

ORCHESTRATION & SCHEDULING ISSUES

Real Interview Scenarios
& How to Handle Them

A practical guide for Data Engineers to answer
real-world orchestration and scheduling questions with confidence

By Ankita Gulati

Senior Data Engineer | Interview Mentor

Interview Edition • Practical • Real Scenarios



Table Of Content

Scenario 1.....	5
Airflow DAG Tasks Stuck in Queued State.....	5
Problem Statement.....	5
Expected vs Actual Behavior.....	5
Why This Happens Frequently.....	6
Clarifying Questions.....	6
Confirmed Facts & Assumptions.....	6
What Teams Often Assume vs Reality.....	7
Root Cause Analysis.....	7
Final Resolution.....	9
Key Learnings.....	9
Core Principle Reinforced.....	9
Scenario 2.....	10
Airflow DAG Fails Due to Circular Dependencies.....	10
Problem Statement.....	10
Expected vs Actual Behavior.....	10
Why This Happens Often.....	11
Clarifying Questions.....	11
Confirmed Facts & Assumptions.....	11
What Teams Often Assume vs Reality.....	12
Root Cause Analysis.....	12
Final Resolution.....	14
Key Learnings.....	14
Core Principle Reinforced.....	14
Scenario 3.....	15
Timezone Misconfiguration Delays Airflow DAG Runs.....	15
Problem Statement.....	15
Expected vs Actual Behavior.....	15
Why This Happens Frequently.....	16
Clarifying Questions.....	16
Confirmed Facts & Assumptions.....	16
What Teams Often Assume vs Reality.....	17
Root Cause Analysis.....	17
Final Resolution.....	19
Key Learnings.....	19
Core Principle Reinforced.....	19

Scenario 4.....	20
Airflow Sensor Tasks Time Out and Delay DAGs.....	20
Problem Statement.....	20
Expected vs Actual Behavior.....	20
Why This Happens Frequently.....	21
Clarifying Questions.....	21
Confirmed Facts & Assumptions.....	21
What Teams Often Assume vs Reality.....	22
Root Cause Analysis.....	22
Final Resolution.....	24
Key Learnings.....	24
Core Principle Reinforced.....	24
Scenario 5.....	25
Resource Contention Between Multiple Airflow DAGs.....	25
Problem Statement.....	25
Expected vs Actual Behavior.....	25
Why This Happens Frequently.....	26
Clarifying Questions.....	26
Confirmed Facts & Assumptions.....	26
What Teams Often Assume vs Reality.....	27
Root Cause Analysis.....	27
Final Resolution.....	29
Key Learnings.....	29
Core Principle Reinforced.....	29
Scenario 6.....	30
DAG Completes Successfully but Tasks Are Skipped.....	30
Problem Statement.....	30
Expected vs Actual Behavior.....	30
Why This Problem Is Dangerous.....	31
Clarifying Questions.....	31
Confirmed Facts & Assumptions.....	31
What Teams Often Assume vs Reality.....	32
Root Cause Analysis.....	32
Final Resolution.....	34
Key Learnings.....	34
Core Principle Reinforced.....	34

Scenario 7.....	35
DAG Breaks After Airflow Version Upgrade.....	35
Problem Statement.....	35
Expected vs Actual Behavior.....	35
Why This Happens Frequently.....	36
Clarifying Questions.....	36
Confirmed Facts & Assumptions.....	36
What Teams Often Assume vs Reality.....	37
Root Cause Analysis.....	37
Final Resolution.....	39
Key Learnings.....	39
Core Principle Reinforced.....	39
Scenario 8.....	40
Airflow DAG Fails Intermittently Due to Transient API Errors.....	40
Problem Statement.....	40
Expected vs Actual Behavior.....	40
Why This Happens Frequently.....	41
Clarifying Questions.....	41
Confirmed Facts & Assumptions.....	41
What Teams Often Assume vs Reality.....	42
Root Cause Analysis.....	42
Final Resolution.....	44
Key Learnings.....	44
Core Principle Reinforced.....	44
Scenario 9.....	45
DAG Fails Due to Misconfigured Airflow Connection.....	45
Problem Statement.....	45
Expected vs Actual Behavior.....	45
Why This Happens Frequently.....	46
Clarifying Questions.....	46
Confirmed Facts & Assumptions.....	46
What Teams Often Assume vs Reality.....	47
Root Cause Analysis.....	47
Final Resolution.....	49
Key Learnings.....	49
Core Principle Reinforced.....	49

Scenario 10.....	50
DAG Backfill Saturates Workers and Delays Daily Jobs.....	50
Problem Statement.....	50
Expected vs Actual Behavior.....	50
Why This Happens Frequently.....	51
Clarifying Questions.....	51
Confirmed Facts & Assumptions.....	51
What Teams Often Assume vs Reality.....	52
Root Cause Analysis.....	52
Final Resolution.....	54
Key Learnings.....	54
Core Principle Reinforced.....	54



Scenario 1

Airflow DAG Tasks Stuck in Queued State

Problem Statement

Multiple tasks in an **Airflow DAG** remain in the “queued” state for hours, even though they are **independent**. As a result, **downstream jobs are delayed**, putting the **1-hour SLA at risk**.

Key Details

- Tasks stuck in **queued** state
- DAG tasks are independent (no true dependency bottleneck)
- Workers are limited
- SLA: 1 hour
- Downstream pipelines blocked

Expected vs Actual Behavior

Expected	Actual
Tasks picked up quickly	Tasks wait in queue
Independent tasks run in parallel	Parallelism unused
Workers actively executing	Workers saturated
SLA met	SLA breached

This is a **scheduler and capacity issue**, not a DAG logic problem.

