| RESEARCH ARTICLE

# Taming Spark Data Skew with Practical Solutions

**Srihari Babu Godleti**
*Roku Inc., USA*
**Corresponding Author:** Srihari Babu Godleti, **E-mail**: godleti.srihari@gmail.com

| ABSTRACT

Data skew represents one of the most critical performance challenges in Apache Spark applications, occurring when data is unevenly distributed across partitions and causing significant processing inefficiencies. This technical article explores the nature of data skew in distributed computing environments, its impact on job execution times, and presents three practical solutions for data engineers to implement. Beginning with examining how skew manifests in real-world datasets like e-commerce transactions and social media analytics, the article progresses through increasingly sophisticated mitigation strategies: basic repartitioning for moderate skew cases, key salting techniques for severe distribution imbalances, and broadcast joins for optimizing operations between tables of disparate sizes. Each solution is presented with implementation considerations, performance implications, and appropriate use cases, providing data practitioners with actionable techniques to optimize their Spark jobs regardless of technical background.

| KEYWORDS

Data skew, Apache Spark optimization, partition balancing, key salting, broadcast joins

1. Introduction

Data skew is one of the most common performance bottlenecks in Apache Spark applications, affecting approximately 30-40% of large-scale Spark jobs. This phenomenon occurs when data is unevenly distributed across partitions, causing some executors to become overloaded while others remain underutilized. According to Pramit Marattha et al.'s analysis of Spark performance tuning, skewed data can cause certain partitions to process significantly more data than others, sometimes handling gigabytes while others process only megabytes, resulting in job completion times extending from minutes to hours [1]. This imbalance creates a classic "straggler" problem where the entire job must wait for the slowest executor to complete its oversized workload.

For data engineers and analysts working with big data platforms like Databricks or orchestration tools like Airflow, understanding and resolving data skew is crucial for maintaining efficient processing pipelines. As highlighted in recent research on cloud-based big data processing costs, organizations implementing anti-skew strategies reported an average of 18-22% reduction in their AWS EMR instance hours, translating to significant cost savings in their data lake environments [2]. The study also notes that companies processing financial or retail analytics workloads, which are particularly prone to data skew due to the natural concentration of transactions, have seen processing windows shrink by up to 35% after implementing skew mitigation techniques.

The impact becomes particularly pronounced in production environments where SLAs must be maintained. When processing social media engagement data or e-commerce transactions, popular entities might account for a disproportionate share of records. Pramit Marattha et al. notes that in real-world scenarios, data skew frequently occurs during shuffle-heavy operations

like joins, aggregations, and repartitioning, creating significant resource imbalance [1]. This article breaks down the concept of data skew and provides practical solutions that even those with limited technical backgrounds can implement to optimize their Spark jobs, helping to ensure balanced resource utilization and improved job throughput.

*2. Understanding Data Skew in Apache Spark*

Spark data skew occurs when the distribution of data across worker nodes becomes unbalanced. In a properly functioning Spark cluster, work should be evenly distributed among all executors (the worker nodes that process tasks in parallel). However, when skew occurs, some executors receive significantly more data than others, creating bottlenecks where a small number of executors must process a disproportionate amount of the workload while others sit idle.

Performance bottlenecks in distributed systems like Spark can be particularly challenging to diagnose and resolve. As explained in Computer Systems Performance Analysis literature, these bottlenecks typically manifest as resource contention issues where certain components become overutilized while others remain underutilized [3]. In Spark's distributed computing paradigm, this imbalance creates what performance analysts call "resource starvation patterns" where processing capacity exists in the cluster but cannot be effectively leveraged due to workload distribution issues. Research in performance engineering shows that such bottlenecks can reduce overall system throughput by 65-75% compared to optimally balanced systems, even when the total processing capacity remains unchanged [3].

The root cause of data skew often relates directly to Spark's partitioning mechanism. According to detailed analysis of partition distribution in production environments, the default hash partitioning strategy in Spark can lead to significant imbalances when data contains natural skew [4]. For instance, when partitioning by user ID in social media analytics, some partitions might contain millions of records while others contain only thousands, creating a partition disparity ratio that can exceed 1000:1 in extreme cases. This phenomenon becomes especially pronounced during shuffle operations like joins, aggregations, and repartitioning, which rely heavily on key-based data distribution.

