

Flink TPC-H Q3 Performance Benchmarking Project based on Cquirrel

Name: CAI, Huaiyu Email: hcaiai@connect.ust.hk StuID: 21139725

GitHub Link: <https://github.com/PeaceChoy/Flink-TPC-H-Q3-Performance-Benchmarking-Project>

1. Introduction

This report presents the results of performance benchmarking conducted on Apache Flink with Cquirrel's BasedProcessFunction, a powerful stream processing framework that also supports batch/stream processing. The benchmarking efforts focused on evaluating Flink's performance characteristics using the TPC-H benchmark through the minicluster mode, a standard decision support benchmark that simulates data warehousing and analytics environments. The primary goal was to assess Flink's processing capabilities under different configurations, optimization strategies and scales of datasets, providing insights into its performance scalability and optimization potential.

2. Benchmark Environment

The benchmarking was conducted in a controlled environment to ensure consistent and reliable results. I utilized the TPC-H data generator to create datasets of varying sizes (scale = 0.1, 1, 10, 20), then processed these datasets using Apache Flink in minicluster environment with different configuration parameters. This project used Apache Flink version 1.18.1 running on Java 17.0.12 with a Windows 10 platform. All tests were executed on the same hardware to maintain consistency across different test scenarios.

3. Methodology

The benchmark methodology focused on implementing TPC-H Query 3 using Cquirrel's BasedProcessFunction, in the following contents, I will directly call the Cquirrel's BasedProcessFunction as ProcessFunction. This query performs a three-way join between the lineitem, orders, and customer tables, followed by grouping, aggregation, and sorting operations - representing a complex analytical workload that tests multiple aspects of a data processing system.

This project developed three distinct implementations, the first implementation (flink_origin.java) represents a baseline single-threaded execution model using ProcessFunction. This implementation establishes a performance baseline without minimal parallelism.

The second implementation (flink_4threads.java) extends the baseline by increasing parallelism to 4 threads, allowing the program to assess Flink's scaling characteristics with increased computational resources. The third implementation (flink_optimize.java) builds

upon the multi-threaded approach by incorporating various optimization techniques, including memory management adjustments, broadcast hints for joins, mini-batch processing, and tuned ProcessFunction implementations.

For each implementation, I measured several performance metrics:

- Total Job Execution Time (ms)
- Table Creation Time
- Query Planning Time
- Data Processing Time
- Sink Processing Time
- Framework Overhead
- Overall Throughput (records per second)
- Processing Throughput
- Sink Throughput
- Total Records Processed
- Total Records Output
- Late Records Number

To ensure robust results, I executed each test configuration multiple times across different dataset sizes generated using scale factors of 0.1, 1, 10, and 20.

4. Implementation Details

All three implementations share a common structure for executing TPC-H Query 3:

```
SELECT
  l_orderkey,
  SUM(l_extendedprice * (1 - l_discount)) AS revenue,
  o_orderdate,
  o_shippriority
FROM lineitem
JOIN orders ON l_orderkey = o_orderkey
JOIN customer ON o_custkey = c_custkey
WHERE
  c_mktsegment = 'BUILDING' AND
  o_orderdate < DATE '1995-03-15' AND
  l_shipdate > DATE '1995-03-15'
GROUP BY l_orderkey, o_orderdate, o_shippriority
ORDER BY revenue DESC, o_orderdate
```

The implementations differ primarily in parallelism settings and optimization techniques. The optimized version includes several key improvements. The memory configuration was enhanced by adjusting network memory fraction and enabling off-heap memory to reduce garbage collection overhead. I incorporated broadcast hints for the customer table, which is typically smaller than the other tables, to optimize join operations. The implementation also

enables mini-batch processing to reduce the overhead of processing individual records, particularly beneficial for high-throughput scenarios. Furthermore, I refined the ProcessFunction implementations to increase the rate limit for better throughput. Buffer timeout settings were tuned to balance latency and throughput considerations.

5. Results and Analysis

The performance results from the benchmark tests reveal significant insights into Flink's behavior under different configurations.

5.1 Single-Thread Implementation

Metric	Scale = 0.1	Scale = 1	Scale = 10	Scale = 20
Execution Time (ms)				
Total Job Execution	18,089	24,853	85,357	158,851
Table Creation	196 (1.1%)	180 (0.7%)	183 (0.2%)	367 (0.2%)
Query Planning	14,714 (81.3%)	18,593 (74.8%)	59,039 (69.2%)	107,496 (67.7%)
Data Processing	455 (2.5%)	3,273 (13.2%)	23,967 (28.1%)	48,802 (30.7%)
Sink	115 (25.3% of processing)	1,093 (33.4% of processing)	2,367 (9.9% of processing)	3,433 (7.0% of processing)
Framework Overhead	2,724 (15.1%)	2,807 (11.3%)	2,168 (2.5%)	2,186 (1.4%)
Throughput (records/sec)				
Overall Throughput	7.63	467.55	1,335.60	1,432.86
Processing Throughput	303.30	3,550.26	4,756.67	4,663.97
Sink Throughput	1,200.00	10,631.29	48,163.50	66,300.90
Data Statistics				
Total Records Processed	138	11,620	114,003	227,611
Total Records Output	138	11,620	114,003	227,611
Late Records	0	2	22	45

The single-thread implementation established the baseline performance. Despite running on a single thread, Flink demonstrated robust processing capabilities. However, as expected, this implementation showed limitations in throughput when processing larger datasets.

