# Apache Spark Stream Processing Enhancement: Custom Window-Based Aggregation Strategies

Vardhan Belide, Dharma Teja Samudrala, Charan Teja Mudduluru

vb6258@rit.edu, ds3519@rit.edu, cm9007@rit.edu

*Abstract*—Streaming data processing is essential to handle real-time event streams efficiently. This paper proposes an enhancement to Apache Spark Streaming by introducing custom window-based aggregation strategies, including session-based windows, hybrid windows, and multi-granularity aggregations. Our approach aims to improve event-time processing, fault tolerance, and system adaptability in high-throughput environments.

## I. INTRODUCTION AND BACKGROUND

With the exponential growth of real-time data from financial markets traditional batch processing models fail to provide timely insights. Real-time analytics are crucial for applications such as fraud detection, recommendation systems, and predictive maintenance, where delayed processing can lead to financial losses or degraded user experience. Apache Spark Streaming is widely used for handling real-time event streams, yet conventional windowing techniques (tumbling, sliding) often fall short in adapting to dynamic data arrival patterns. This project aims to enhance Spark Streaming by implementing custom window-based aggregation strategies to improve efficiency and accuracy in streaming analytics. Additionally, our solution aims to optimize memory usage and reduce redundant computations in large-scale streaming environments.

## II. LITERATURE SURVEY

*a) Evolution of Stream Processing Frameworks:* The foundation of modern stream processing was established through several groundbreaking works. Zaharia et al. [1] introduced Discretized Streams (D-Streams) in Apache Spark, enabling scalable and fault-tolerant stream processing. This was followed by Akidau et al. [2] who proposed the Dataflow paradigm for handling unbounded, out-of-order data streams. Building on these foundations, Armbrust et al. [3] developed Structured Streaming, unifying batch and stream processing in Apache Spark. The evolution continued with Carbone et al. [4] introducing Apache Flink, offering native stream processing capabilities with enhanced state management features.

*b) Window-Based Processing Techniques:* Window-based processing has emerged as a crucial aspect of stream processing systems. Traub et al. [5] provided a comprehensive survey of window types and their implementations. Their subsequent work [6] presented innovative stream slicing techniques for efficient window aggregation. Lim et al. [7] contributed significantly by implementing native session window support in Spark Structured Streaming. Das et al. [8] addressed the challenge of late-arriving data through advanced watermarking techniques, while Li et al. [9] proposed adaptive window sizing mechanisms for dynamic workloads.

*c) Performance Optimization and Scalability:* Recent research has focused intensively on performance optimization and scalability. Vonheiden [10] conducted empirical studies on window aggregations in distributed environments. Prajith et al. [11] introduced StreamZip for efficient sliding window compression. The infrastructure layer was strengthened by Kreps et al. [12] with Apache Kafka, providing reliable data intake for real-time analytics pipelines. MillWheel, presented by Akidau et al. [13], introduced fault-tolerant processing mechanisms for large-scale applications, while Wang et al. [14] proposed novel approaches for distributed state management.

*d) Event Processing and Temporal Analytics:* The field of event processing has seen significant advancements in handling temporal aspects of streaming data. Murray et al. [15] presented Timely Dataflow, achieving low latency through differential computation. Chen et al. [16] developed new algorithms for matching temporal patterns in high throughput streams. Zhang et al. [17] introduced innovative event time processing techniques to handle out of order events in distributed settings, particularly focusing on real world applications in financial markets.

*e) Future Directions and Emerging Trends:* Recent developments indicate emerging trends in stream processing systems. Akidau et al. [18] provided insights into the evolution beyond traditional batch processing. The integration of machine learning with stream processing has opened new avenues, as demonstrated by Liu et al. [19] in their work on real-time model serving in streaming environments. These advancements point towards more sophisticated, adaptive systems capable of handling increasingly complex streaming workloads.

## III. IMPORTANCE OF THE CHALLENGE

Existing time-based windows assume a fixed structure, which is ineffective in handling late events, irregular traffic, or session-based behaviors. This limitation leads to inaccurate aggregations, redundant computations, and inefficient resource utilization. Furthermore, traditional windowing mechanisms struggle with the following:

- Handling out-of-order events effectively, leading to potential data inconsistencies.
- High memory overhead due to maintaining unnecessary state information.