The technical challenge stems from Spark's execution model, which divides work into stages and tasks. Each stage must complete all its tasks before the next stage can begin, creating a synchronization point in the processing pipeline. When examining Spark's execution DAG (Directed Acyclic Graph), the skewed partitions create what performance engineers call "long-tail latency" where a small number of tasks take significantly longer to complete than the majority [3]. This synchronization requirement means that even if 99 out of 100 tasks complete quickly, the entire stage must wait for that single slow task to finish.

Practical measurements from production Spark applications show that optimal partition sizes typically range between 100MB to 200MB of data [4]. However, in skewed scenarios, partition sizes can range from mere kilobytes to multiple gigabytes, creating extreme processing time differences. Since a Spark job can only complete when all tasks finish processing, these overloaded executors effectively become the limiting factor for the entire job, forcing the cluster to operate at the speed of its slowest component rather than leveraging its aggregate capacity.

*3. Common Causes of Data Skew*

Several scenarios typically lead to data skew in Spark applications, each creating distinct performance challenges in big data processing:

**Natural data imbalances**: Real-world datasets often have inherent skews that reflect natural phenomena. For example, in a dataset with 1 million sales records, 80% might belong to high-population cities like "New York" or "Los Angeles." This type of skew is particularly common in enterprise data architectures where analytical workloads encounter what data architects call the "long-tail distribution problem." According to lakeFS's analysis of enterprise data architecture patterns, modern data lakes typically contain a mix of structured and unstructured data with natural distribution imbalances that mirror business realities [5]. In retail analytics, geographic concentration of customers, seasonal buying patterns, and product popularity all create natural hotspots in the data. These patterns become particularly problematic when designing data lakes that need to handle both transactional data (OLTP) and analytical processing (OLAP) workloads simultaneously.

**Key distribution in joins**: When performing joins between tables, certain join keys may appear much more frequently than others, causing the partitions handling those keys to receive substantially more data. Research on data skewness methods has classified this phenomenon as "join key skew," which occurs when the frequency distribution of join key values exhibits high variance [6]. This pattern emerges frequently when joining fact tables with dimension tables in star schemas, especially when dimension tables contain hierarchical data like geographic regions or product categories. The join imbalance becomes particularly severe when combining historical data with current data, as certain identifiers may have accumulated significantly more records over time.

**Aggregation operations**: Functions like *groupByKey* and *reduceByKey* can create skew when some keys have significantly more associated values than others. This scenario, referred to as "aggregation key skew" in the literature, creates computation hotspots during the reduced phase of MapReduce-style operations [6]. The challenge is compounded in time-series data processing, where temporal aggregations (like daily, weekly, or monthly summaries) can create severe imbalances if data volumes vary significantly across time periods. In modern enterprise data architectures supporting both batch and streaming workloads, these aggregation patterns often become performance bottlenecks in dashboards and reporting pipelines [5].

**Transformations on skewed columns**: Operations on columns with highly imbalanced value distributions can perpetuate or worsen existing skew. When applying complex transformations to already skewed data, the computational workload becomes further imbalanced. This pattern is particularly evident in what researchers have termed "computational skew," where even with balanced data partitioning, the processing requirements vary significantly across partitions [6]. This occurs frequently in enterprise data lakes when performing resource-intensive operations like natural language processing on text fields, complex mathematical calculations, or machine learning feature engineering. The modern data mesh architecture approach attempts to address these challenges by distributing ownership of data products, but still requires careful design to prevent transformation bottlenecks [5].

| Skew Pattern | Description | Example Scenario | Impact |
|---|---|---|---|
| Natural Data Imbalances | Inherent distribution patterns in real-world data | E-commerce: Top products generate majority of sales | Certain partitions receive disproportionate data volume |
| Join Key Skew | Uneven distribution of values in join columns | Customer data: Popular regions have more records | Executors processing frequent keys become bottlenecks |
| Aggregation Key Skew | Some keys have significantly more values to aggregate | Social media: Viral content generates more interactions | Reduce operations become imbalanced across partitions |
| Computational Skew | Processing complexity varies across otherwise balanced data | Text analytics: Variable-length descriptions require different processing time | Tasks complete at different rates despite similar data volumes |

Table 1: Common Data Skew Patterns and Their Characteristics [5, 6]

*4. Solution 1: Repartitioning*

The simplest approach to addressing data skew is repartitioning, which redistributes data more evenly across the cluster. By explicitly specifying the number of partitions (100 in this example) and the column to partition by ("city"), we force Spark to redistribute the data across executors. This technique initiates a shuffle operation that creates a new set of partitions with a more balanced distribution of records.

