

Open in app ↗

Medium

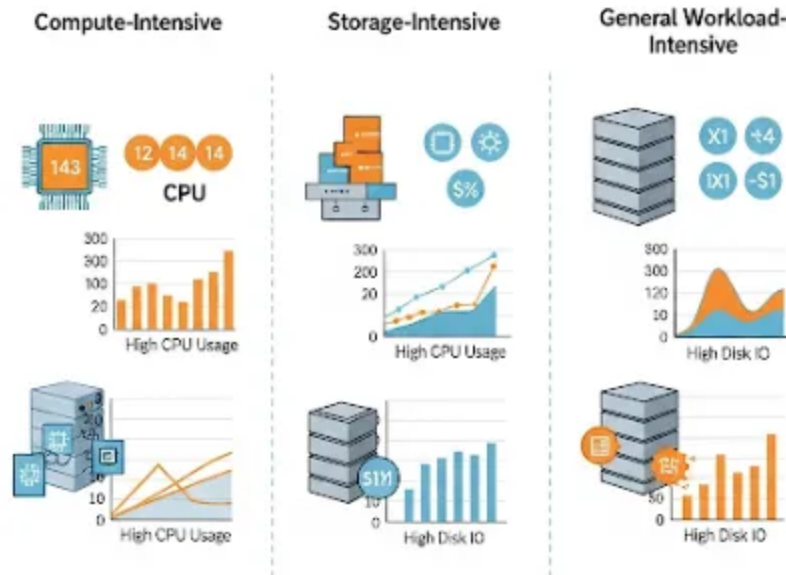
Search

Write

3



★ Get unlimited access to the best of Medium for less than \$1/week. [Become a member](#)



Understanding Whether Your Spark Job Is Compute-Intensive, Storage-Intensive, or General Workload-Intensive



Anurag Sharma 6 min read · Jun 18, 2025

13



When running Apache Spark jobs, optimizing performance and resource allocation hinges on understanding the nature of your workload. Spark jobs can generally be classified as **compute-intensive**, **storage-intensive**, or **general workload-intensive** based on their resource demands. Knowing which category your job falls into helps you allocate resources effectively, avoid bottlenecks, and reduce costs. In this blog, we'll explore how to identify whether your Spark job is compute-intensive, storage-intensive, or general workload-intensive, and provide practical optimization tips.

What Are Compute-Intensive, Storage-Intensive, and General Workload-Intensive Jobs?

Before diving into how to identify these workload types, let's define them:

- **Compute-Intensive:** The job's performance is bottlenecked by CPU usage, often due to complex computations, transformations, or machine learning algorithms.
- **Storage-Intensive:** The job is constrained by disk I/O, memory limitations, or shuffle operations, typically involving large data reads/writes or spills to disk.
- **General Workload-Intensive:** The job has balanced resource demands, requiring moderate CPU, memory, and storage resources without a single dominant bottleneck.

Understanding the type of workload in your job is critical for tuning Spark configurations and allocating resources efficiently.

How to Identify Your Spark Job's Workload Type

To determine whether your Spark job is compute-intensive, storage-intensive, or general workload-intensive, monitor key metrics using the Spark UI, cluster monitoring tools (e.g., Prometheus, Datadog), or logs. Below are the primary indicators for each workload type.

1. Compute-Intensive Jobs

Compute-intensive jobs consume significant CPU resources due to complex processing tasks. Common examples include machine learning model training, heavy data transformations (e.g., map, reduceByKey), or iterative algorithms.

Indicators:

- **High CPU Usage:** Executors show sustained CPU utilization (>80%) across cores, visible in the Spark UI's "Executors" tab or cluster monitoring tools.
- **Low Disk/IO Activity:** Minimal disk spills or shuffle I/O, as the job focuses on in-memory computations.
- **Long Task Execution Times:** Tasks take longer due to computational complexity rather than data movement.
- **High Garbage Collection (GC) Time for CPU-Bound Tasks:** If memory is sufficient but GC is frequent, it may indicate CPU-intensive operations triggering memory churn.
- **Example Workloads:** K-means clustering, PageRank, or jobs with heavy groupBy or join operations on small datasets.

