

WSPPD Presentation

A Pre-Aggregation Strategy Towards Reduction on Job Completion Time in Stream Process Systems

Breno Fanchiotti Zanchetta - 00240494

Paulo R. Souza Jr, Kassiano J. Matteussi, Vinícius P. Perego

Julio C. S. Anjos, Claudio F. R. Geyer, Edison P. de Freitas

Agenda

3. Context

4. Aggregation

5. Window

6. Research Question

7. Model

8. Experiments

9. Results

10. Conclusion and Future Works

11. References

12. Final Slides

Aggregation

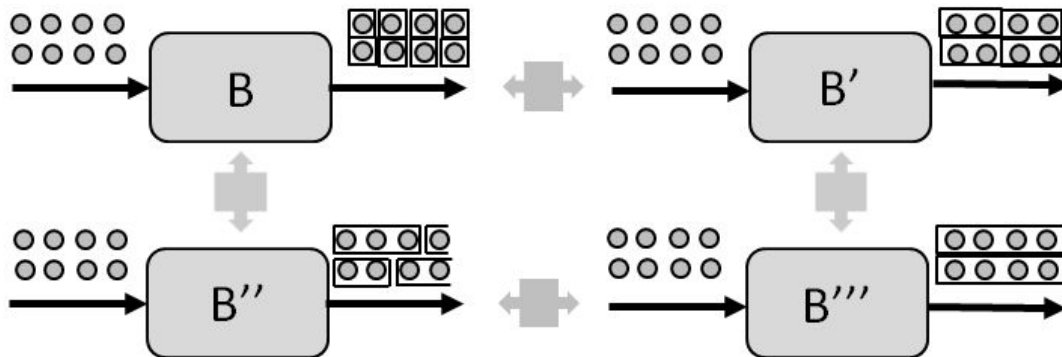
Static value is configured before experiments (Python Node).

For every input stream:

For every N aggregations:

Aggregate the Stream

Send Aggregation



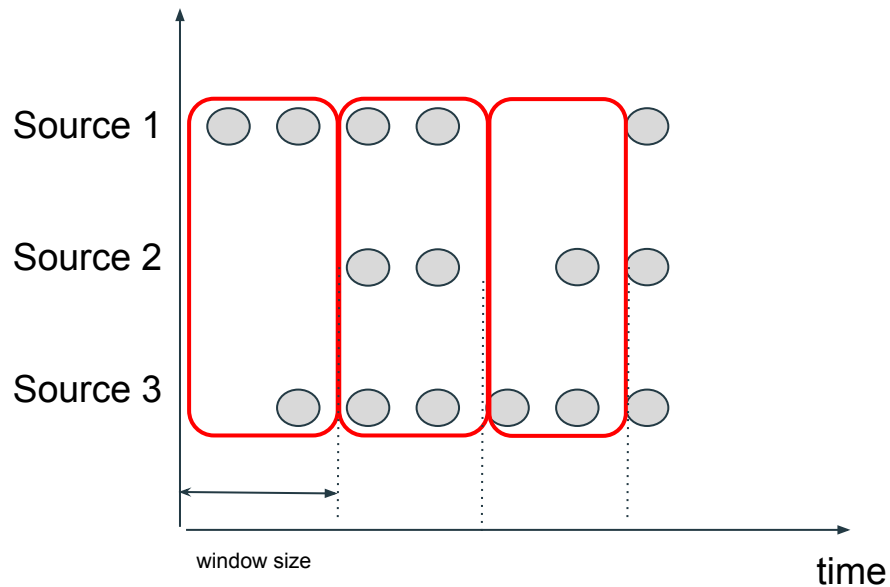
Window

Streaming engine tool.

Set by a time or size variable.

Has different subtypes: Tumbling x Sliding.

“Buckets” that apply computation.



Research Question

Is there an optimal value -- or group of values -- for batching technique that reduces Flink's Processing time?

