

- Increase **spark.sql.shuffle.partitions**
- Re-run the job
- Compare task duration distribution

Outcome:

Shuffle tasks become smaller and more evenly distributed.

Step 6 : Corrective Actions

- Increase shuffle partitions to improve parallelism
- Monitor partition size distribution
- Watch for skewed keys in joins or aggregations
- Tune shuffle settings incrementally

This reduces stragglers without increasing cluster cost.

Step 7 : Result After Fix

| Before Fix | After Fix |
|---------------------------|--------------------------|
| Few long-running tasks | Tasks complete uniformly |
| Executors idle | Executors fully utilized |
| Shuffle dominates runtime | Shuffle time reduced |
| SLA at risk | SLA met |

Final Resolution

- **Root Cause:** Shuffle imbalance due to insufficient partitioning
- **Fix Applied:** Increased shuffle partitions to improve parallelism

Key Learnings

- Shuffle stages expose data distribution issues
- Memory scaling does not fix skew
- Task time variance is a critical signal
- Parallelism tuning is a core Spark skill

Core Principle Reinforced

In Spark, performance problems at shuffle are usually about data distribution, not compute size.



Scenario 4

Streaming Job Slows Down During Peak Traffic

Problem Statement

A Kafka → Spark Structured Streaming pipeline that normally stays near real time starts lagging by **~30 minutes during peak traffic**. Downstream dashboards are real-time, infrastructure limits are fixed, and lag spikes unpredictably.

Key Details

- Source: Kafka
- Processing: Spark Structured Streaming
- Lag: ~30 minutes during peaks
- Real-time dashboards depend on this data
- Memory and CPU limits are fixed

Expected vs Actual Behavior

| Expected | Actual |
|---------------------------|--------------------------|
| Near real-time processing | 30-minute lag |
| Stable consumer lag | Lag spikes unpredictably |
| Dashboards up-to-date | Dashboards delayed |
| SLA met | SLA breached |

This clearly points to a **throughput mismatch**, not a correctness issue.

Why This Problem Is Misleading

The job is:

- Still running
- Not failing
- Eventually catching up

This often leads teams to:

- Increase micro-batch interval
- Blame temporary traffic spikes
- Accept lag as “normal during peaks”

In reality, lag growth means the system **cannot keep up with incoming data**.

Clarifying Questions

Before acting, a senior engineer asks:

- Is Kafka lag increasing faster than it is being processed?
- Are all consumers equally busy?
- How many partitions exist on the topic?
- Are executors fully utilized?
- Does lag disappear after traffic drops?

These questions distinguish **capacity issues** from **processing inefficiencies**.

Confirmed Facts & Assumptions

After investigation:

- Lag increases rapidly during traffic spikes
- Consumers are saturated
- Micro-batch processing time exceeds batch interval
- No data loss or parsing errors
- Infrastructure limits cannot be increased

Interpretation:

This is a **consumer throughput problem**, not a latency configuration issue.

What the System Assumes vs Reality

| System Assumption | Reality |
|---------------------------------|---------------------------------------|
| Consumers can keep up | Incoming rate exceeds processing rate |
| Batch interval tuning is enough | Processing capacity is insufficient |
| Lag will self-correct | Lag compounds during peaks |

Streaming systems fall behind when **ingest rate > processing rate**.

Root Cause Analysis

Step 1: Analyze Streaming Metrics

Observed:

- Growing Kafka consumer lag
- Executors busy but unable to drain backlog
- Micro-batches consistently running long

Conclusion:

The system lacks enough parallel consumers to handle peak load.

Step 2: Understand Streaming Throughput

In Kafka + Spark:

- Each partition can be consumed by only one task at a time
- Throughput scales with number of consumers and partitions
- Batch interval tuning does not increase capacity

If consumers are underscaled, lag is inevitable.

Step 3: Conceptual Root Cause

The root cause is **insufficient consumer parallelism**:

- Too few consumers for the partition count
- Processing cannot match peak ingest rate
- Lag accumulates rapidly

This is a scaling issue, not a configuration mistake.

Step 4 : Wrong Approach vs Right Approach

Wrong Approach

- Increase micro-batch interval
- Ignore lag until traffic drops
- Drop late-arriving messages

Right Approach

- Scale consumers
- Ensure enough Kafka partitions
- Balance load across executors

Senior engineers scale **throughput**, not delay processing.

Step 5 : Validation of Root Cause

To confirm:

- Increase consumer parallelism
- Monitor Kafka lag trend
- Compare processing rate vs ingest rate

Outcome:

Lag stabilizes and starts decreasing even during peak traffic.

Step 6 : Corrective Actions

- Scale Kafka consumers (increase parallelism)
- Ensure topic has sufficient partitions
- Align Spark tasks with partition count
- Monitor lag as a first-class metric

This restores real-time behavior without data loss.

Step 7 : Result After Fix

| Before Fix | After Fix |
|---------------------|----------------------------|
| Lag grows to 30 min | Lag remains near real time |
| Dashboards delayed | Dashboards up-to-date |
| Consumers saturated | Load balanced |
| SLA breached | SLA met |

Final Resolution

- **Root Cause:** Insufficient consumer throughput during peak load
- **Fix Applied:** Scaled consumers to match ingest rate

Key Learnings

- Streaming lag indicates throughput mismatch
- Batch interval tuning does not increase capacity
- Kafka partitioning limits parallelism
- Lag metrics are early warning signals

Core Principle Reinforced

In streaming systems, you fix lag by increasing throughput, not by waiting longer.



Scenario 5

High Garbage Collection (GC) Time in Spark

Problem Statement

A production Spark job intermittently fails because **garbage collection consumes more than 50% of total runtime**. The SLA is two hours, cluster memory is limited, and data volume has remained stable.

Key Details

- GC time exceeds 50% of runtime
- Job fails intermittently
- Data volume unchanged
- Cluster memory is constrained
- SLA: 2 hours

Expected vs Actual Behavior

| Expected | Actual |
|------------------------------------|---------------------------|
| Majority time spent on computation | Majority time spent in GC |
| Stable job execution | Intermittent failures |
| Predictable runtime | Unstable runtime |
| SLA achievable | SLA at risk |

This indicates a **memory management problem**, not a data growth problem.