Why This Happens Frequently

Because:

- Worker capacity is underestimated
- Default parallelism limits are left unchanged
- More DAGs are added over time without scaling workers

Teams often misdiagnose this as:

- A DAG dependency issue
- A task failure
- A transient scheduler glitch

But **queued tasks usually mean “no worker available,” not “task problem.”**

Clarifying Questions

A senior data engineer asks:

- How many workers are available?
- What is the current worker utilization?
- Are multiple DAGs competing for the same pool?
- Is task parallelism limited by config?
- Do task priorities exist?

These questions focus on **execution capacity**, not task retries.

Confirmed Facts & Assumptions

After investigation:

- Tasks are independent
- No task failures
- Workers fully occupied
- Scheduler functioning normally
- Increasing workers is feasible

Interpretation:

This is a **resource bottleneck**, not an orchestration bug.

What Teams Often Assume vs Reality

Assumption	Reality
Reordering DAG will fix it	Workers are still limited
Retrying helps	Retries re-enter the same queue
Scheduler is slow	Scheduler is waiting for workers
Ignore for now	SLA impact increases

Airflow can only schedule what **workers can execute**.

Root Cause Analysis

Step 1: Observe Task States

- Tasks stuck in **queued**
- No failures or retries

Conclusion:

Scheduler is healthy; workers are the constraint.

Step 2: Analyze Worker Capacity

- Worker slots exhausted
- Multiple DAGs competing

This confirms **insufficient execution capacity**.

Step 3: Conceptual Root Cause

The root cause is **under-provisioned Airflow workers**:

- DAG growth without scaling
- No pool or priority management

This is a **scaling oversight**, not a design flaw.

Step 4 :Wrong Approach vs Right Approach

Wrong Approach

- Retry tasks repeatedly
- Reorder DAG unnecessarily
- Ignore queue buildup

Right Approach

- Increase number of workers
- Scale executor appropriately
- Use pools/priorities if needed

Senior engineers **scale execution, not workarounds.**

Step 5 : Validation of the Fix

To validate:

- Add workers
- Monitor queue length
- Track task start latency

Outcome:

Tasks move from queued → running immediately.

Step 6 : Corrective Actions

- Increase Airflow worker count
- Tune parallelism and concurrency
- Use task pools for isolation
- Monitor queue wait time
- Capacity-plan as DAGs grow

These steps ensure **SLA stability at scale.**

Step 7 : Result After Fix

Before	After
Tasks stuck queued	Tasks execute promptly
Downstream delays	Smooth DAG execution
SLA breaches	SLA met
Reactive firefighting	Predictable orchestration

Final Resolution

- **Root Issue:** Insufficient Airflow worker capacity
- **Action Taken:** Scaled workers to reduce queueing

Key Learnings

- Queued ≠ failed
- Airflow performance is capacity-driven
- Independent tasks still need workers
- Scaling workers is often the simplest fix

Core Principle Reinforced

Schedulers don't run jobs—workers do.

■ ■ ■

Scenario 2

Airflow DAG Fails Due to Circular Dependencies

Problem Statement

A newly deployed **Airflow DAG fails to start** because tasks **indirectly reference each other**, triggering **cycle detection errors**. Since the DAG supports **multiple dependent pipelines** and must run **daily**, the failure threatens a **2-hour SLA**.

Key Details

- New DAG deployment
- Tasks form an indirect dependency loop
- Airflow cycle detection error
- DAG supports multiple downstream pipelines
- SLA: 2 hours

Expected vs Actual Behavior

Expected	Actual
DAG loads successfully	DAG fails at parse time
Tasks execute in order	Scheduler blocks execution
Daily runs proceed	No runs triggered
Downstream pipelines run	Pipelines blocked

This is a **DAG design error**, not an execution or resource issue.

Why This Happens Often

Because:

- DAGs grow incrementally over time
- New dependencies are added without re-evaluating the full graph
- Indirect cycles are harder to spot visually

Common misconceptions:

- “Retry will fix it”
- “Scheduler is unstable”
- “Splitting DAG will magically help”

But **Airflow enforces strict acyclic graphs at parse time.**

Clarifying Questions

A senior data engineer asks:

- Which tasks reference each other?
- Is the dependency logical or historical?
- Can tasks be decoupled or reordered?
- Are control dependencies mixed with data dependencies?
- Can sensors or external triggers replace direct links?

These questions focus on **graph correctness**, not runtime fixes.

Confirmed Facts & Assumptions

After inspection:

- Tasks form an indirect loop ($A \rightarrow B \rightarrow C \rightarrow A$)
- No runtime execution occurs
- Retry has no effect
- DAG splitting without logic change won't help
- Dependency cycle can be removed safely

Interpretation:

This is a **logical design flaw in the DAG structure.**

What Teams Often Assume vs Reality

Assumption	Reality
Scheduler bug caused failure	DAG graph is invalid
Retrying will help	DAG never starts
Splitting DAG fixes cycles	Logic still cycles
Ignoring is acceptable	DAG remains broken

Airflow fails **fast and correctly** when DAGs are cyclic.

Root Cause Analysis

Step 1: Inspect DAG Graph

Observed:

- Tasks indirectly depend on each other
- Cycle detected at parse time

Conclusion:

DAG violates acyclic requirement.

Step 2: Identify Dependency Intent

Observed:

- Some dependencies added for sequencing, not necessity
- Logical separation possible

This confirms **cycle is unnecessary**.