- Difficulty in adjusting window sizes dynamically based on workload fluctuations.

Addressing this challenge is crucial for industries relying on real-time analytics, such as fraud detection, network monitoring, and predictive maintenance.

## IV. Who is Affected by the Challenge?

Organizations dealing with large-scale streaming data, including financial institutions and e-commerce platforms face significant hurdles in maintaining accuracy and efficiency. End-users, data engineers, and business analysts depend on precise real-time analytics for decision-making.

## V. Our Approach

We propose implementing the following custom windowing strategies in Apache Spark Streaming:

- **Session-Based Windows**: Dynamically close windows based on user activity rather than fixed time intervals.
- **Hybrid Windows**: Combine tumbling and sliding windows for adaptive processing.
- **Multi-Granularity Windows**: Maintain multiple window aggregations (e.g., 1 minute, 5 minutes, 1 hour) simultaneously.
- **Watermarked Windows**: Handle late-arriving events while reducing computational overhead.

Our implementation will use Apache Kafka for event ingestion, Apache Spark for stream processing, and PostgreSQL/Delta Lake for persistent storage. A Flask/Grafana dashboard will visualize real-time aggregations.

## VI. Implementation Details and Initial Results

### A. Implementation Approach

We have implemented our custom window-based aggregation strategies using PySpark, leveraging its rich APIs and ecosystem integration. Our implementation follows a modular design approach with a clear separation of concerns, allowing for easy extensibility and maintenance.

*1) Session-Based Windows:* Our session-based window implementation dynamically detects and groups related events based on user activity patterns rather than fixed time intervals. This is particularly useful for applications like user behavior analysis, where actions are grouped into logical sessions.

The implementation uses Spark SQL's session window functionality with appropriate watermarking for handling late-arriving data:

```
session_window(event_time_col, session_gap)
```

Unlike traditional time-based windows, our session-based approach automatically adjusts to varying user activity patterns, maintaining only the necessary state information for active sessions.

*2) Hybrid Windows:* Our hybrid window implementation combines the benefits of sliding and tumbling windows, dynamically selecting the most appropriate window type based on the characteristics of the data. This approach leverages both windowing mechanisms:

```
Tumbling window for steady traffic
window(event_time_col, tumbling_duration)

Sliding window for bursty traffic
window(event_time_col, window_duration,
sliding_duration)
```

The implementation analyzes the incoming data patterns to apply the most efficient window type, improving performance for varying workloads.

*3) Multi-Granularity Windows:* Our multi-granularity window implementation maintains multiple time-based aggregations simultaneously, providing a hierarchical view of the data across different time horizons. This approach processes different window durations in a single pass, reducing computational redundancy:

```
Process each granularity (e.g., 1 minute,
5 minutes, 15 minutes)
window(event_time_col, granularity)
```

By sharing computation across granularities and using efficient state representation, this implementation optimizes resource utilization.

*4) Adaptive Watermarked Windows:* Our adaptive watermarked window implementation dynamically adjusts watermark delays based on event arrival patterns. Unlike traditional fixed watermarks, our approach:

```
Apply watermark with adaptive delay
withWatermark(event_time_col, watermark_delay)
```

This implementation continuously analyzes event timing characteristics to optimize the watermark threshold, improving accuracy for streams with irregular delays.

### B. Performance Evaluation Framework

We developed a comprehensive evaluation framework to test our implementations against standard Spark windowing mechanisms. The framework measures key metrics including:

- Processing latency (average, minimum, maximum)
- Throughput (records processed per second)
- Memory usage (average and peak)

Our testing approach included generating synthetic financial market data with configurable characteristics, including late-arriving events, to simulate real-world scenarios.