Why This Problem Is Misleading

At first glance, high GC time looks like a simple memory shortage.

This often leads teams to:

- Increase executor memory
- Reduce partitions blindly
- Restart jobs repeatedly

But GC pressure is often caused by **how objects are created and serialized**, not just how much memory exists.

Clarifying Questions

Before acting, a senior engineer asks:

- Are many short-lived objects being created?
- Which serialization format is being used?
- Is GC time high even when memory is available?
- Are wide transformations creating large object graphs?
- Does GC spike correlate with shuffle stages?

These questions focus on **object lifecycle**, not raw memory size.

Confirmed Facts & Assumptions

After investigation:

- GC time spikes without corresponding data growth
- Executor memory is sufficient but poorly utilized
- Job uses default Java serialization
- Failures correlate with shuffle-heavy stages
- Restarting does not change behavior

Interpretation:

This is an **inefficient object serialization issue**, not a capacity issue.

What Spark Expects vs Reality

| Spark Expectation | Reality |
|--------------------------------|---------------------------|
| Efficient object serialization | Excessive object creation |
| Short GC pauses | Long GC cycles |
| Memory used for computation | Memory churned by GC |
| Stable execution | Intermittent failures |

When serialization is inefficient, GC dominates execution.

Root Cause Analysis

Step 1: Inspect GC Metrics

Observed:

- Frequent full GC cycles
- Long GC pause times
- High object allocation rate

Conclusion:

The JVM is spending more time cleaning memory than doing work.

Step 2: Understand Serialization Impact

Default Java serialization:

- Produces many intermediate objects
- Increases memory pressure
- Triggers frequent GC

More efficient serializers reduce object creation dramatically.

Step 3: Conceptual Root Cause

The root cause is **inefficient serialization and excessive object creation**:

- Too many temporary objects
- High allocation rate
- GC overwhelms execution

This is a **data representation issue**, not a memory size issue.

Step 4 : Wrong Approach vs Right Approach

Wrong Approach

- Increase executor memory
- Reduce partitions blindly
- Restart job

Right Approach

- Optimize serialization (e.g., Kryo)
- Reduce object creation
- Tune memory usage patterns

Senior engineers reduce **memory churn**, not just increase memory.

Step 5 : Validation of Root Cause

To confirm:

- Switch to a more efficient serializer
- Re-run the job
- Monitor GC time and pause duration

Outcome:

GC time drops significantly and job stabilizes.

Step 6 : Corrective Actions

- Use Kryo serialization
- Register commonly used classes
- Avoid unnecessary object creation
- Review data structures in transformations
- Monitor GC metrics continuously

These steps reduce GC pressure without increasing cluster size.

Step 7 : Result After Fix

| Before Fix | After Fix |
|-----------------------|------------------------|
| GC > 50% runtime | GC < 15% runtime |
| Intermittent failures | Stable execution |
| SLA at risk | SLA met |
| Memory churn | Efficient memory usage |

Final Resolution

- **Root Cause:** Inefficient serialization causing excessive GC
- **Fix Applied:** Optimized serialization to reduce object creation

Key Learnings

- High GC time is usually a symptom, not the problem
- Memory size does not fix poor object management
- Serialization choice matters at scale
- GC metrics are critical performance signals

Core Principle Reinforced

In Spark, reducing object churn is often more effective than adding memory.



Scenario 6

Query Performance Degrades After Column Addition

Problem Statement

After adding a new column to a **1 TB table**, several business-critical queries start running **5× slower**, even though the query logic itself was not changed. The SLA is one hour.

Key Details

- Table size: ~1 TB
- New column added recently
- Existing queries slowed down significantly
- Queries are business-critical
- SLA: 1 hour

Expected vs Actual Behavior

| Expected | Actual |
|-------------------------------|--------------------------|
| Queries run within SLA | Queries run 5× slower |
| Indexes / sort keys effective | Optimizer ignores them |
| Stable execution plans | Degraded execution plans |
| Predictable performance | SLA at risk |

This indicates a **storage or optimization issue**, not a query logic issue.

Why This Problem Is Deceptive

Adding a column feels like a *non-breaking* change:

- No query changes
- No data growth in existing columns
- No errors

This often leads teams to:

- Scale the cluster
- Rewrite queries
- Blame the database engine

In reality, schema changes can **invalidate physical optimizations**.

Clarifying Questions

Before reacting, a senior engineer asks:

- Are indexes or sort keys affected by the new column?
- Did table statistics change?
- Has data distribution or storage layout shifted?
- Are queries still filtering on the same keys?
- Did execution plans change after the column addition?

These questions focus on **how data is stored and accessed**, not SQL syntax.

Confirmed Facts & Assumptions

After investigation:

- New column was added to an existing large table
- Indexes / sort keys were not rebuilt
- Query execution plans changed
- Optimizer no longer uses optimal access paths
- Infrastructure remains healthy

Interpretation:

The schema change disrupted **physical data organization**.

What the Optimizer Expects vs Reality

| Optimizer Expectation | Reality |
|-------------------------|------------------------------|
| Indexed / sorted layout | Layout partially invalidated |
| Efficient data pruning | More data scanned |
| Stable access paths | Suboptimal scans |
| Predictable performance | Query slowdown |

When physical structures are outdated, the optimizer cannot work efficiently.

Root Cause Analysis

Step 1: Inspect Execution Plans

Observed:

- Indexes or sort keys no longer used
- Increased scan volume
- Higher I/O cost

Conclusion:

The optimizer lacks updated physical structures.

Step 2: Understand Impact of Column Addition

In large analytical systems:

- Adding columns can change row layout
- Indexes and sort keys may become ineffective
- Statistics may become inaccurate

Without rebuilding, performance silently degrades.

Step 3: Conceptual Root Cause

The root cause is **outdated indexes / sort keys after a schema change**.

This is a **storage-level optimization issue**, not a compute or SQL issue.

Step 4 : Wrong Approach vs Right Approach

Wrong Approach

- Increase cluster size
- Ignore the slowdown
- Drop the column

Right Approach

- Rebuild indexes or sort keys
- Refresh statistics
- Validate execution plans

Senior engineers restore **optimizer efficiency** before adding compute.