Step 3: Conceptual Root Cause

The root cause is **improper DAG dependency design**:

- Data and control dependencies mixed
- No full-graph review after changes

This is a **design governance gap**, not a tooling issue.

Step 4 :Wrong Approach vs Right Approach

Wrong Approach

- Retry the DAG
- Ignore errors
- Split DAG without fixing logic

Right Approach

- Break the dependency cycle
- Redesign task order
- Use sensors or external triggers if needed

Senior engineers **fix graph logic, not symptoms.**

Step 5 : Validation of the Fix

To validate:

- Remove cyclic dependency
- Reload DAG
- Confirm scheduler parses successfully
- Trigger manual run

Outcome:

DAG loads and executes normally.

Step 6 : Corrective Actions

- Review DAG dependency graph end-to-end
- Enforce DAG design reviews
- Separate data vs control dependencies
- Use sensors instead of direct cycles
- Add DAG validation checks in CI/CD

These steps prevent **future cycle-related outages.**

Step 7 : Result After Fix

Before	After
DAG fails to load	DAG runs successfully
Scheduler blocked	Scheduler healthy
Downstream pipelines stuck	Pipelines unblocked
Repeated failures	Stable daily execution

Final Resolution

- **Root Issue:** Circular dependencies in DAG design
- **Action Taken:** Broke dependency cycle and redesigned DAG

Key Learnings

- Airflow DAGs must be acyclic
- Cycles fail at parse time, not runtime
- Retry never fixes design errors
- DAG design discipline matters at scale

Core Principle Reinforced

If your DAG has a cycle, Airflow will stop you—and it's right to do so.



Scenario 3

Timezone Misconfiguration Delays Airflow DAG Runs

Problem Statement

An Airflow DAG is configured to run in **UTC**, but the team expects execution based on **IST business hours**. As a result, DAGs run later than expected, causing **downstream processing delays** and **SLA risk**.

Key Details

- DAG timezone set to UTC
- Business expectation based on IST
- Daily reports tied to business hours
- Multiple DAGs affected
- SLA: 1 hour

Expected vs Actual Behavior

Expected	Actual
DAG runs during IST business hours	DAG runs late due to UTC offset
Reports available on time	Reports delayed
SLA consistently met	SLA frequently breached
Predictable scheduling	Confusing execution times

This is a **scheduling configuration issue**, not a task or capacity failure.

Why This Happens Frequently

Because:

- Airflow defaults to UTC
- Teams assume local timezone implicitly
- Cron expressions look correct but shift with timezone
- Timezone validation is often skipped

Common reactions:

- Manually triggering DAGs
- Retrying jobs
- Ignoring delays as “minor”

But **manual workarounds hide a systemic scheduling bug.**

Clarifying Questions

A senior data engineer asks:

- What timezone is the DAG explicitly set to?
- What timezone do stakeholders expect?
- Are **start_date** and **schedule_interval** aligned?
- Are multiple DAGs sharing the same assumption?
- Are reports tied to business or system time?

These questions focus on **time semantics**, not execution logic.

Confirmed Facts & Assumptions

After review:

- DAG timezone is UTC
- Business expects IST-aligned execution
- No explicit timezone override configured
- Manual triggers used as workaround
- Timezone correction is safe

Interpretation:

This is a **misalignment between system defaults and business expectations**.

What Teams Often Assume vs Reality

Assumption	Reality
Cron schedule is self-explanatory	Cron depends on timezone
Retry fixes delays	Retry runs in same timezone
Manual trigger is fine	Manual work is not scalable
One DAG affected	All similar DAGs affected

Time errors scale silently across pipelines.

Root Cause Analysis

Step 1: Inspect DAG Configuration

Observed:

- No timezone explicitly set
- Airflow defaulting to UTC

Conclusion:

Execution time offset introduced unintentionally.

Step 2: Map Business Expectation

Observed:

- Reports expected during IST business hours
- Downstream dependencies assume earlier completion

This confirms **timezone mismatch**.

Step 3: Conceptual Root Cause

The root cause is **implicit timezone assumptions in DAG design**:

- UTC default not documented

- Business time not encoded in configuration

This is a **configuration governance gap**.

Step 4 :Wrong Approach vs Right Approach

Wrong Approach

- Ignore delays
- Manually trigger DAGs
- Retry jobs

Right Approach

- Explicitly set DAG timezone
- Align cron schedule with business hours
- Validate `start_date` and schedule

Senior engineers **encode time expectations explicitly**.

Step 5 : Validation of the Fix

To validate:

- Update DAG timezone to IST
- Trigger test run
- Confirm execution aligns with business hours
- Monitor downstream SLAs

Outcome:

DAG runs predictably at expected times.

Step 6 : Corrective Actions

- Set timezone explicitly in all DAGs
- Standardize timezone conventions
- Document business-time expectations
- Validate schedules during DAG reviews
- Monitor execution-time drift

These steps prevent **recurring time-based SLA failures**.

Step 7 : Result After Fix

Before	After
Late DAG runs	On-time execution
Manual triggers	Automated scheduling
SLA breaches	SLA stability
Confusing schedules	Predictable pipelines

Final Resolution

- **Root Issue:** DAG timezone misconfiguration
- **Action Taken:** Adjusted DAG timezone to match business expectations

Key Learnings

- Timezones matter in orchestration
- Cron schedules are timezone-dependent
- Defaults are rarely business-friendly
- Explicit configuration prevents confusion

Core Principle Reinforced

If time matters to the business, encode it explicitly in your DAG.