Metric: Aggregation size x Flink's time

Model

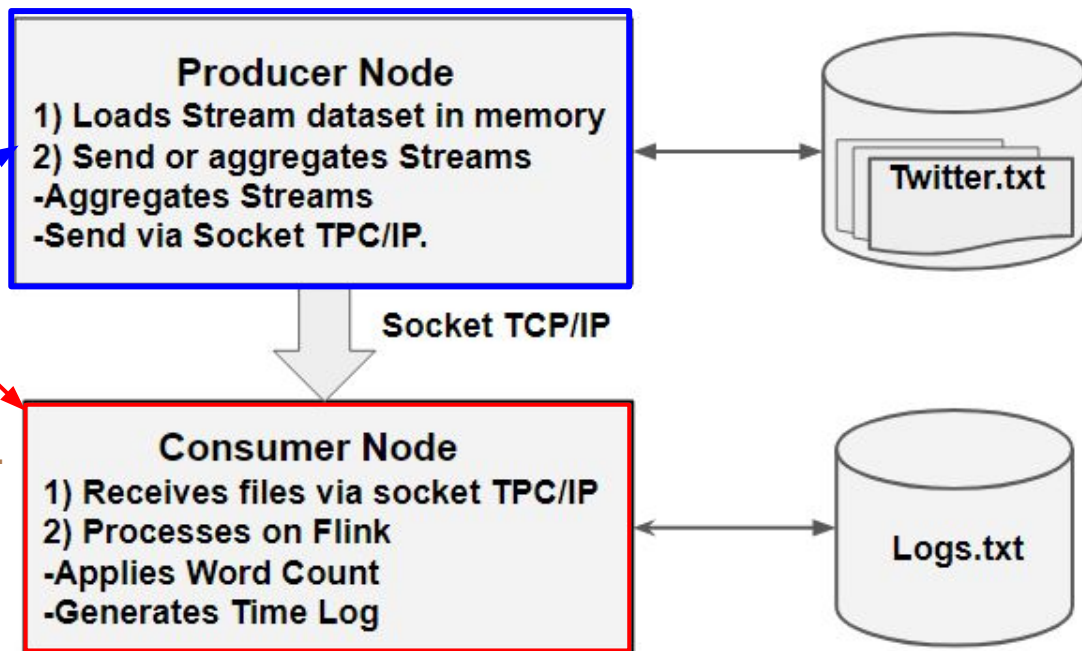
Prototype contains Python and Flink.

Python node aggregates streams.

Flink node consumes Aggregation.

Logs with time are generated.

Compares: Aggregation size x time logs.



Experiments

33 repetitions of each test , different aggregation sizes, 5 second window or not

6 strategies of aggregation, 2 datasets, 2 types of window.

had $6 \times 2 \times 2 = 24$ possibilities, 33 experiments for each configuration (so $24 \times 33 = \underline{792}$ total).

Tests executed on private Cluster: 5x Power Edge 1950

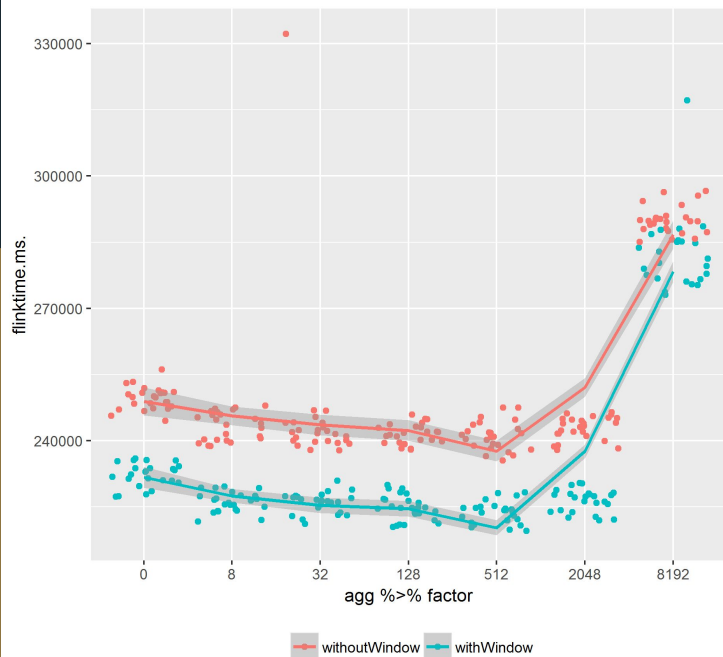
8 Cores (4 Real and 4 Virtual) - 16GB of RAM

Ubuntu Server version 16.04 LTS.