1. Query Planning Dominance

Regardless of input scale, query planning consumes 67.7–81.3% of total runtime.

Because the engine performs this phase on one thread, the system becomes CPU-bound early and stays that way; every doubling of data volume simply lengthens the planning step linearly. The fixed planning cost therefore caps overall throughput growth and prevents the job from scaling beyond 1.5k records/sec even at the largest scale.

2. Serial Sink Bottleneck

The sink phase, though small in relative terms (7–33 % of processing time), executes strictly in sequence after processing. Its throughput tops out at 66.3k records/sec, which is two orders of magnitude below the potential disk or network bandwidth. This serialization point creates pressure that keeps the processing threads idle while data is flushed, further limiting end-to-end speed.

5.2 Multi-Thread Implementation

Metric	Scale = 0.1	Scale = 1	Scale = 10	Scale = 20
Execution Time (ms)				
Total Job Execution	14,629	24,682	72,359	134,060
Table Creation	189 (1.3 %)	174 (0.7 %)	212 (0.3 %)	194 (0.1 %)
Query Planning	12,151 (83.1 %)	19,041 (77.1 %)	45,380 (62.7 %)	81,845 (61.1 %)
Data Processing	326 (2.2 %)	3,249 (13.2 %)	24,644 (34.1 %)	49,562 (37.0 %)
Sink	6 (1.8 % of processing)	27 (0.8 % of processing)	93 (0.4 % of processing)	104 (0.2 % of processing)
Framework Overhead	1,963 (13.4 %)	2,218 (9.0 %)	2,123 (2.9 %)	2,459 (1.8 %)
Throughput (records/sec)				
Overall Throughput	9.43	470.79	1,575.52	1,697.83
Processing Throughput	423.31	3,576.49	4,625.99	4,592.45
Sink Throughput	23,000.00	430,370.37	1,225,838.71	2,188,567.31
Data Statistics				

Metric	Scale = 0.1	Scale = 1	Scale = 10	Scale = 20
Total Records Processed	138	11,620	114,003	227,611
Total Records Output	138	11,620	114,003	227,611
Late Records	0	2	22	45

The transition to a 4-thread implementation resulted in substantial performance improvements.

1. Parallel Sink Eliminates I/O Stall

A dedicated thread pool performs output writes concurrently. Sink latency collapses to < 1 % of processing time and throughput jumps to over 2 M records/sec at Scale 20. By removing the serial sink bottleneck, the pipeline keeps all processing threads saturated; this single change delivers the entire 15–20 % wall-clock speed-up observed at large scales.

2. Planning Phase Still Single-Threaded

Query planning remains on one thread and still represents 61.1–83.1 % of total runtime. Because this phase cannot be parallelized, the overall speed-up is bounded: once the sink bottleneck is gone, further gains would require algorithmic improvements in planning, not additional threads.

5.3 Optimized Implementation

Metric	Scale = 0.1	Scale = 1	Scale = 10	Scale = 20
Execution Time (ms)				
Total Job Execution	10,487	14,253	41,999	72,459
Table Creation	116 (1.1 %)	133 (0.9 %)	132 (0.3 %)	118 (0.2 %)
Query Planning	7,516 (71.7 %)	10,567 (74.1 %)	27,759 (66.1 %)	46,564 (64.3 %)
Data Processing	1,334 (12.7 %)	2,288 (16.1 %)	12,897 (30.7 %)	24,573 (33.9 %)
Sink	7 (0.5 % of processing)	11 (0.5 % of processing)	47 (0.4 % of processing)	78 (0.3 % of processing)
Framework Overhead	1,521 (14.5 %)	1,265 (8.9 %)	1,211 (2.9 %)	1,204 (1.7 %)
Throughput (records/sec)				

Metric	Scale = 0.1	Scale = 1	Scale = 10	Scale = 20
Overall Throughput	13.16	815.27	2,714.42	3,141.24
Processing Throughput	103.45	5,078.67	8,839.50	9,262.65
Sink Throughput	19,714.29	1,056,363.64	2,425,595.74	2,918,089.74
Data Statistics				
Total Records Processed	138	11,620	114,003	227,611
Total Records Output	138	11,620	114,003	227,611
Late Records	0	1	11	22

The optimized implementation delivered the most significant performance gains. It attacks those single-thread costs with broadcast joins, mini-batch execution, and better memory management, doubling overall performance on the same 4-thread setup.