Scenario 4

Airflow Sensor Tasks Time Out and Delay DAGs

Problem Statement

Airflow **S3 sensor tasks frequently time out**, even though the expected data eventually arrives. These timeouts **delay downstream ETL jobs** and put the **1-hour SLA** at risk. The data arrival pattern is **unpredictable**, but the sensor is critical for correctness.

Key Details

- S3 sensors timing out
- Data arrival unpredictable
- Sensor blocks downstream ETL
- SLA: 1 hour
- Retries cause repeated delays

Expected vs Actual Behavior

Expected	Actual
Sensor waits until data arrives	Sensor times out early
Downstream jobs triggered on data	Jobs delayed by sensor failure
SLA consistently met	SLA breached
Sensor failures rare	Frequent timeouts

This is a **sensor configuration issue**, not a data availability problem.

Why This Happens Frequently

Because:

- Default sensor timeouts are conservative
- Data arrival times vary across days
- Sensors are configured without historical arrival analysis

Common mistakes:

- Retrying the DAG
- Ignoring intermittent failures
- Reducing poke frequency without adjusting timeout

But **sensor defaults rarely match real-world data latency.**

Clarifying Questions

A senior data engineer asks:

- What is the historical data arrival distribution?
- How long does data usually take to arrive?
- Is the sensor running in poke or reschedule mode?
- What is the current timeout vs expected delay?
- Is the sensor blocking critical paths?

These questions focus on **alignment with real data behavior.**

Confirmed Facts & Assumptions

After analysis:

- Data sometimes arrives later than default timeout
- Sensor configuration unchanged since pipeline creation
- Increasing timeout does not break logic
- Poke interval is reasonable
- Sensor failure causes unnecessary retries

Interpretation:

This is a mismatch between sensor configuration and data arrival patterns.

What Teams Often Assume vs Reality

Assumption	Reality
Sensor timeout equals SLA	Timeout should match arrival variance
Retry will eventually work	Retries repeat the same config
Polling frequency is the issue	Timeout is the real bottleneck
Sensor failures indicate data issues	Often config issues

Sensors must be **tuned to the data, not the clock.**

Root Cause Analysis

Step 1: Inspect Sensor Config

Observed:

- Timeout shorter than max arrival delay
- Sensor configured long before traffic growth

Conclusion:

Sensor is timing out prematurely.

Step 2: Evaluate Impact of Timeout Change

Observed:

- Increasing timeout allows sensor to wait safely
- No adverse impact on resource usage

This confirms **timeout tuning is safe and effective.**

Step 3: Conceptual Root Cause

The root cause is **untuned sensor parameters:**

- Timeout not aligned with data latency
- Historical arrival patterns ignored

This is a **configuration maturity gap**.

Step 4 :Wrong Approach vs Right Approach

Wrong Approach

- Ignore intermittent timeouts
- Retry jobs repeatedly
- Reduce polling blindly

Right Approach

- Increase sensor timeout
- Tune poke interval based on arrival data
- Use reschedule mode if needed

Senior engineers **tune sensors, not firefight failures.**

Step 5 : Validation of the Fix

To validate:

- Increase sensor timeout
- Monitor success rate
- Track downstream SLA adherence

Outcome:

Sensor waits correctly, downstream jobs start on time.

Step 6 : Corrective Actions

- Analyze historical data arrival times
- Tune sensor timeout and poke interval
- Prefer reschedule mode for long waits
- Alert on abnormal delay
- Review sensor configs periodically

These steps prevent **chronic sensor-induced delays**.

Step 7 : Result After Fix

Before	After
Frequent sensor timeouts	Stable sensor behavior
Downstream delays	On-time execution
SLA breaches	SLA compliance
Reactive retries	Predictable orchestration

Final Resolution

- **Root Issue:** Sensor timeout misconfiguration
- **Action Taken:** Increased sensor timeout and tuned configuration

Key Learnings

- Sensors are configuration-sensitive
- Defaults rarely fit production data
- Timeouts should reflect data behavior
- Tuning prevents false failures

Core Principle Reinforced

A sensor should wait for data—not give up before it arrives.



Scenario 5

Resource Contention Between Multiple Airflow DAGs

Problem Statement

Multiple Airflow DAGs are running on a **shared worker pool**. A **business-critical DAG** is delayed because **lower-priority DAGs consume available workers**, causing the **1-hour SLA** to be at risk.

Key Details

- Multiple DAGs sharing the same workers
- Critical DAG delayed
- Jobs are time-sensitive
- Shared Airflow cluster
- SLA: 1 hour

Expected vs Actual Behavior

Expected	Actual
Critical DAG runs on time	Critical DAG waits
Resources allocated by priority	First-come tasks dominate
SLA protected	SLA breached
Predictable execution order	Uncontrolled contention

This is a **resource allocation and prioritization issue**, not a task or dependency failure.

Why This Happens Frequently

Because:

- All DAGs share the same worker pool by default
- Task priorities are not explicitly set
- New DAGs are added without revisiting capacity planning

Teams often react by:

- Retrying delayed DAGs
- Ignoring delays temporarily
- Blaming scheduler instability

But **Airflow executes what workers are free to run—not what's most important.**

Clarifying Questions

A senior data engineer asks:

- Are task priorities defined across DAGs?
- Do critical DAGs have dedicated pools?
- Is worker capacity sized for peak load?
- Are SLAs aligned with priority levels?
- Which DAGs can tolerate delays?

These questions focus on **intentional resource governance**, not firefighting.