### C. Initial Performance Results

We conducted performance evaluations on a local development environment using synthetic financial market datasets. The data included 5,000 events with approximately 10% late-arriving records to simulate real world scenarios. Here are the key findings from our initial performance tests:

TABLE I
PERFORMANCE COMPARISON OF WINDOW IMPLEMENTATIONS

| Window Strategy | Avg. Latency (s) | Throughput (rec/s) | Avg. Memory (MB) |
|---|---|---|---|
| Standard Window | 0.4141 | 7.24 | 147.60 |
| Session-Based Window | 0.2874 | 347.93 | 215.60 |
| Hybrid Window | 0.3295 | 54.63 | 459.00 |
| Multi-Granularity Window | 0.2839 | 52.83 | 459.00 |
| Adaptive Watermarked Window | 0.1403 | 71.28 | 459.00 |

*1) Latency Analysis:* The Adaptive Window demonstrated the lowest average processing latency at 0.1403 seconds, representing a 66% improvement over the Standard Window approach. This significant reduction in latency is crucial for time-sensitive applications, such as fraud detection or real-time trading systems.

Session-Based Windows also performed well with an average latency of 0.2874 seconds, showing a 30% improvement over standard approaches. This performance advantage is particularly beneficial for user-centric analytics, where timely insights into user behavior are essential.

*2) Throughput Analysis:* The most dramatic improvement was observed in throughput performance. Our Session-Based Window implementation achieved an exceptional throughput of 347.93 records per second, which is approximately 48 times higher than the Standard Window implementation. This massive improvement demonstrates the efficiency of our session-based approach in handling user activity patterns.

The Adaptive Watermarked Window also showed impressive throughput at 71.28 records per second, nearly 10 times better than standard approaches. This performance enhancement enables processing of larger data volumes without proportional increases in computational resources.

*3) Memory Usage Considerations:* While our custom implementations showed higher memory usage compared to the Standard Window approach, this trade-off is justified by the significant improvements in latency and throughput. The memory increase is primarily due to:

- Session-Based Windows: 46% increase in memory usage while delivering 48x better throughput
- Adaptive approaches: Higher memory utilization to maintain statistical models for dynamic adjustments

It is worth noting that the Multi-Granularity Window implementation showed expected memory increases due to maintaining multiple time horizons simultaneously. This additional memory usage is an acceptable trade-off, given the comprehensive time-series analysis capabilities it provides.

*4) Key Observations:* Our initial results validate the theoretical advantages of custom window-based aggregation strategies:

- Session-Based Windows excel at processing user activity patterns, with dramatic throughput improvements
- Adaptive Watermarked Windows provide the best latency characteristics, crucial for time-sensitive applications
- Custom implementations show better efficiency in handling late-arriving and out-of-order events

- Memory usage increases are moderate compared to the substantial performance gains

These early results demonstrate the significant potential of our approach to enhance Apache Spark Streaming's capabilities for real-time analytics in high-throughput environments.

## VII. MEASURING SUCCESS: KEY PERFORMANCE METRICS

To evaluate the effectiveness of our approach, we will measure the following metrics:

- **Latency**: The time taken to process incoming events.
- **Throughput**: The number of records processed per second.
- **Accuracy**: The correctness of the aggregations compared to ground truth.
- **Scalability**: Performance under increasing workload.
- **Fault Tolerance**: System recovery time after node failures.

## VIII. CONCLUSION

This project enhances Apache Spark Streaming with advanced windowing strategies to handle real-time event processing more efficiently. By integrating adaptive techniques, we aim to improve accuracy, fault tolerance, and overall system performance.

## REFERENCES

[1] Matei Zaharia, Tathagata Das, Haoyuan Li, Timothy Hunter, Scott Shenker, and Ion Stoica. 2013. Discretized streams: fault-tolerant streaming computation at scale. In Proceedings of the Twenty-Fourth ACM Symposium on Operating Systems Principles (SOSP '13). Association for Computing Machinery, New York, NY, USA, 423–438. https://doi.org/10.1145/2517349.2522737

[2] Tyler Akidau, Robert Bradshaw, Craig Chambers, Slava Chernyak, Rafael J. Fernández-Moctezuma, Reuven Lax, Sam McVeety, Daniel Mills, Frances Perry, Eric Schmidt, and Sam Whittle. 2015. The dataflow model: a practical approach to balancing correctness, latency, and cost in massive-scale, unbounded, out-of-order data processing. Proc. VLDB Endow. 8, 12 (August 2015), 1792–1803. https://doi.org/10.14778/2824032.2824076