How to Check:

- In the Spark UI, check the "Tasks" section for stages with high CPU time and low input/output bytes.
- Use cluster monitoring tools to confirm high CPU usage and low disk/network I/O.

2. Storage-Intensive Jobs

Storage-intensive jobs are bottlenecked by disk I/O, memory constraints, or network I/O (e.g., during shuffles). These jobs often involve reading/writing large datasets or spilling data to disk due to insufficient memory.

Indicators:

- **High Disk Spills:** Data spills to disk because the allocated memory is insufficient, visible in the Spark UI's "Stages" tab under "Spill (Memory)" or

“Spill (Disk).”

- **High Shuffle I/O:** Large shuffle read/write operations during operations like join, groupByKey, or repartition.
- **Out-of-Memory (OOM) Errors:** Executors crash with OOM errors, indicating insufficient memory for the dataset.
- **High Network I/O:** Significant data transfer between executors during shuffles, visible in the Spark UI or cluster monitoring tools.
- **Example Workloads:** ETL jobs processing terabytes of data, wide transformations (e.g., groupByKey on large datasets), or jobs with heavy disk-based persistence.

How to Check:

- In the Spark UI, look for high “Shuffle Read/Write” or “Disk Spill” metrics in the “Stages” tab.
- Check logs for OOM errors or warnings about excessive shuffle data.
- Monitor disk and network I/O using tools like Ganglia or Prometheus.

3. General Workload-Intensive Jobs

General workload-intensive jobs have balanced resource demands, with no single resource (CPU, memory, or storage) dominating. These jobs require moderate resources across the board and are often less sensitive to specific tuning.

Indicators:

- **Balanced Resource Usage:** CPU, memory, and disk I/O usage are moderate (e.g., 40–60% utilization) with no extreme bottlenecks.
- **No Significant Spills or Errors:** Minimal disk spills and no OOM errors, indicating sufficient memory and storage resources.
- **Moderate Shuffle Activity:** Shuffle operations occur but don’t dominate execution time.
- **Example Workloads:** Simple ETL jobs, lightweight aggregations, or queries on moderately sized datasets.

How to Check:

- In the Spark UI, observe balanced metrics across CPU, memory, and I/O in the “Executors” and “Stages” tabs.
- Use cluster monitoring tools to confirm that no single resource is consistently maxed out.



Practical Example: Diagnosing Workload Type

Let's walk through an example to illustrate how to identify a job's workload type.

Scenario:

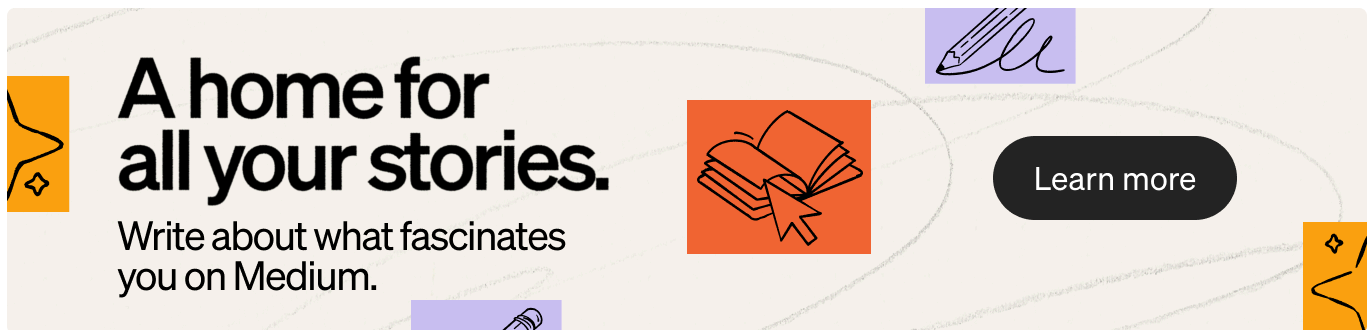
A Spark job processes a 1 TB dataset and takes 90 minutes to complete. The Spark UI shows:

- CPU usage is consistently at 95% across all executors.

- Minimal disk spills (1 GB) and low shuffle read/write (500 MB).
- No OOM errors, and memory usage is at 50%.
- Tasks in one stage take significantly longer due to a complex reduceByKey operation.