Confirmed Facts & Assumptions

After investigation:

- Workers are fully utilized
- No task priorities or pools configured
- Critical DAG has no reserved capacity
- Scaling workers is feasible
- Airflow supports priority weights and pools

Interpretation:

This is a **lack of prioritization strategy**, not a scaling failure alone.

What Teams Often Assume vs Reality

Assumption	Reality
Airflow knows what's critical	Airflow needs explicit priorities
Retrying helps	Retries re-enter the same queue
All DAGs are equal	Business impact differs
Scaling is optional	Capacity must match priority

Schedulers enforce rules—not business importance.

Root Cause Analysis

Step 1: Observe Execution Order

Observed:

- Low-priority tasks occupy workers
- Critical DAG remains queued

Conclusion:

Worker slots are consumed indiscriminately.

Step 2: Evaluate Resource Controls

Observed:

- No task priorities
- No pools for isolation

This confirms **missing resource governance**.

Step 3: Conceptual Root Cause

The root cause is **absence of prioritization and capacity isolation**:

- All DAGs treated equally
- No safeguards for critical workflows

This is a **multi-tenant orchestration design gap**.

Step 4 :Wrong Approach vs Right Approach

Wrong Approach

- Ignore contention
- Retry delayed DAGs
- Hope scheduling improves

Right Approach

- Increase worker pool capacity
- Assign task priorities
- Use pools to isolate critical DAGs

Senior engineers **encode business importance into orchestration**.

Step 5 : Validation of the Fix

To validate:

- Increase workers or configure pools
- Assign higher priority to critical DAG
- Monitor queue wait time

Outcome:

Critical DAG executes on time even under load.

Step 6 : Corrective Actions

- Define DAG and task priority standards
- Create dedicated pools for critical workflows
- Scale workers for peak concurrency
- Monitor contention metrics
- Review priorities as pipelines grow

These steps ensure **SLAs are protected in shared environments**.

Step 7 : Result After Fix

Before	After
Critical DAG delayed	Critical DAG prioritized
Worker contention	Controlled allocation
SLA breaches	SLA compliance
Reactive retries	Predictable scheduling

Final Resolution

- **Root Issue:** Resource contention between DAGs
- **Action Taken:** Increased worker capacity and applied task prioritization

Key Learnings

- Shared clusters need explicit prioritization
- Airflow schedules based on capacity, not importance
- Pools and priorities are essential at scale
- SLA protection requires proactive design

Core Principle Reinforced

If everything is high priority, nothing is. Encode priority into your DAGs.

■ ■ ■

Scenario 6

DAG Completes Successfully but Tasks Are Skipped

Problem Statement

An Airflow DAG **shows a successful completion**, but several downstream tasks are **skipped** because upstream tasks failed and the **trigger rules were not aligned with business logic**. As a result, **data is incomplete**, even though the DAG appears healthy. Retries are limited, and downstream pipelines are **business-critical**.

Key Details

- Tasks skipped due to trigger rules
- DAG marked as successful
- Data output incomplete
- Downstream pipelines impacted
- SLA: 2 hours

Expected vs Actual Behavior

Expected	Actual
DAG success implies full data	DAG success hides skipped tasks
Business-critical tasks always run	Tasks skipped silently
Downstream pipelines get complete data	Partial data consumed
SLA reflects correctness	SLA met with incorrect output

This is a **logic and configuration problem**, not a scheduling or capacity issue.

Why This Problem Is Dangerous

Because:

- Skipped tasks do **not raise obvious failures**
- Dashboards and alerts may show green
- Downstream teams assume data is complete
- Errors surface much later in analytics

Common misconceptions:

- “DAG succeeded, so data must be fine”
- “Retries will handle it”
- “Skipping is acceptable”

But **skipped tasks are silent data failures**.

Clarifying Questions

A senior data engineer asks:

- What trigger rule is applied (**all_success, one_success, all_done**)?
- Are skipped tasks actually optional?
- What business logic requires these tasks to run?
- Should downstream tasks tolerate upstream failures?
- How is data completeness validated?

These questions focus on **business intent**, not just DAG status.

Confirmed Facts & Assumptions

After investigation:

- Default trigger rules caused skips
- Tasks should have executed despite upstream failure
- DAG status misled monitoring
- Retrying repeats the same behavior
- Trigger rules can be safely adjusted

Interpretation:

This is a misalignment between Airflow semantics and business requirements.

What Teams Often Assume vs Reality

Assumption	Reality
Green DAG = correct data	Green DAG can hide skips
Default trigger rules are safe	Defaults may not fit logic
Skipped tasks are harmless	Skips cause data loss
Retries fix logic issues	Retries repeat the same logic

Airflow does exactly what you tell it—not what you mean.

Root Cause Analysis

Step 1: Inspect Trigger Rules

Observed:

- Tasks configured with default **all_success**
- Upstream failure caused downstream skips

Conclusion:

Trigger rules are too strict for business needs.

Step 2: Map Business Expectations

Observed:

- Certain tasks must run even if upstream partially fails
- Data completeness more important than strict dependency success

This confirms **trigger rules must be adjusted**.

Step 3: Conceptual Root Cause

The root cause is **misconfigured trigger rules**:

- Technical defaults override business intent

- DAG success not aligned with data correctness

This is a **workflow design gap**, not a runtime issue.

Step 4 :Wrong Approach vs Right Approach

Wrong Approach

- Ignore skipped tasks
- Retry the DAG
- Skip failed tasks manually

Right Approach

- Adjust trigger rules to match business logic
- Use `all_done` or conditional branching where appropriate
- Explicitly encode correctness rules

Senior engineers **design DAGs for correctness, not just completion.**

Step 5 : Validation of the Fix

To validate:

- Update trigger rules
- Run DAG with simulated upstream failure
- Confirm all required tasks execute
- Verify downstream data completeness

Outcome:

DAG behavior matches business expectations.