Step 5 : Validation of Root Cause

To confirm:

- Rebuild indexes / sort keys
- Refresh table statistics
- Compare execution plans before and after

Outcome:

Queries return to expected performance.

Step 6 : Corrective Actions

- Rebuild affected indexes or sort keys
- Refresh statistics after schema changes
- Review schema evolution impact on performance
- Add post-schema-change validation checks

These steps restore performance without extra cost.

Step 7 : Result After Fix

| Before Fix | After Fix |
|-------------------|--------------------|
| Queries 5× slower | Queries within SLA |
| Indexes ignored | Indexes used |
| High scan cost | Optimized access |
| SLA at risk | SLA met |

Final Resolution

- **Root Cause:** Indexes / sort keys not rebuilt after column addition
- **Fix Applied:** Rebuilt physical optimizations and refreshed stats

Key Learnings

- Schema changes can impact performance without breaking queries
- Physical optimizations matter at scale
- Scaling compute is not the first response
- Execution plans should be reviewed after schema changes

Core Principle Reinforced

After schema changes, always realign physical optimizations before scaling infrastructure.



Scenario 7

Network Latency Causes Job Delay

Problem Statement

A production ETL job reads data from **S3 into a Spark cluster**. Due to network throttling, the job's runtime doubles, putting the **1-hour SLA** at risk. The job is **I/O-bound**, and network capacity cannot be increased immediately.

Key Details

- Source: S3 → Spark
- Runtime increased ~2×
- Job is I/O-bound
- Network capacity fixed
- SLA: 1 hour

Expected vs Actual Behavior

| Expected | Actual |
|------------------------|--------------------------|
| Steady read throughput | Throttled reads |
| Predictable runtime | Runtime doubles |
| Executors kept busy | Executors waiting on I/O |
| SLA met | SLA at risk |

This pattern points to a **read-path bottleneck**, not a compute bottleneck.

Why This Problem Is Misleading

Because the job slows down uniformly:

- It's tempting to retry
- Or to add memory/CPU

But when executors are idle and waiting on data, adding compute **does not increase throughput**. The constraint is the **networked read path**.

Clarifying Questions

Before acting, a senior engineer asks:

- Are executors spending time waiting on input?
- Is read throughput capped during the job?
- Are files few and large, or many and small?
- Can reads be parallelized further?
- Do retries change throughput?

These questions isolate **I/O saturation** from compute issues.

Confirmed Facts & Assumptions

After investigation:

- Executors are underutilized
- Read throughput flattens during execution
- Data is stored as a small number of large objects
- Retrying does not improve speed
- Memory and CPU are not limiting factors

Interpretation:

This is a **network-bound read bottleneck**.

What the System Assumes vs Reality

| System Assumption | Reality |
|---------------------------|---------------------------|
| Reads scale automatically | Throughput is capped |
| More memory helps | I/O is the limiter |
| Retries may succeed | Same throttling applies |
| Compute is the bottleneck | Network is the bottleneck |

When I/O is saturated, compute sits idle.

Root Cause Analysis

Step 1: Inspect Executor Utilization

Observed:

- Low CPU usage
- High task wait time on input
- Stable but capped read throughput

Conclusion:

The job is waiting on data, not processing it.

Step 2: Understand S3 Read Behavior

- Each task reads independently
- Fewer, larger files limit parallel reads
- Network throttling caps per-stream throughput

Without enough parallelism, read speed plateaus.

Step 3: Conceptual Root Cause

The root cause is **insufficient parallel reads under network throttling**:

- Large objects limit concurrency
- Read throughput cannot scale
- Job runtime increases predictably

This is an **I/O parallelism issue**, not a Spark tuning issue.

Step 4 : Wrong Approach vs Right Approach

Wrong Approach

- Retry the job
- Increase executor memory
- Ignore the slowdown

Right Approach

- Split reads into smaller chunks
- Increase read parallelism
- Align tasks with available bandwidth

Senior engineers optimize **data access patterns** before scaling compute.

Step 5 : Validation of Root Cause

To confirm:

- Split large inputs into smaller chunks
- Increase concurrent read tasks
- Monitor aggregate read throughput

Outcome:

Throughput increases and runtime drops.

Step 6 : Corrective Actions

- Split large files to enable parallel reads
- Increase input partitioning where possible
- Align Spark tasks with S3 read concurrency
- Monitor read throughput as a primary metric

These changes improve performance without changing network capacity.

Step 7 : Result After Fix

| Before Fix | After Fix |
|--------------------------|-----------------------------|
| Runtime 2× slower | Runtime within SLA |
| Executors waiting on I/O | Executors fully utilized |
| Read throughput capped | Higher aggregate throughput |
| SLA at risk | SLA met |

Final Resolution

- **Root Cause:** Network-bound reads with insufficient parallelism
- **Fix Applied:** Increased parallel reads by splitting input

Key Learnings

- I/O-bound jobs require I/O-focused fixes
- Compute scaling does not fix network limits
- Parallelism applies to reads, not just tasks
- Executor idle time is a critical signal

Core Principle Reinforced

When a job is I/O-bound, optimize data access before adding compute.



Scenario 8

Skewed Join Causes Spark Job Slowness

Problem Statement

A Spark job joins two large tables, but execution becomes extremely slow because **one executor processes nearly 10× more data than others**. The SLA is two hours, the cluster is shared, and scaling the cluster is not an option.

Key Details

- Join between two large tables
- Severe data skew on join key
- One executor overloaded
- Cluster shared; nodes cannot be increased
- SLA: 2 hours

Expected vs Actual Behavior

| Expected | Actual |
|--------------------------|---------------------------|
| Even data distribution | One executor overloaded |
| Balanced task runtime | One task runs much longer |
| Executors fully utilized | Most executors idle |
| SLA achievable | SLA at risk |

This pattern clearly indicates a **skewed join**, not a general performance issue.

Why This Problem Is Deceptive

The job:

- Does not fail
- Appears to use sufficient resources
- Slows down only at the join stage

This often leads teams to:

- Increase executor memory
- Retry the job
- Assume cluster contention

However, when **one executor dominates runtime**, the issue is almost always **key skew**.

Clarifying Questions

Before acting, a senior engineer asks:

- Are a few join keys extremely frequent?
- Does Spark UI show one long-running task?
- Is data distribution uneven after shuffle?
- Can the join key be transformed safely?
- Is one side of the join small enough to broadcast?