According to recent research on effective partitioning strategies for large-scale analytics, repartitioning serves as a fundamental technique for addressing data imbalance issues in distributed processing frameworks [7]. When implemented with proper parameters, this approach can significantly reduce execution time for skewed workloads by ensuring more uniform resource utilization across the cluster. One key finding from partitioning research is that the optimal partitioning strategy should consider both data volume and processing complexity. For moderately skewed datasets, researchers observed performance improvements ranging from 20% to 40% when appropriate repartitioning strategies were applied, particularly for operations involving wide transformations like joins and aggregations.

The effectiveness of repartitioning depends largely on understanding the relationship between data characteristics and cluster configuration. Chandrakanth Lekkala emphasize that determining the ideal partition count and size requires careful consideration of available executor memory, core count, and the nature of subsequent transformations [7]. Their research indicates that partition sizes should generally remain within the 100-200MB range to balance parallelism with task scheduling overhead. Exceeding this range in either direction leads to performance degradation, with excessively small partitions creating scheduling bottlenecks and overly large partitions risking out-of-memory errors during shuffle operations.

However, as comparative analyses of big data frameworks have shown, simple repartitioning has inherent limitations when dealing with severely skewed data distributions [8]. The fundamental challenge is that Spark's default hash partitioner redistributes data based on the hash values of keys, which doesn't change the underlying frequency distribution of values within the dataset. Madiha Khalid and Muhammad Murtaza Yousaf observed that for datasets exhibiting extreme Zipfian distributions (where certain values appear disproportionately often), repartitioning alone provides only marginal improvements once certain threshold conditions are met [8]. Their framework comparison study demonstrated that Spark's performance with highly skewed data depends not just on partitioning strategy but also on memory management and shuffle optimization techniques.

Another consideration highlighted in framework adoption research is the trade-off between repartitioning benefits and the additional shuffle cost [8]. While repartitioning helps balance subsequent processing stages, it introduces its own shuffle operation that can be expensive, particularly for large datasets. The cost-benefit analysis provided by Narayanan and Kumar suggests that repartitioning is most effective for multi-stage pipelines where the initial shuffling cost is amortized across numerous downstream operations. For single-stage transformations, the overhead may outweigh the benefits unless skew is particularly severe.

This solution works well for moderate skew but may not be sufficient for extreme cases where certain values dominate the dataset. More sophisticated techniques are required when dealing with severe data imbalances.

| Parameter | Recommended Value | Considerations | Effect on Performance |
|---|---|---|---|
| Partition Count | 2-3× total core count | Higher values increase parallelism but add scheduling overhead | Balances task granularity with resource utilization |
| Optimal Partition Size | 100-200MB | Varies based on operation complexity and memory requirements | Prevents out-of-memory errors while maximizing throughput |
| Skew Threshold for Effectiveness | Skewness factor < 20:1 | Diminishing returns beyond this ratio | Most effective for moderate skew scenarios |
| Shuffle Cost | One-time data redistribution | Consider amortization across downstream operations | May not be justified for one-time transformations |

Table 2: Repartitioning [7, 8]

## 5. Solution 2: Key Salting

For severe skew where certain keys have disproportionately large data volumes, key salting offers a more targeted solution. This technique adds a random number (0-9) as a suffix to each value in the skewed column. The effect is that a single value like "New York" becomes multiple values ("New York_0", "New York_1", etc.), spreading records with that value across multiple partitions. After processing, you can remove the salt to restore the original values.

Key salting addresses fundamental partitioning challenges that simple repartitioning cannot solve. According to Anh Trần Tuấn's analysis of effective data partitioning strategies, key salting represents an advanced technique that thrives in situations where data exhibits natural hot spots that would otherwise create processing bottlenecks [9]. When implementing this approach, we

artificially increase the cardinality of frequently occurring values, creating what partitioning experts call "virtual keys" that distribute processing more evenly. This technique proves particularly valuable in scenarios where business requirements make it impossible to entirely redesign data structures, allowing teams to maintain logical data organization while improving physical distribution characteristics.

The implementation complexity of key salting must be balanced against its significant performance benefits. As Anh Trần Tuấn points out, the technique requires modifications to both the partitioning strategy and subsequent aggregation logic to ensure accurate results [9]. The pre-processing overhead is typically negligible compared to the performance gains, especially for iterative algorithms where the same data is processed multiple times. For ad-hoc queries or one-time transformations, the additional complexity might not justify the implementation effort, but for production ETL pipelines or regularly scheduled batch jobs, the consistent performance improvements make key salting a valuable addition to the optimization toolkit.