Step 6 : Corrective Actions

- Review trigger rules during DAG design
- Add data completeness checks
- Alert on skipped critical tasks
- Document business intent in DAG code
- Avoid relying on defaults blindly

These steps prevent **silent data corruption**.

Step 7 : Result After Fix

Before	After
DAG green, data wrong	DAG behavior correct
Skipped tasks unnoticed	Tasks run as intended
Downstream issues	Downstream stability
False confidence	Reliable pipelines

Final Resolution

- **Root Issue:** Incorrect trigger rules causing skipped tasks
- **Action Taken:** Updated trigger rules to align with business logic

Key Learnings

- DAG success ≠ data correctness
- Trigger rules are critical design decisions
- Defaults are rarely sufficient
- Silent failures are the most dangerous

Core Principle Reinforced

A DAG that “succeeds” with missing data has failed its real purpose.



Scenario 7

DAG Breaks After Airflow Version Upgrade

Problem Statement

After an **Airflow version upgrade**, a **business-critical DAG fails to run** because it uses **deprecated operators** that are no longer supported. The upgrade **cannot be rolled back immediately**, and the **2-hour SLA** is at risk.

Key Details

- Recent Airflow upgrade
- DAG fails at parse or runtime
- Deprecated operators in use
- Rollback not immediately possible
- SLA: 2 hours

Expected vs Actual Behavior

Expected	Actual
DAG runs after upgrade	DAG fails due to incompatibility
Backward compatibility maintained	Deprecated operators removed
SLA preserved	SLA breached
Smooth upgrade	Production outage

This is a **compatibility and upgrade readiness issue**, not a scheduling or data problem.

Why This Happens Frequently

Because:

- Airflow deprecates operators across major versions
- DAGs are not tested against new versions in advance
- Deprecation warnings are ignored during previous runs
- Upgrade focuses on infra, not DAG code

Common reactions:

- Retry jobs
- Blame scheduler instability
- Attempt risky downgrades

But **Airflow upgrades change APIs, not just binaries.**

Clarifying Questions

A senior data engineer asks:

- Which operators were deprecated in this version?
- Did logs show deprecation warnings earlier?
- Is the failure at parse time or runtime?
- Are multiple DAGs using the same operators?
- Is there a compatibility matrix available?

These questions focus on **version-aware DAG design**.

Confirmed Facts & Assumptions

After investigation:

- DAG uses deprecated operators
- Failure reproducible across retries
- Rollback not feasible immediately
- Updated operators are available
- Code change is safe and scoped

Interpretation:

This is a **DAG code compatibility issue introduced by the upgrade**.

What Teams Often Assume vs Reality

Assumption	Reality
Upgrade only affects infra	DAG code must also be compatible
Retry will fix it	Code incompatibility persists
Downgrade is safest	Downgrades carry risk
Deprecation warnings are optional	They are early failure signals

Upgrades punish ignored warnings.

Root Cause Analysis

Step 1: Inspect Failure Logs

Observed:

- Errors referencing deprecated operators
- DAG fails to load or execute

Conclusion:

Operator removal caused DAG failure.

Step 2: Review DAG Code

Observed:

- Operators deprecated in previous versions
- No migration performed

This confirms **missed upgrade preparation**.

Step 3: Conceptual Root Cause

The root cause is **lack of upgrade readiness**:

- DAGs not validated against new Airflow version

- Deprecation warnings ignored

This is a **release hygiene gap**, not a runtime issue.

Step 4 :Wrong Approach vs Right Approach

Wrong Approach

- Ignore failure
- Retry jobs
- Downgrade Airflow hastily

Right Approach

- Update DAGs to supported operators
- Follow Airflow migration guides
- Test DAGs pre-upgrade

Senior engineers **treat upgrades as code changes, not infra events.**

Step 5 : Validation of the Fix

To validate:

- Replace deprecated operators
- Reload DAG
- Trigger test run
- Monitor SLA compliance

Outcome:

DAG runs successfully on the upgraded Airflow version.

Step 6 : Corrective Actions

- Update DAGs to supported operators
- Review deprecation warnings regularly
- Add DAG compatibility tests
- Maintain upgrade checklists
- Validate DAGs in staging before prod upgrades

These steps prevent **upgrade-induced outages**.

Step 7 : Result After Fix

Before	After
DAG broken post-upgrade	DAG compatible
SLA breaches	SLA restored
Firefighting	Predictable upgrades
Risky rollbacks	Controlled migrations

Final Resolution

- **Root Issue:** Deprecated operators after Airflow upgrade
- **Action Taken:** Updated DAG to supported operators

Key Learnings

- Airflow upgrades affect DAG code
- Deprecation warnings matter
- Retries don't fix compatibility
- Testing DAGs pre-upgrade is essential

Core Principle Reinforced

Infrastructure upgrades fail pipelines only when code is not upgrade-ready.



Scenario 8

Airflow DAG Fails Intermittently Due to Transient API Errors

Problem Statement

An Airflow DAG **fails unpredictably** because it depends on an **upstream external API** that occasionally returns transient errors (timeouts, 5xx). The failures are **non-deterministic**, downstream jobs are **business-critical**, and the **1-hour SLA** is at risk.