These questions focus on **data distribution**, not compute size.

Confirmed Facts & Assumptions

After investigation:

- A small number of keys dominate the join
- One executor handles most shuffled data
- Other executors complete early and wait
- Increasing memory does not fix imbalance
- Cluster capacity is otherwise sufficient

Interpretation:

This is a **classic skewed join problem**.

What Spark Assumes vs Reality

| Spark Assumption | Reality |
|-------------------------|-----------------------------|
| Keys evenly distributed | Few keys extremely frequent |
| Similar partition sizes | One partition much larger |
| Parallel execution | Parallelism collapses |
| Fast join stage | One task dominates runtime |

Spark waits for the slowest task, regardless of how fast others finish.

Root Cause Analysis

Step 1: Inspect Shuffle and Task Metrics

Observed:

- Large variance between average and max task duration
- One shuffle partition significantly larger
- Idle executors during join stage

Conclusion:

Data skew is breaking parallelism.

Step 2: Understand Join Execution

In Spark joins:

- Data is shuffled by join key
- Each key goes to a single partition
- Hot keys create oversized partitions

This makes one executor the bottleneck.

Step 3: Conceptual Root Cause

The root cause is **uneven key distribution during the join**:

- Skewed keys concentrate data
- One task becomes a straggler
- SLA is breached despite sufficient resources

This is a **data modeling and distribution issue**.

Step 4 : Wrong Approach vs Right Approach

Wrong Approach

- Increase executor memory
- Retry the job
- Ignore the imbalance

Right Approach

- Salt the skewed key
- Broadcast smaller table (when applicable)
- Redistribute data before the join

Senior engineers fix **work distribution**, not just resource size.

Step 5 : Validation of Root Cause

To confirm:

- Apply key salting
- Re-run the join
- Observe partition sizes and task runtimes

Outcome:

Data is evenly distributed and task runtimes converge.

Step 6 : Corrective Actions

- Apply key salting for hot keys
- Broadcast the smaller table if size permits
- Monitor join key distributions
- Add skew detection metrics

These actions restore parallelism without scaling the cluster.

Step 7 : Result After Fix

| Before Fix | After Fix |
|-----------------------------|--------------------------|
| One executor overloaded | Load evenly distributed |
| Long-running straggler task | Tasks complete uniformly |
| Executors idle | Executors fully utilized |
| SLA at risk | SLA met |

Final Resolution

- **Root Cause:** Skewed join key causing uneven data distribution
- **Fix Applied:** Key salting to rebalance join workload

Key Learnings

- Skewed joins are a leading cause of Spark slowness
- Memory scaling does not fix data skew
- Task duration variance is the key signal
- Data distribution matters more than raw compute

Core Principle Reinforced

In distributed joins, fixing data distribution is more effective than adding resources.



Scenario 9

Slow Downstream Writes to Redshift

Problem Statement

A Spark job completes all transformations on time, but **writing the final results to Amazon Redshift takes 10× longer than expected**. The SLA is one hour, the Redshift cluster is shared, and multiple tables are updated in this step.

Key Details

- Spark processing completes successfully
- Write phase is the bottleneck
- Redshift cluster is shared
- Multiple target tables
- SLA: 1 hour

Expected vs Actual Behavior

| Expected | Actual |
|-----------------------|-------------------------------|
| Fast write completion | Writes dominate total runtime |
| Balanced job stages | Job stuck at sink stage |
| SLA met | SLA at risk |
| Predictable load | Heavy write contention |

This pattern indicates a **sink-side throughput issue**, not a Spark compute issue.

Why This Problem Is Misleading

Because:

- Spark processing finishes correctly
- No failures or retries occur
- Only the final stage is slow

Teams often assume:

- Redshift is “slow today”
- Scaling the cluster is required
- Retrying might help

In reality, **how data is written** matters more than **how fast Spark processes it**.

Clarifying Questions

Before acting, a senior engineer asks:

- Are writes row-based or batched?
- Are multiple writers competing for Redshift resources?
- Are COPY-based loads being used?
- Is network or commit overhead dominating?
- Does write time scale linearly with data size?

These questions isolate **write pattern inefficiency** from infrastructure limits.

Confirmed Facts & Assumptions

After investigation:

- Data is written using small, frequent writes
- Many concurrent writers hit Redshift
- COPY-style bulk loads are not used
- Retrying does not improve speed
- Cluster capacity is otherwise stable

Interpretation:

This is a **write amplification and commit overhead issue**.

What the System Expects vs Reality

| System Expectation | Reality |
|--------------------------|---------------------------|
| Writes are batched | Writes are fragmented |
| Few large commits | Many small commits |
| High throughput | Commit overhead dominates |
| Sink keeps up with Spark | Sink becomes bottleneck |

Databases handle **bulk loads** far better than **row-level inserts**.

Root Cause Analysis

Step 1: Analyze Write Pattern

Observed:

- High number of small write operations
- Significant commit and network overhead
- Redshift waiting on frequent transactions

Conclusion:

The write path is inefficient.

Step 2: Understand Redshift Write Characteristics

Redshift is optimized for:

- Large batch inserts
- COPY-based bulk loads from S3

It performs poorly with:

- Row-by-row inserts
- Many concurrent small writes

Step 3: Conceptual Root Cause

The root cause is **inefficient write strategy**:

- Too many small writes
- Excessive commit overhead
- Poor alignment with Redshift's ingestion model

This is a **sink optimization problem**, not a Spark problem.

Step 4 : Wrong Approach vs Right Approach

Wrong Approach

- Increase Redshift cluster size
- Retry writes
- Ignore slow sink

Right Approach

- Batch writes
- Use bulk loading (COPY)
- Reduce commit frequency

Senior engineers optimize **how data enters the database** before scaling it.

Step 5 : Validation of Root Cause

To confirm:

- Switch to batched writes or COPY-based ingestion
- Reduce number of write operations
- Measure write throughput

Outcome:

Write time drops dramatically and SLA is met.

Step 6 : Corrective Actions

- Batch records before writing
- Stage data to S3 and use Redshift COPY
- Limit concurrent writers
- Monitor write throughput separately from Spark runtime

These changes improve performance without increasing cluster cost.