The effectiveness of key salting depends greatly on understanding the degree of skew in your dataset. According to Tyler garrett's research on operationalizing data skew detection, successful implementation requires first establishing monitoring to quantify skew patterns [10]. Their recommended approach involves calculating and tracking metrics like coefficient of variation (CV) across partition sizes, with CV values exceeding 1.0 indicating severe skew that would benefit from key salting techniques. This monitoring-first approach enables data engineers to apply key salting selectively where it provides the greatest benefit, rather than adding complexity to all data pipelines regardless of skew characteristics.

For join operations, which are particularly vulnerable to skew effects, key salting requires special consideration. Tyler garrett recommends a technique called "broadcast salt join" where the smaller dataset is replicated for each salt value and joined with the corresponding partition of the larger dataset [10]. This approach maintains join correctness while distributing the computational load more evenly across executors. Their performance analysis showed that for datasets with identified hot keys making up more than 40% of the total records, this technique consistently outperformed both standard joins and repartitioning approaches, especially as data volumes scaled into hundreds of gigabytes.

After processing with salted keys is complete, the original values must be reconstructed by stripping away the salt component, typically through a simple string manipulation operation that removes the suffix.

| Salt Factor | Appropriate Skew Scenario | Implementation Complexity | Memory Impact |
|---|---|---|---|
| Low (2-4) | Moderate skew with few dominant keys | Simple transformation, minimal downstream changes | Minor increase in metadata overhead |
| Medium (8-16) | Severe skew with identifiable hotspots | Requires consistent salting across operations | Balanced partition sizes reduce GC pressure |
| High (32+) | Extreme skew with single dominant value | Complex aggregation logic for results recombination | Potential for excessive small partitions |
| Adaptive | Dynamic skew patterns | Requires monitoring and adjustment framework | Optimizes resource utilization across varied workloads |

Table 3: Key Salting [9, 10]

*6. Solution 3: Broadcast Joins*

When joining a large table with a small one, broadcast joins can eliminate shuffling of the large dataset. This approach sends a complete copy of the smaller table to all executors, allowing each to perform the join locally without needing to shuffle the larger dataset. Broadcast joins are ideal when one table is small enough to fit in memory (typically under 10GB).

Broadcast joins represent one of the most effective techniques for addressing skew in join operations, particularly for star schema queries common in business intelligence applications. According to research published in Applied Sciences by Thanda Shwe and Masayoshi Aritsugi broadcast joins can dramatically reduce both computation time and resource utilization when dealing with asymmetric datasets [11]. Their experimental evaluation using real-world workloads demonstrated that broadcast joins significantly outperform standard shuffle-based joins when dimension tables are substantially smaller than fact tables. The researchers observed that as the ratio between table sizes increases, the benefits of broadcasting become more pronounced, with performance improvements scaling almost linearly with the size difference between tables.

The fundamental advantage of broadcast joins stems from their ability to eliminate data movement of the larger dataset. When using standard shuffle joins, both tables must be redistributed across the cluster based on join keys, creating network bottlenecks and potential memory pressure on executors handling frequently occurring keys. Thanda Shwe and Masayoshi Aritsugi's team found that in typical data warehouse workloads, join operations consume a disproportionate amount of execution time compared to other operations, with network transfer during shuffles representing a significant portion of this overhead [11]. By broadcasting the smaller table instead, this network congestion is avoided entirely for the larger dataset, substantially reducing both execution time and cluster resource consumption.

The decision to implement broadcast joins should be guided by a thorough understanding of the query execution process. As explained by CelerData's analysis of cost-based optimization techniques, modern query optimizers evaluate multiple join strategies and select the one with the lowest estimated cost based on statistics like table sizes, data distribution, and available system resources [12]. Their research indicates that while Spark can automatically determine when to use broadcast joins based on table size thresholds, manual optimization often yields better results for complex queries with multiple join operations. This is particularly true when dealing with skewed data distributions that may not be accurately captured by standard statistics.

Memory management becomes a critical consideration when implementing broadcast joins. CelerData's performance engineering team notes that broadcast tables must fit entirely within executor memory, with additional overhead required for deserialization and processing [12]. Their guidelines recommend broadcasting tables that occupy no more than 20-25% of available executor memory to avoid garbage collection issues and out-of-memory errors. For larger dimension tables that exceed this threshold but still exhibit significant size disparity compared to fact tables, techniques like bucketing or salting may provide better performance than standard shuffle joins.

Spark's broadcast join implementation includes several optimizations that enhance its effectiveness for skewed datasets. The broadcast variable uses a BitTorrent-like protocol to distribute data efficiently across the cluster, reducing network congestion at the driver node and enabling faster distribution even for larger broadcast tables.