Key Details

- Intermittent DAG failures
- External API dependency
- Errors are transient, not persistent
- Downstream jobs blocked
- SLA: 1 hour

Expected vs Actual Behavior

Expected	Actual
DAG handles temporary API issues	DAG fails immediately
Transient errors auto-recovered	Manual intervention required
Downstream jobs run reliably	Downstream jobs blocked
SLA consistently met	SLA intermittently breached

This is a **resiliency and error-handling issue**, not a logic or scheduling problem.

Why This Happens Frequently

Because:

- External APIs are inherently unreliable
- Network hiccups and rate limits are common
- DAGs are often written assuming “happy path” execution

Common reactions:

- Manually retry DAGs
- Ignore intermittent failures
- Skip failing tasks

But **intermittent failures are predictable in distributed systems.**

Clarifying Questions

A senior data engineer asks:

- Are failures transient or deterministic?
- What error codes are returned (timeouts, 429, 5xx)?
- Is retry logic already configured?
- Should retries block downstream tasks?
- Is exponential backoff implemented?

These questions focus on **fault tolerance**, not blame.

Confirmed Facts & Assumptions

After analysis:

- API failures are transient
- No retry logic configured
- Manual retries usually succeed
- Skipping tasks breaks data correctness
- Retry configuration is safe to add

Interpretation:

This is a **missing resiliency pattern**, not an unstable pipeline.

What Teams Often Assume vs Reality

Assumption	Reality
API should be stable	External systems fail
Retry is manual work	Retry should be automated
Skipping is acceptable	Skipping breaks correctness
Random failures are unavoidable	Resilience can be engineered

Intermittent ≠ unavoidable.

Root Cause Analysis

Step 1: Inspect Failure Pattern

Observed:

- Failures vary run to run
- Same task succeeds on retry

Conclusion:

Failures are transient.

Step 2: Review DAG Error Handling

Observed:

- No retries configured
- Immediate task failure on first error

This confirms **lack of retry strategy**.

Step 3: Conceptual Root Cause

The root cause is **absence of fault tolerance for external dependencies**:

- No retries
- No exponential backoff

- No graceful recovery

This is a **resilience design gap**.

Step 4 :Wrong Approach vs Right Approach

Wrong Approach

- Ignore intermittent failures
- Retry DAG manually
- Skip failing tasks

Right Approach

- Implement task-level retries
- Use exponential backoff
- Set retry limits based on SLA

Senior engineers **design for failure, not perfection**.

Step 5 : Validation of the Fix

To validate:

- Add retries with backoff
- Simulate transient API failures
- Observe automatic recovery
- Monitor SLA compliance

Outcome:

DAG absorbs transient failures and completes successfully.

Step 6 : Corrective Actions

- Configure retries on API-dependent tasks
- Use exponential backoff
- Set sensible retry limits
- Add alerts only after retries exhausted
- Document external dependency behavior

These steps ensure **predictable DAG execution under unreliable conditions**.

Step 7 : Result After Fix

Before	After
Random DAG failures	Stable execution
Manual retries	Automatic recovery
SLA breaches	SLA protected
Operator intervention	Self-healing pipeline

Final Resolution

- **Root Issue:** No retry handling for transient API failures
- **Action Taken:** Implemented retries with exponential backoff

Key Learnings

- External dependencies always fail sometimes
- Retries are a first-class design pattern
- Manual recovery does not scale
- Reliability is engineered, not hoped for

Core Principle Reinforced

If a failure is transient, your system should be too.



Scenario 9

DAG Fails Due to Misconfigured Airflow Connection

Problem Statement

A business-critical **ETL DAG fails to access S3/Redshift** because the **Airflow connection contains invalid or expired credentials**. Since **multiple DAGs share this connection**, the failure cascades quickly and threatens the **1-hour SLA**.

Key Details

- DAG fails at connection step
- Invalid or expired credentials
- Connection shared across multiple DAGs
- Jobs are business-critical
- SLA: 1 hour

Expected vs Actual Behavior

Expected	Actual
DAG connects to S3/Redshift	Authentication failure
Shared connection works for all DAGs	Multiple DAGs fail simultaneously
SLA met	SLA breached
Minimal operational noise	Widespread pipeline failures

This is a **configuration and credential management issue**, not a data or orchestration problem.

Why This Happens Frequently

Because:

- Credentials expire or rotate silently
- Same connection reused across environments
- Manual updates introduce inconsistencies
- No alerting on credential validity

Common reactions:

- Retrying DAGs
- Manually copying data
- Blaming infra instability

But **no retry can fix invalid credentials.**

Clarifying Questions

A senior data engineer asks:

- Did credentials recently rotate or expire?
- Is this connection shared across DAGs?
- Is the failure authentication-related or network-related?
- Are secrets managed centrally?
- Do non-prod and prod use separate connections?

These questions isolate **configuration scope**, not task logic.

Confirmed Facts & Assumptions

After investigation:

- Credentials are invalid/expired
- All DAGs using this connection fail
- Retrying does not help
- Updating credentials is safe
- Centralized secret storage is available

Interpretation:

This is a **single-point configuration failure with wide blast radius**.

What Teams Often Assume vs Reality

Assumption	Reality
Retry might fix it	Credentials remain invalid
One DAG issue	All dependent DAGs affected
Manual data copy helps	Breaks automation and trust
Ignore briefly	SLA impact compounds

Connection failures fail **fast and loudly**.

Root Cause Analysis

Step 1: Inspect Error Logs

Observed:

- Authentication/authorization errors
- Consistent failures across DAGs

Conclusion:

Connection credentials are invalid.

Step 2: Trace Connection Usage

Observed:

- Same Airflow connection referenced by multiple DAGs

This confirms **shared credential dependency**.