Step 7 : Result After Fix

| Before Fix | After Fix |
|----------------------|----------------------|
| Writes 10× slower | Writes within SLA |
| High commit overhead | Efficient bulk loads |
| Sink bottleneck | Balanced pipeline |
| SLA at risk | SLA met |

Final Resolution

- **Root Cause:** Inefficient row-level writes to Redshift
- **Fix Applied:** Batched / bulk writes aligned with Redshift ingestion

Key Learnings

- Databases are often the slowest part of pipelines
- Write patterns matter as much as read patterns
- Bulk ingestion beats row inserts at scale
- Sink performance must be monitored explicitly

Core Principle Reinforced

In data pipelines, optimize the sink before scaling the system.



Scenario 10

Spark Job Slows Down After Dataset Growth

Problem Statement

A Spark batch job that previously processed **500 GB** of data now processes **1.2 TB** and runs **3× slower**, even though the transformation logic and cluster size remain unchanged. The **2-hour SLA** is now at risk.

Key Details

- Dataset growth: 500 GB → 1.2 TB
- Transformation logic unchanged
- Cluster size fixed
- Runtime increased ~3×
- SLA: 2 hours

Expected vs Actual Behavior

| Expected | Actual |
|----------------------------|----------------------|
| Moderate runtime increase | Runtime increases 3× |
| Even workload distribution | Oversized partitions |
| Balanced task execution | Straggler tasks |
| SLA met | SLA breached |

This pattern points to a **parallelism breakdown**, not inefficient logic or cluster failure.

Why This Problem Is Misleading

Because:

- Code did not change
- Infrastructure did not change
- Job still completes successfully

Teams often assume:

- The cluster is too small
- Joins suddenly became inefficient
- Retrying might help

In reality, **data growth silently breaks partition assumptions.**

Clarifying Questions

Before acting, a senior engineer asks:

- How many partitions exist after data growth?
- How large is each partition now?
- Do task runtimes show a long tail?
- Are executors idle while few tasks run long?
- Is repartitioning applied before heavy stages?

These questions focus on **work distribution**, not resource size.

Confirmed Facts & Assumptions

After investigation:

- Partition count remained unchanged
- Individual partitions grew significantly larger
- A few tasks dominate total runtime
- Executors remain underutilized
- No join or transformation inefficiency detected

Interpretation:

The job is **under-parallelized for the new data volume.**

What Spark Assumes vs Reality

| Spark Assumption | Reality |
|-------------------------------|----------------------------|
| Partitions scale with data | Partition count is static |
| Tasks complete uniformly | Few tasks dominate runtime |
| Parallelism remains effective | Parallelism collapses |
| Runtime scales linearly | Runtime degrades sharply |

Spark does not automatically adjust partitioning as data grows.

Root Cause Analysis

Step 1: Inspect Task and Partition Metrics

Observed:

- High variance between average and max task duration
- Oversized partitions
- Executors idle during late stages

Conclusion:

Partition sizing no longer matches dataset scale.

Step 2: Understand Impact of Dataset Growth

As data grows:

- Fixed partitions become heavier
- Each task does more work
- Effective parallelism decreases

Without repartitioning, slowdown is inevitable.

Step 3: Conceptual Root Cause

The root cause is **static partitioning in a growing dataset**:

- Partition count unchanged
- Load unevenly distributed
- Straggler tasks extend runtime

This is a **scalability design issue**, not a compute issue.

Step 4 : Wrong Approach vs Right Approach

Wrong Approach

- Increase cluster size
- Retry the job

Right Approach

- Repartition data to restore parallelism

Senior engineers scale **work distribution** before scaling infrastructure.

Step 5 : Validation of Root Cause

To confirm:

- Increase partition count appropriately
- Re-run the job
- Compare task runtime distribution

Outcome:

Tasks finish uniformly and runtime drops significantly.

Step 6 : Corrective Actions

- Repartition data based on current volume
- Periodically review partition strategy
- Monitor partition sizes as data grows
- Align partition count with cluster capacity

These steps restore performance without increasing cost.

Step 7 : Result After Fix

| Before Fix | After Fix |
|-------------------|--------------------|
| Runtime 3× slower | Runtime within SLA |
| Straggler tasks | Balanced execution |
| Idle executors | Full utilization |
| SLA breached | SLA met |

Final Resolution

- **Root Cause:** Partitions did not scale with dataset growth
- **Fix Applied:** Repartitioned data to restore parallelism

Key Learnings

- Data growth can silently break performance
- Partition strategy must evolve with scale
- Scaling compute is not the first response
- Task duration variance is a key signal

Core Principle Reinforced

As data grows, partitioning must scale with it—or performance will degrade.



Scenario 11

DAG Task Queueing Causes Latency (Airflow)

Problem Statement

An Airflow DAG experiences significant delays because tasks remain **queued for long periods**, even though the tasks are independent. This queueing delays downstream processing and puts the **1-hour SLA** at risk.

Key Details

- Orchestration tool: Airflow
- Tasks are independent
- Workers are limited
- Downstream processing delayed
- SLA: 1 hour

Expected vs Actual Behavior

| Expected | Actual |
|-----------------------------------|-------------------------------------|
| Tasks start soon after scheduling | Tasks remain queued |
| DAG progresses smoothly | Artificial waiting time |
| Parallel execution | Underutilized potential parallelism |
| SLA met | SLA breached |

This points to an **orchestration capacity issue**, not a task execution issue.

Why This Problem Is Misleading

Because:

- Tasks themselves run fast once started
- No task failures are observed
- DAG logic is correct

Teams often focus on:

- Reordering tasks
- Reducing task size

But queueing happens **before tasks even start**, which indicates a **worker bottleneck**.

Clarifying Questions

Before acting, a senior engineer asks:

- How many tasks are in the queue vs running?
- Are tasks truly independent?
- Is worker utilization consistently high?
- Are task slots the limiting factor?
- Is the delay happening at scheduling or execution?

These questions isolate **scheduler/worker capacity** from task logic.

Confirmed Facts & Assumptions

After investigation:

- Tasks spend most of their time in the queued state
- Workers are fully utilized
- Tasks do not depend on each other
- DAG structure is not the limiting factor
- Execution time is small compared to queue time

Interpretation:

The bottleneck is **insufficient worker capacity**.