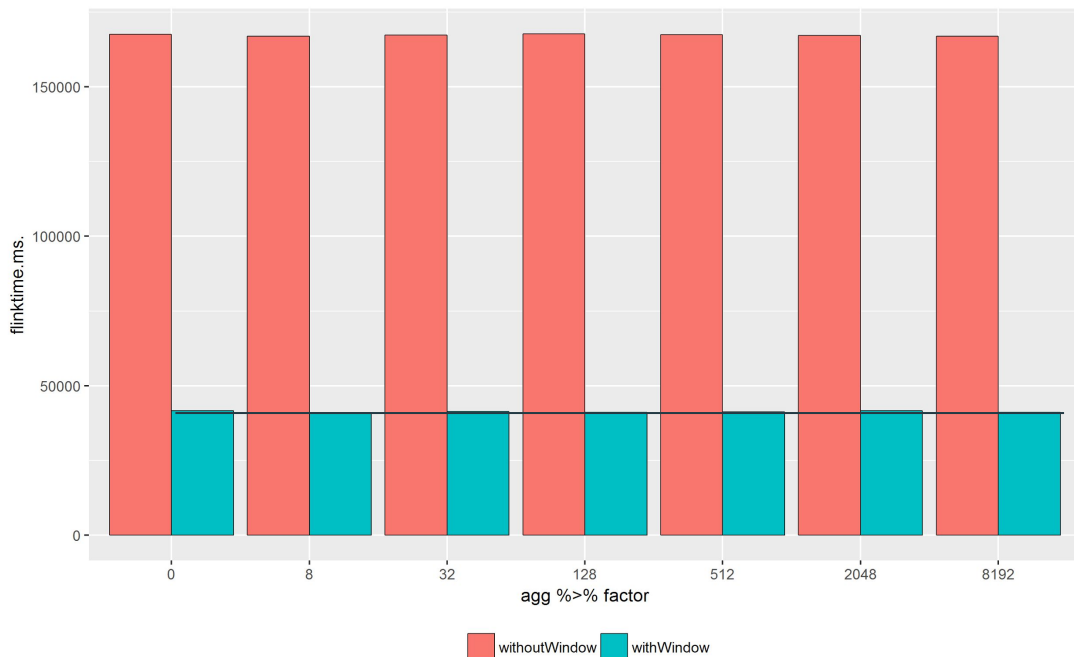
LAN Network bandwidth of 100 Mbps.

Results

Large Datasets - up to 5,806% gains
(Windowed)



Small Datasets - no significant gains.



Conclusions and Future Works

Window is always good for performance. Batching is not very significant for small sized datasets, however demonstrates performance gain in big datasets (e.g. high input flow for a long time).

Future works will add latency and throughput as new metrics and further develop the prototype so it might be able to perform dynamic readjustment of aggregation.

References

- [1] P. Basanta-Val, N. Fernandez-Garcia, L. Sanchez-Fernandez, and J. A. Fisteus. Patterns for distributed real-time stream processing. *IEEE Transactions on Parallel and Distributed Systems*, 2017.
- [2] N. Duffield, Y. Xu, L. Xia, N. Ahmed, and M. Yu. Stream aggregation through order sampling. arXiv preprint arXiv:1703.02693, 2017.
- [3] R. Guerraoui, E. Le Merrer, R. Patra, and B.-D. Tran. Frugal topology construction for stream aggregation in the cloud. In *Computer Communications, IEEE INFOCOM 2016-The 35th Annual IEEE International Conference on*, pages 1–9. Ieee, 2016.
- [4] Q. Huang and P. P. Lee. Ld-sketch: A distributed sketching design for accurate and scalable anomaly detection in network data streams. In *INFOCOM, 2014 Proceedings IEEE pages 1420–1428*. IEEE, 2014.
- [5] R. Huebsch, M. Garofalakis, J. M. Hellerstein, and I. Stoica. Sharing aggregate computation for distributed queries. In *Proceedings of the 2007 ACM SIGMOD international conference on Management of data*, pages 485–496. ACM, 2007.
- [6] S. Imai, E. Blasch, A. Galli, W. Zhu, F. Lee, and C. A. Varela. Airplane flight safety using error-tolerant data stream processing. *IEEE Aerospace and Electronic Systems Magazine* 32(4):4–17, 2017.
- [7] D. Noyes. The top 20 valuable facebook statistics—updated february 2015. Retrieved from Zephoria: <https://zephoria.com/social-media/top-15-valuable-facebookstatistics> , 2015.
- [8] O. Papapetrou, M. Garofalakis, and A. Deligiannakis. Sketch-based querying of distributed sliding-window data streams. *Proceedings of the VLDB Endowment* , 5(10):992–1003, 2012.
- [9] R. Tudoran, A. Costan, O. Nano, I. Santos, H. Soncu, and G. Antoniu. Jetstream: Enabling high throughput live event streaming on multi-site clouds. *Future Generation Computer Systems* , 54:274–291, 2016.
- [10] K. Wahner. Wahner, kai. "real-time stream processing as game changer in a big data world with hadoop and data warehouse." internet. 2014. [Online]. Accessed em: 30-06-2017.

Thanks!

A Pre-Aggregation Strategy Towards Reduction on Job
Completion Time in
Stream Process Systems

Breno Fanchiotti Zanchetta - 00240494

Paulo R. Souza Jr, Kassiano J. Matteussi, Vinícius P. Perego

Julio C. S. Anjos, Claudio F. R. Geyer, Edison P. de Freitas