1. Planning phase shaved by SQL-level tuning

By adding a broadcast hint for the customer dimension and removing the LIMIT clause, while the scale of TPC-H dataset being larger, the optimizer is steered toward a cheaper join strategy and a smaller search space. This trims the Query Planning time from 81,845ms to 46,564ms, 43% reduction that translates directly into the observed 46% shorter total job runtime.

2. Runtime operators keep 4-way parallelism, but memory & object-reuse cut CPU stalls

With the plan fixed, the heavy operators still execute in parallel across four sub-tasks, but the explicit enableObjectReuse() call and a larger managed-memory fraction eliminate the per-record allocation storms and minor garbage-collection pauses that lingered in the multi-thread run. These JVM-level tweaks cut Data Processing time from 49,562ms to 24,573ms, a further 50% win, while the single-threaded sink is left untouched yet contributes only 78ms, confirming that the pipeline is now CPU-bound rather than I/O-bound.

5.4 Scaling Characteristics

The tests across different dataset sizes revealed important insights into Flink's scaling characteristics. As the dataset size increased, we observed that the optimized implementation maintained consistently better performance compared to the other implementations. Also, the performance gap widened with larger datasets, highlighting the increasing importance of optimization techniques as data volume grows.

6. Performance Bottlenecks

The benchmarking identified several performance bottlenecks in Flink processing:

1. CPU Saturation Ceiling

With parallelism fixed at four slots, the heavy hash-join and aggregation operators already fully occupy the available cores; doubling the slot count would yield diminishing returns. The job is therefore CPU-bound during Data Processing, and any future gain must come from algorithmic improvements (e.g., better join order, Bloom filters) rather than more threads.

2. Memory Pressure from Object Reuse

Although object reuse reduces GC churn, the query still materializes large intermediate hash tables. With managed memory capped at 2GB, occasional spilling to disk becomes noticeable at Scale = 20, adding tens of milliseconds per partition and widening the tail latency of the processing phase.

3. Single-threaded Sink

The final sink is deliberately pinned to one sub-task to preserve result order. While its absolute cost is only 78ms, it creates a back-pressure window where the upstream pipeline must buffer the last batch of records, lengthening the tail of the Data Processing phase by a few hundred milliseconds. However, with adding sink threads to 4, the sink time multiplied, It shows a regressive phenomenon of "parallel expenses exceeding parallel income", which may only occur in the minicluster mode.

7. Recommendations

Based on the benchmark findings, we recommend the following optimization strategies for Flink deployments:

1. Parallelism should be adjusted based on available resources and dataset characteristics. The results indicate significant benefits from increased parallelism, but the optimal setting will depend on specific workload characteristics and hardware configuration.
2. Memory configuration, particularly network buffer settings and heap/off-heap allocation, should be carefully tuned. The default settings may not be optimal for complex analytical queries with large intermediate results.
3. Join optimization should be prioritized for queries involving multiple large tables. Broadcast hints for smaller tables can provide substantial performance benefits, and careful consideration of join order can significantly impact performance.
4. Mini-batch processing should be enabled for high-throughput scenarios, as it reduces per-record processing overhead while maintaining acceptable latency characteristics for analytical workloads.
5. Buffer timeout settings should be tuned based on throughput and latency requirements, finding the appropriate balance for specific use cases.

8. Conclusion

The performance benchmarking of Apache Flink on TPC-H Query 3 reveals that while Flink is capable of efficiently processing complex analytical workloads, its performance is highly sensitive to configuration and optimization strategies. The evaluation across three implementations—baseline single-threaded, multi-threaded, and optimized—demonstrates a clear progression in performance gains driven by parallelism, memory management, and query optimization techniques.

Key findings indicate that the query planning phase is a major bottleneck, consuming up to 81% of total execution time in the baseline setup and remaining single-threaded across all configurations. This limits scalability regardless of increased parallelism in data processing. However, optimization techniques such as broadcast joins for small dimension tables and removal of restrictive clauses significantly reduce planning overhead, cutting planning time by over 40% in the optimized version.

The introduction of 4-thread parallelism eliminates the serial sink bottleneck, increasing sink throughput by more than 30x and ensuring that processing resources remain fully utilized. Further improvements from mini-batch processing, object reuse, and managed memory tuning reduce CPU stalls and garbage collection pressure, effectively halving data processing time at larger scales.

Despite these gains, the system eventually becomes CPU-bound, with diminishing returns from additional parallelism. Memory pressure and single-threaded stages—particularly the ordered sink—introduce tail latency and limit end-to-end scalability. These bottlenecks highlight the importance of holistic optimization that balances parallelism, memory usage, and algorithmic efficiency.

In conclusion, Flink's performance on complex batch analytics can be dramatically improved—by up to 54% in total execution time through targeted optimizations—even within a minicluster environment. For production deployments, these results underscore the necessity of careful configuration tuning, including parallelism settings, memory allocation, join strategies, and batch processing parameters. Future work could explore adaptive query planning, operator fusion, or out-of-core execution to further mitigate current limitations and enhance Flink's competitiveness in large-scale analytical processing.