What the Scheduler Assumes vs Reality

| Scheduler Assumption | Reality |
|-----------------------------|------------------------------------|
| Enough workers available | Workers are saturated |
| Tasks can start immediately | Tasks wait in queue |
| DAG parallelism utilized | Parallelism constrained by workers |
| Execution dominates runtime | Queueing dominates runtime |

Orchestration latency is often **invisible but costly**.

Root Cause Analysis

Step 1: Analyze Task States

Observed:

- Large number of tasks in **queued** state
- Few tasks actively running
- Stable but capped throughput

Conclusion:

Workers are the limiting resource.

Step 2: Understand Airflow Execution Model

In Airflow:

- Each task requires a worker slot
- Independent tasks can run in parallel
- Worker limits directly cap throughput

If workers are insufficient, queueing is inevitable.

Step 3: Conceptual Root Cause

The root cause is **worker under-provisioning**:

- DAG has parallelizable tasks
- Execution capacity is too low
- Artificial latency accumulates

This is an **orchestration scaling issue**, not a DAG design flaw.

Step 4 : Wrong Approach vs Right Approach

Wrong Approach

- Ignore queueing
- Only reorder DAG
- Reduce task size

Right Approach

- Increase number of workers
- Align worker count with task parallelism
- Prioritize critical tasks

Senior engineers scale **orchestration capacity** before restructuring workflows.

Step 5 : Validation of Root Cause

To confirm:

- Increase worker count
- Observe queue depth and task start time
- Measure end-to-end DAG runtime

Outcome:

Queue time drops and SLA is met.

Step 6 : Corrective Actions

- Increase Airflow workers
- Monitor queued vs running task metrics
- Set appropriate task concurrency limits
- Prioritize SLA-critical DAGs

These actions reduce latency without changing task logic.

Step 7 : Result After Fix

| Before Fix | After Fix |
|-------------------------------|-------------------------|
| Long queue times | Tasks start immediately |
| Underutilized DAG parallelism | Full parallel execution |
| SLA breached | SLA met |
| Artificial latency | Predictable execution |

Final Resolution

- **Root Cause:** Insufficient Airflow worker capacity
- **Fix Applied:** Scaled workers to match task parallelism

Key Learnings

- Orchestration can be a hidden bottleneck
- Queue time matters as much as execution time
- Independent tasks should run in parallel
- Worker capacity defines DAG throughput

Core Principle Reinforced

In orchestration systems, artificial latency comes from queues—scale workers before redesigning workflows



Scenario 12

Slow Join on Large Tables in Production

Problem Statement

A daily production Spark job performs a join between **two very large tables (~2 TB each)**. The job previously completed in about **1 hour**, but now takes **nearly 3 hours**, breaching the **2-hour SLA**. The cluster is shared, and the tables cannot be split.

Key Details

- Join between two ~2 TB tables
- Runtime increased from 1 hour → 3 hours
- Cluster is shared
- Tables cannot be split
- SLA: 2 hours

Expected vs Actual Behavior

| Expected | Actual |
|---------------------------|---------------------------|
| Join completes within SLA | Join takes 3 hours |
| Balanced shuffle workload | Heavy shuffle overhead |
| Executors fully utilized | Executors wait on shuffle |
| Predictable runtime | SLA breached |

This strongly suggests a **join execution inefficiency**, not general cluster slowness.

Why This Problem Is Misleading

Because:

- The job still succeeds
- No code changes were made
- Data volume appears “known”

Teams often jump to:

- Increasing cluster size
- Retrying the job
- Accepting longer runtimes

However, large joins are often slow due to **shuffle behavior**, not lack of compute.

Clarifying Questions

Before acting, a senior engineer asks:

- Is one side of the join significantly smaller after filters?
- How much data is being shuffled?
- Are join keys evenly distributed?
- Can broadcast join be applied safely?
- Is the join causing large intermediate data?

These questions focus on **join strategy**, not infrastructure.

Confirmed Facts & Assumptions

After investigation:

- One side of the join is much smaller after filtering
- Spark performs a full shuffle join
- Large shuffle stages dominate runtime
- Cluster resources are otherwise stable
- Retrying does not improve performance

Interpretation:

The job is using an **inefficient join strategy**.

What Spark Assumes vs Reality

| Spark Assumption | Reality |
|--------------------------------|---------------------------------------|
| Shuffle join is acceptable | Shuffle dominates runtime |
| Both sides need repartitioning | One side is small enough to broadcast |
| Parallelism sufficient | Network and shuffle bottleneck |
| Runtime scales linearly | Runtime explodes |

Shuffle-heavy joins are one of the most expensive Spark operations.

Root Cause Analysis

Step 1: Inspect Join and Shuffle Metrics

Observed:

- Very large shuffle read/write volumes
- Long-running shuffle stages
- Executors waiting on network I/O

Conclusion:

Shuffle is the primary bottleneck.

Step 2: Understand Broadcast Join Optimization

In Spark:

- Broadcast joins avoid shuffling the larger table
- The smaller table is sent to all executors
- Join becomes local and much faster

When applicable, broadcast joins drastically reduce runtime.

Step 3: Conceptual Root Cause

The root cause is **using a shuffle join when a broadcast join is feasible**:

- Unnecessary data movement
- Excessive network and disk I/O
- Extended execution time

This is a **join strategy issue**, not a capacity issue.

Step 4 : Wrong Approach vs Right Approach

Wrong Approach

- Increase cluster size
- Retry the job
- Ignore the slowdown

Right Approach

- Broadcast the smaller table
- Reduce shuffle volume
- Optimize join strategy

Senior engineers optimize **execution plans** before scaling clusters.

Step 5 : Validation of Root Cause

To confirm:

- Force or enable broadcast join
- Re-run the job
- Measure shuffle volume and runtime

Outcome:

Shuffle is eliminated or reduced significantly, and runtime drops below SLA.

Step 6: Corrective Actions

- Apply broadcast join where safe
- Validate table size thresholds for broadcasting
- Monitor shuffle metrics for joins
- Review join strategies as data evolves

These actions improve performance without increasing infrastructure cost.