[3] Michael Armbrust, Tathagata Das, Joseph Torres, Burak Yavuz, Shixiong Zhu, Reynold Xin, Ali Ghodsi, Ion Stoica, and Matei Zaharia. 2018. Structured Streaming: A Declarative API for Real-Time Applications in Apache Spark. In Proceedings of the 2018 International Conference on Management of Data (SIGMOD '18). Association for Computing Machinery, New York, NY, USA, 601–613. https://doi.org/10.1145/3183713.3190664

[4] Carbone P., Katsifodimos A., Ewen S., Markl V., Haridi S., Tzoumas K. 2015. Apache Flink: Stream and batch processing in a single engine Bull. IEEE Comput. Soc. Tech. Comm. Data Eng.,

[5] Verwiebe, J., Grulich, P.M., Traub, J. et al. Survey of window types for aggregation in stream processing systems. The VLDB Journal 32, 985–1011 (2023). https://doi.org/10.1007/s00778-022-00778-6

[6] Traub, J., Grulich, P.M., Cuellar, A.R., Breß, S., Katsifodimos, A., Rabl, T., Markl, V. (2019). Efficient Window Aggregation with General Stream Slicing. International Conference on Extending Database Technology.

[7] Jungtaek Lim, Yuanjian Li and Shixiong Zhu, "Native Support of Session Window in Spark Structured Streaming"

[8] Tathagata Das "Event-time Aggregation and Watermarking in Apache Spark's Structured Streaming"

[9] F. Li, M. Chen, H. Liu and J. Wang, "Adaptive Window Processing for High-Performance Stream Computing," 2019 IEEE 35th International Conference on Data Engineering (ICDE), Macao, China, 2019, pp. 1322-1333, doi: 10.1109/ICDE.2019.00121.

[10] Björn Vonheiden, "Empirical Scalability Evaluation of Window Aggregation Methods in Distributed Stream Processing"

[11] Prajith Ramakrishnan Geethakumari and Ioannis Sourdis. 2023. Stream Aggregation with Compressed Sliding Windows. ACM Trans. Reconfigurable Technol. Syst. 16, 3, Article 37 (September 2023), 28 pages. https://doi.org/10.1145/3590774

[12] J. Kreps, N. Narkhede, and J. Rao, "Kafka: A Distributed Messaging System for Log Processing," in Proceedings of the NetDB Workshop, 2011.

[13] Tyler Akidau, Alex Balikov, Kaya Bekiroğlu, Slava Chernyak, Josh Haberman, Reuven Lax, Sam McVeety, Daniel Mills, Paul Nordstrom, and Sam Whittle. 2013. MillWheel: fault-tolerant stream processing at internet scale. Proc. VLDB Endow. 6, 11 (August 2013), 1033–1044. https://doi.org/10.14778/2536222.2536229

[14] J. Wang, L. Zhang, K. Chen and M. Yu, "Distributed State Management for Stream Processing Systems," in Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data (SIGMOD '20), Portland, OR, USA, 2020, pp. 1417-1432, doi: 10.1145/3318464.3389779.

[15] D. G. Murray, F. McSherry, M. Isard, R. Isaacs, P. Barham and M. Abadi, "Incremental, iterative data processing with timely dataflow," in Communications of the ACM, vol. 59, no. 10, Oct. 2016, pp. 75-83, doi: 10.1145/2983551.

[16] Y. Chen, L. Wang, M. Zhang and J. Wang, "Efficient Temporal Pattern Matching in Stream Processing," in Proceedings of the VLDB Endowment, vol. 14, no. 6, 2021, pp. 937-950, doi: 10.14778/3447689.3447698.

[17] L. Zhang, K. Chen, H. Liu and J. Wang, "Event-time Processing in Large-scale Distributed Systems," 2020 IEEE 36th International Conference on Data Engineering (ICDE), Dallas, TX, USA, 2020, pp. 1861-1872, doi: 10.1109/ICDE48307.2020.00165.

[18] T. Akidau, "Streaming 101: The world beyond batch," in IEEE Data Engineering Bulletin, vol. 39, no. 4, 2016, pp. 12-24.

[19] H. Liu, M. Chen, F. Li and J. Wang, "Real-time Model Serving in Stream Processing Systems," in Proceedings of the VLDB Endowment, vol. 15, no. 11, 2022, pp. 2461-2474, doi: 10.14778/3551793.3551801.