Step 3: Conceptual Root Cause

The root cause is **mismanaged credentials**:

- No proactive rotation validation
- Central connection used without safeguards

This is a **secrets governance gap**, not a DAG bug.

Step 4 :Wrong Approach vs Right Approach

Wrong Approach

- Retry jobs
- Manually copy data
- Ignore failures

Right Approach

- Update credentials in Airflow connection
- Validate access immediately
- Centralize and monitor credential rotation

Senior engineers **fix shared configuration first.**

Step 5 : Validation of the Fix

To validate:

- Update credentials
- Test connection
- Trigger one DAG
- Confirm dependent DAGs recover

Outcome:

All DAGs regain access and execute successfully.

Step 6 : Corrective Actions

- Centralize credential management (Secrets Manager / Vault)
- Separate connections per environment
- Monitor credential expiry
- Alert on authentication failures
- Document shared connection dependencies

These steps reduce **blast radius from configuration errors.**

Step 7 : Result After Fix

Before	After
Multiple DAG failures	DAGs restored
SLA breaches	SLA met
Manual workarounds	Automated recovery
High operational noise	Stable pipelines

Final Resolution

- **Root Issue:** Invalid credentials in shared Airflow connection
- **Action Taken:** Updated credentials and restored access

Key Learnings

- Shared connections amplify failures
- Retries don't fix auth issues
- Centralized secrets reduce risk
- Configuration errors can be more dangerous than code bugs

Core Principle Reinforced

In orchestration systems, shared credentials create shared failure domains. Manage them carefully.

■ ■ ■

Scenario 10

DAG Backfill Saturates Workers and Delays Daily Jobs

Problem Statement

A **backfill of historical DAG runs** is triggered to recover past data. However, the backfill **consumes most of the available workers**, causing **current daily production jobs to queue** and threatening the **1-hour SLA** for business-critical pipelines.

Key Details

- Large historical backfill triggered
- Shared Airflow cluster
- Worker pool saturated
- Daily production DAGs delayed
- SLA for daily jobs: 1 hour

Expected vs Actual Behavior

Expected	Actual
Backfill runs without impacting production	Backfill starves daily DAGs
Daily jobs prioritized	Daily jobs queued
SLA maintained	SLA breached
Controlled resource usage	Unbounded concurrency

This is a **capacity and scheduling issue**, not a DAG logic failure.

Why This Happens Frequently

Because:

- Backfills default to high concurrency
- Teams treat backfill like normal runs
- No separation between historical and live workloads
- Shared worker pools lack safeguards

Common mistakes:

- Running backfill during peak hours
- Allowing backfill to use full cluster capacity
- Retrying backfill aggressively

But **backfills compete for the same workers as production jobs.**

Clarifying Questions

A senior data engineer asks:

- What is the backfill concurrency?
- Are pools or priorities configured?
- Can backfill run off-peak?
- Are daily jobs protected with higher priority?
- What is the acceptable recovery window?

These questions focus on **protecting production SLAs.**

Confirmed Facts & Assumptions

After investigation:

- Backfill uses default (high) concurrency
- No priority or pool isolation
- Daily jobs are time-sensitive
- Backfill timing is flexible
- Airflow supports concurrency limits and scheduling controls

Interpretation:

This is a **lack of backfill governance**, not insufficient infrastructure.

What Teams Often Assume vs Reality

Assumption	Reality
Backfill is just another DAG run	Backfill multiplies load
More parallelism is faster	It starves critical jobs
Retry helps	Retries worsen contention
Ignore temporarily	SLA damage compounds

Backfills are **production-impacting operations**.

Root Cause Analysis

Step 1: Observe Worker Utilization

Observed:

- Workers fully occupied by backfill tasks
- Daily DAG tasks stuck in **queued**

Conclusion:

Backfill consumed all execution capacity.

Step 2: Inspect Backfill Configuration

Observed:

- No concurrency limits
- No pool or priority separation

This confirms **uncontrolled backfill execution**.

Step 3: Conceptual Root Cause

The root cause is **missing backfill controls**:

- No concurrency throttling
- No off-peak scheduling

- No isolation from production workloads

This is a **capacity planning gap**.

Step 4 :Wrong Approach vs Right Approach

Wrong Approach

- Ignore backfill impact
- Retry backfill during peak
- Let backfill run unrestricted

Right Approach

- Limit backfill concurrency
- Schedule backfill during off-peak hour
- Use pools/priorities to protect daily jobs

Senior engineers **treat backfills as controlled recovery operations**.

Step 5 : Validation of the Fix

To validate:

- Limit backfill concurrency
- Schedule backfill off-peak
- Monitor worker utilization
- Confirm daily DAGs start on time

Outcome:

Backfill progresses steadily without affecting SLAs.

Step 6 : Corrective Actions

- Cap backfill concurrency explicitly
- Run backfills during low-traffic windows
- Use dedicated pools for backfill
- Assign higher priority to daily jobs
- Document backfill runbooks

These steps prevent **production starvation during recovery**.

Step 7 : Result After Fix

Before	After
Workers saturated	Balanced utilization
Daily jobs delayed	Daily jobs on time
SLA breaches	SLA protected
Reactive firefighting	Controlled recovery

Final Resolution

- **Root Issue:** Uncontrolled backfill consuming worker capacity
- **Action Taken:** Limited backfill concurrency and scheduled off-peak execution

Key Learnings

- Backfills are not free operations
- Production SLAs must be protected
- Concurrency controls matter
- Recovery workloads need governance

Core Principle Reinforced

Backfill should recover the past—never break the present.