Sep 7 : Result After Fix

| Before Fix | After Fix |
|--------------------|-----------------|
| Runtime ~3 hours | Runtime ~1 hour |
| Heavy shuffle | Minimal shuffle |
| Network bottleneck | Local joins |
| SLA breached | SLA met |

Final Resolution

- **Root Cause:** Inefficient shuffle join on large tables
- **Fix Applied:** Broadcast join to eliminate unnecessary shuffle

Key Learnings

- Joins are among the most expensive Spark operations
- Shuffle volume is a critical performance metric
- Broadcast joins can drastically reduce runtime
- Scaling clusters should be the last step

Core Principle Reinforced

Optimize join strategy before adding compute—shuffle is often the real enemy.



Scenario 13

Slow Job After Code Refactor

Problem Statement

After a recent code refactor, a Spark job's runtime **doubles**, even though the **data volume remains unchanged**. The SLA is two hours, cluster resources are limited, and the slowdown was not anticipated.

Key Details

- Recent refactor introduced additional transformations
- Data volume unchanged
- Cluster resources limited
- Runtime doubled
- SLA: 2 hours

Expected vs Actual Behavior

| Expected | Actual |
|---------------------------|---------------------|
| Similar runtime as before | Runtime doubled |
| Efficient transformations | Additional overhead |
| Stable execution | SLA at risk |
| Predictable performance | Unexpected slowdown |

This indicates a **logic-level inefficiency**, not a scale or data issue.

Why This Problem Is Misleading

Because:

- The refactor was intended to improve code quality
- No data growth occurred
- Infrastructure did not change

Teams often assume:

- The cluster is underpowered
- The slowdown is temporary
- Retrying might help

In reality, **additional transformations can silently increase execution cost.**

Clarifying Questions

Before acting, a senior engineer asks:

- Which new transformations were added?
- Are there extra shuffles or wide transformations?
- Did the refactor introduce unnecessary passes over data?
- Are any transformations redundant or repeated?
- Has execution plan complexity increased?

These questions focus on **execution cost**, not infrastructure.

Confirmed Facts & Assumptions

After investigation:

- New transformations introduced extra stages
- Additional shuffles appear in the execution plan
- CPU usage increased without data growth
- Memory and cluster limits remain unchanged
- Retrying does not improve runtime

Interpretation:

The refactor introduced **inefficient transformations**.

What the Developer Intended vs Reality

| Intended | Reality |
|-----------------------------|---------------------------|
| Cleaner, more readable code | Heavier execution plan |
| Same performance | Runtime doubled |
| Minimal overhead | Extra stages and shuffles |
| Safe refactor | Performance regression |

Refactors improve maintainability—but can harm performance if not profiled.

Root Cause Analysis

Step 1: Compare Execution Plans

Observed:

- More stages than before
- New shuffle boundaries
- Higher task execution time

Conclusion:

The refactor increased computational and shuffle overhead.

Step 2: Understand Transformation Cost

In Spark:

- Each wide transformation adds a shuffle
- Multiple passes over data multiply cost
- Small logical changes can have large physical impact

Step 3: Conceptual Root Cause

The root cause is **inefficient transformation design introduced during refactor**:

- Redundant transformations
- Unnecessary shuffles
- Poor stage consolidation

This is a **code-level performance regression**.

Step 4 : Wrong Approach vs Right Approach

Wrong Approach

- Increase cluster size
- Retry the job

Right Approach

- Profile the refactored job
- Optimize or remove unnecessary transformations
- Reduce shuffles and stages

Senior engineers fix **execution efficiency** before adding resources.

Step 5 : Validation of Root Cause

To confirm:

- Remove or optimize added transformations
- Re-run the job
- Compare stage count and runtime

Outcome:

Runtime drops back within SLA.

Step 6 : Corrective Actions

- Profile refactors using Spark UI
- Eliminate redundant transformations
- Cache intermediate results where appropriate
- Consolidate multiple transformations
- Add performance regression checks

These steps restore performance without increasing cost.

Step 7 : Result After Fix

| Before Fix | After Fix |
|----------------------------|----------------------------|
| Runtime doubled | Runtime within SLA |
| Extra stages | Optimized execution plan |
| SLA at risk | SLA met |
| Limited resources stressed | Resources used efficiently |

Final Resolution

- **Root Cause:** Inefficient transformations introduced during refactor
- **Fix Applied:** Optimized transformation logic and execution plan

Key Learnings

- Refactors can introduce hidden performance costs
- Transformation count and type matter
- Execution plans must be reviewed after refactors
- Readability and performance must be balanced

Core Principle Reinforced

Every refactor should be profiled—clean code is not always fast code.



Scenario 14

Slow ETL Due to Excessive Logging

Problem Statement

A production ETL job starts running **50% slower** after extensive logging is enabled. Logs are required for compliance, storage throughput is limited, and the **1-hour SLA** is now at risk.

Key Details

- Significant increase in log volume
- Logs required for compliance
- Storage throughput limited
- Runtime increased by ~50%
- SLA: 1 hour

Expected vs Actual Behavior

| Expected | Actual |
|------------------------------|---------------------------|
| Logs captured without impact | Logging dominates runtime |
| Stable ETL execution | Job slowed by I/O waits |
| SLA met | SLA breached |
| Compute-bound workload | I/O-bound due to logging |

This points to a **logging-induced I/O bottleneck**, not a compute issue.

Why This Problem Is Misleading

Logging is often considered “cheap”:

- No data logic changed
- No errors introduced
- Job still completes

Teams may assume:

- The slowdown is temporary
- Retrying might help
- More storage throughput is required

In reality, **excessive synchronous logging can silently throttle pipelines.**

Clarifying Questions

Before acting, a senior engineer asks:

- What log level is enabled in production?
- Are logs written synchronously?
- How much data is logged per record?
- Are logs required at this granularity?
- Does logging scale with data volume?

These questions isolate **logging overhead** from core processing.

Confirmed Facts & Assumptions

After investigation:

- Debug-level logging enabled
- Logs written frequently to external storage
- Executors spend time waiting on I/O
- Compute resources are underutilized
- Storage throughput is the limiting factor

Interpretation:

The job is **I/O-bound due to excessive logging.**

What the System Assumes vs Reality

| Assumption | Reality |
|-------------------------------|-------------------------------|
| Logging cost is negligible | Logging dominates I/O |
| Logs don't affect performance | Logs slow execution |
| Compute is the bottleneck | Storage I/O is the bottleneck |
| More storage is required | Less logging is required |

Root Cause Analysis

Step 1: Analyze I/O Patterns

Observed:

- High write volume during execution
- Executors blocked on log writes
- Increased job runtime proportional to log volume

Conclusion:

Logging overhead is the primary cause of slowdown.