*Conclusion*

Addressing data skew in Spark applications requires a strategic approach that matches the mitigation technique to the specific skew pattern encountered. Repartitioning offers a straightforward solution for moderate skew situations, distributing workloads more evenly with minimal implementation complexity. Key salting provides more targeted intervention for severe skew cases by artificially increasing key cardinality, effectively breaking up processing hotspots at the cost of additional transformation steps. Broadcast joins excel in asymmetric join scenarios by eliminating shuffle operations for larger datasets, dramatically improving performance when smaller tables can fit in executor memory. The effectiveness of these techniques depends greatly on understanding underlying data distributions and monitoring partition metrics to identify when and where skew occurs. By incorporating these solutions into their optimization toolkit and applying them selectively based on workload characteristics, data engineers can significantly improve job reliability, resource utilization, and processing times in their Spark applications, ultimately delivering more responsive and cost-effective data pipelines.

**References**

[1] Anh Trần Tuấn, "Strategies for Effective Data Partitioning: Why It Matters and How to Implement It," Medium, 2024. [Online]. Available: https://medium.tuanh.net/strategies-for-effective-data-partitioning-why-it-matters-and-how-to-implement-it-1f9d6da4df68

[2] CelerData, "What is a Cost Based Optimizer?," CelerData, 2024. [Online]. Available: https://celerdata.com/glossary/cost-based-optimizer#:~:text=A%20Cost%2DBased%20Optimizer%20(CBO,with%20the%20lowest%20estimated%20cost.

[3] Chandrakanth Lekkala, "Strategies for Effective Partitioning Data at Scale in Large-scale Analytics," Journal of Big Data Analytics, 2019. [Online]. Available: https://www.researchgate.net/publication/382365695_Strategies_for_Effective_Partitioning_Data_at_Scale_in_Large-scale_Analytics

[4] Iddo Avneri, "Analytical Data: Guide to Enterprise Data Architecture," lakeFS, 2024. [Online]. Available: https://lakefs.io/blog/analytical-data-guide-enterprise-data-architecture/

[5] Jerome H. Saltzer and M. Frans Kaashoek, "Performance Bottleneck," 2009, Principles of Computer System Design, 2009. [Online]. Available: https://www.sciencedirect.com/topics/computer-science/performance-bottleneck

[6] Madiha Khalid and Muhammad Murtaza Yousaf, "A Comparative Analysis of Big Data Frameworks: An Adoption Perspective," Applied Sciences, 2021. [Online]. Available: https://www.researchgate.net/publication/356461768_A_Comparative_Analysis_of_Big_Data_Frameworks_An_Adoption_Perspective

[7] Pramit Marattha, "Apache Spark Performance Tuning: 7 Optimization Tips (2025)," ChaosGenius, 2023. [Online]. Available: http://chaosgenius.io/blog/spark-performance-tuning/#:~:text=To%20mitigate%20this%2C%20one%20of,to%20perform%20the%20join%20locally.

[8] Scaibu, "Spark Partitioning: Unlocking Big Data Performance," Medium, Feb. 2024. [Online]. Available: https://scaibu.medium.com/spark-partitioning-unlocking-big-data-performance-e08f6891135d

[9] Subhankar Mishra et al., "Various Data Skewness Methods in the Hadoop Environment," 2019 International Conference on Recent Advances in Energy-efficient Computing and Communication (ICRAECC), 2019. [Online]. Available: https://www.researchgate.net/publication/339264546_Various_Data_Skewness_Methods_in_the_Hadoop_Environment

[10] Thanda Shwe and Masayoshi Aritsugi, "Optimizing Data Processing: A Comparative Study of Big Data Platforms in Edge, Fog, and Cloud Layers," Applied Sciences, 2024. [Online]. Available: https://www.mdpi.com/2076-3417/14/1/452

[11] Tyler garrett, "Operationalizing Data Skew Detection in Distributed Processing," Dev3lop Engineering, 2025. [Online]. Available: https://dev3lop.com/operationalizing-data-skew-detection-in-distributed-processing/

[12] Vishnu Vardhan Reddy Chilukoori and Srikanth Gangarapu, "Strategic Cost Management in Cloud-Based Big Data Processing: An AWS Case Study from Amazon," Ijraset Journal For Research in Applied Science and Engineering Technology, 2024. [Online]. Available: https://www.ijraset.com/research-paper/strategic-cost-management-in-cloud-based-big-data-processing-an-aws-case-study-from-amazon