Diagnosis:

- **Compute-Intensive:** The high CPU usage, low disk spills, and minimal shuffle activity indicate that the job is bottlenecked by computational complexity (likely the reduceByKey operation).
- **Actionable Insights:**
 - Increase `spark.executor.cores` to provide more CPU resources.
 - Optimize the reduceByKey logic to reduce computational overhead (e.g., pre-aggregate data).
 - Consider increasing `spark.executor.instances` to distribute the workload across more executors.



Alternative Scenario:

If the same job showed high disk spills (50 GB), significant shuffle read/write (100 GB), and OOM errors, it would be classified as **storage-intensive**. In this case, you'd increase `spark.executor.memory`, adjust `spark.memory.fraction`, or optimize shuffle operations.

How to Optimize Based on Workload Type:

Once you've identified the workload type, you can tailor your Spark configuration and job design to optimize performance.

Here are targeted strategies for each type:

1. Optimizing Compute-Intensive Jobs:

- **Increase CPU Resources:** Add more cores per executor (`spark.executor.cores`) or increase the number of executors (`spark.executor.instances`).
- **Optimize Algorithms:** Simplify complex transformations or use more efficient operations (e.g., `reduceByKey` instead of `groupByKey`).
- **Enable Parallelism:** Increase `spark.default.parallelism` or `spark.sql.shuffle.partitions` to distribute computations across more tasks.
- **Use Faster Hardware:** Deploy the job on nodes with higher CPU clock speeds or more cores.

2. Optimizing Storage-Intensive Jobs:

- **Increase Memory:** Boost `spark.executor.memory` to reduce disk spills and OOM errors.
- **Optimize Shuffles:** Increase `spark.shuffle.memoryFraction` or use faster storage (e.g., SSDs) for shuffle data. Reduce shuffles by using operations like `reduceByKey` or broadcasting small datasets.
- **Address Data Skew:** Repartition data evenly (`df.repartition()`) or use salting techniques to balance task workloads.
- **Use Compression:** Enable data compression (`spark.io.compression.codec`) to reduce I/O overhead.

3. Optimizing General Workload-Intensive Jobs:

- **Balance Resources:** Allocate moderate CPU, memory, and storage resources to avoid overprovisioning or underprovisioning.
- **Enable Dynamic Allocation:** Use `spark.dynamicAllocation.enabled` to let Spark adjust executor counts based on workload.
- **Monitor and Fine-Tune:** Regularly check the Spark UI to ensure resources remain balanced as data or job complexity changes.

Tools to Aid Workload Analysis:

To accurately classify and optimize your Spark job's workload, leverage these tools:

- **Spark UI:** Provides detailed metrics on CPU, memory, disk spills, and shuffle activity.
- **Cluster Managers:** YARN, Kubernetes, or Mesos dashboards show cluster-wide resource utilization.
- **Monitoring Tools:** Prometheus, Grafana, or Datadog offer real-time insights into CPU, memory, and I/O metrics.
- **Logs:** Enable verbose logging (`log4j.logger.org.apache.spark=DEBUG`) to capture detailed performance information.

Conclusion:

Determining whether your Spark job is compute-intensive, storage-intensive, or general workload-intensive is essential for optimizing performance and resource allocation.

By monitoring metrics like CPU usage, disk spills, shuffle I/O, and memory consumption in the Spark UI or cluster monitoring tools, you can classify your job's workload and apply targeted optimizations.

Compute-intensive jobs benefit from more CPU resources and optimized algorithms, storage-intensive jobs require increased memory and shuffle tuning, and general workload-intensive jobs need balanced resource allocation.

Regularly analyze your Spark jobs, experiment with configurations, and leverage visualization tools to ensure efficient resource usage.

By understanding your job's workload type, you can achieve faster execution, lower costs, and better cluster utilization.

Data Engineering

Data Engineer

Data Analytics

Automation

Cloud Cost Optimization

**Written by Anurag Sharma**[Edit profile](#)

83 followers · 3 following

Data Engineering Specialist with 10+ exp. Passionate about optimizing pipelines, data lineage, and Spark performance and sharing insights to empower data pros!

No responses yet



Anurag Sharma him/he

What are your thoughts?

More from Anurag Sharma