Step 2: Understand Logging Impact

In production ETL systems:

- Logging often scales with record count
- High verbosity multiplies I/O operations
- Storage throughput becomes a hidden bottleneck

Compliance does not require **maximum verbosity**.

Step 3: Conceptual Root Cause

The root cause is **overly verbose logging in a data-intensive pipeline**:

- Excessive log volume
- I/O saturation
- Reduced processing throughput

This is a **configuration issue**, not a data or compute issue.

Step 4 : Wrong Approach vs Right Approach

Wrong Approach

- Increase storage throughput
- Ignore the slowdown
- Retry the job

Right Approach

- Reduce logging level
- Log only actionable events
- Balance compliance and performance

Senior engineers tune **observability**, not just infrastructure.

Step 5 : Validation of Root Cause

To confirm:

- Reduce logging verbosity
- Re-run the job
- Measure runtime and I/O usage

Outcome:

Runtime improves significantly while compliance needs are still met.

Step 6 : Corrective Actions

- Reduce logging level in production
- Avoid per-record logging
- Aggregate logs where possible
- Separate audit logs from debug logs
- Review logging strategy periodically

These changes restore performance without sacrificing compliance.

Step 7 : Result After Fix

| Before Fix | After Fix |
|-------------------|----------------------|
| Runtime +50% | Runtime within SLA |
| High I/O wait | Balanced I/O usage |
| Executors blocked | Executors productive |
| SLA breached | SLA met |

Final Resolution

- **Root Cause:** Excessive logging causing I/O saturation
- **Fix Applied:** Reduced logging verbosity while meeting compliance

Key Learnings

- Logging can be a performance bottleneck
- Observability must be tuned for scale
- Compliance does not mean maximum verbosity
- I/O issues often hide in non-obvious places

Core Principle Reinforced

Observability is powerful—but uncontrolled logging can silently slow production pipelines.



Scenario 15

Slow Job Due to Skewed Aggregations

Problem Statement

A Spark job performs aggregations by key, but **one key accounts for nearly 90% of the records**, causing a small number of executors to process most of the data. This results in **straggler tasks** and puts the **2-hour SLA** at risk. The cluster is shared, and data volume cannot be reduced.

Key Details

- Aggregation by key
- Extreme key skew (one key ~90% of data)
- Straggler tasks observed
- Cluster shared
- SLA: 2 hours

Expected vs Actual Behavior

| Expected | Actual |
|-------------------------------|----------------------------|
| Even aggregation workload | One task dominates runtime |
| Similar task completion times | Long-running stragglers |
| Executors fully utilized | Most executors idle |
| SLA met | SLA breached |

This pattern clearly indicates **aggregation skew**, not general performance issues.

Why This Problem Is Misleading

Because:

- The job does not fail
- Total data volume is stable
- Infrastructure appears sufficient

Teams often attempt:

- Increasing cluster size
- Retrying the job
- Ignoring a “temporary” slowdown

But aggregation skew **cannot be fixed by adding compute alone.**

Clarifying Questions

Before acting, a senior engineer asks:

- Which keys dominate the aggregation?
- How skewed is the key distribution?
- Do task runtimes show extreme variance?
- Is aggregation happening before or after filtering?
- Can keys be safely transformed?

These questions focus on **data distribution**, not compute size.

Confirmed Facts & Assumptions

After investigation:

- One key holds ~90% of records
- A single task processes most aggregation work
- Executors wait for the straggler to finish
- Retrying does not change behavior
- Cluster capacity is otherwise sufficient

Interpretation:

This is a **classic aggregation skew problem.**

What Spark Assumes vs Reality

| Spark Assumption | Reality |
|-------------------------------|---------------------------|
| Keys evenly distributed | One key dominates |
| Aggregations parallelize well | Parallelism collapses |
| Tasks finish together | One task runs much longer |
| Runtime scales linearly | Runtime spikes |

Spark must wait for the slowest aggregation task.

Root Cause Analysis

Step 1: Inspect Aggregation Metrics

Observed:

- Large variance between average and max task duration
- Single partition much larger than others
- Executors idle during final aggregation stage

Conclusion:

Key skew is breaking parallelism.

Step 2: Understand Aggregation Execution

In Spark:

- Aggregations shuffle data by key
- Each key is processed by one task
- Hot keys concentrate data in one partition

This makes one executor the bottleneck.

Step 3: Conceptual Root Cause

The root cause is **extreme key skew during aggregation**:

- One key dominates input
- One task does most of the work
- SLA is breached despite available resources

This is a **data modeling and distribution issue**.

Step 4 : Wrong Approach vs Right Approach

Wrong Approach

- Increase cluster size
- Retry the job
- Ignore the slowdown

Right Approach

- Salt the aggregation key
- Pre-aggregate where possible
- Redistribute workload

Senior engineers fix **data distribution**, not infrastructure.

Step 5 : Validation of Root Cause

To confirm:

- Apply key salting
- Re-run the aggregation
- Observe task runtime distribution

Outcome:

Aggregation workload is evenly distributed and runtime improves.

Step 6 : Corrective Actions

- Salt heavily skewed keys
- Perform partial aggregations before shuffle
- Monitor key distribution over time
- Add skew detection to pipelines

These steps restore parallelism without scaling the cluster.

Step 7 : Result After Fix

| Before Fix | After Fix |
|-----------------------|--------------------------|
| One straggler task | Balanced task durations |
| Executors idle | Executors fully utilized |
| SLA breached | SLA met |
| Unpredictable runtime | Stable performance |

Final Resolution

- **Root Cause:** Skewed aggregation key causing stragglers
- **Fix Applied:** Key salting to rebalance aggregation workload

Key Learnings

- Aggregation skew is common at scale
- Task runtime variance is a key signal
- Compute scaling does not fix skew
- Data distribution drives performance

Core Principle Reinforced

In distributed aggregations, fixing skew beats adding compute every time.

