

Adaptative performance optimization for distributed Big Data Server – Application in sky surveying

Mariem Brahem, Karine Zeitouni and Lauret Yeh

DAVID Laboratory, University of Versailles, France



Introduction

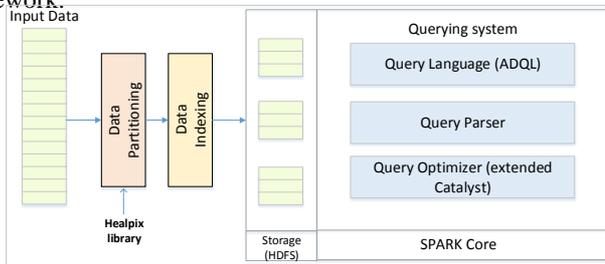
- Large amounts of astronomical data are continuously collected.
- Gaia is set to map our galaxy in three dimensions, to locate and characterize more than a billion of stars.
- Existing distributed framework (Hadoop and Spark) do not provide efficient SQL-like processing of astronomical data.
- There is a need for an efficient framework for large scale astronomical data handling.

Contributions

- We propose AstroSpark, a distributed framework specifically tailored for data intensive applications in astronomy.
- AstroSpark supports data partitioning with Healpix, a structure for the pixelization of data on the sphere, to speed up query processing.
- AstroSpark offers an expressive programming interface using a unified query language ADQL [3], a SQL-Like language improved with geometrical functions.
- AstroSpark implements a query optimizer and provides a cost based optimization module to select the best query execution plans.

AstroSpark Architecture

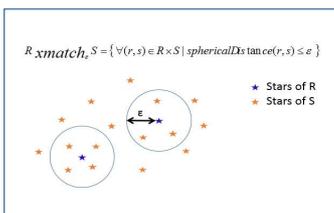
AstroSpark [1] is an extension of Spark towards a scalable, low latency, cost-effective and efficient astronomical query processing framework.



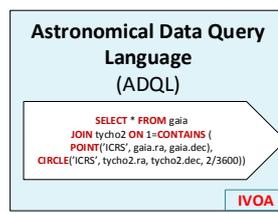
Astronomical operations:

- Cone Search returns a set of stars whose positions lie within a cone on the sky which is defined by a position and an angular distance.
- Cross-Matching aims at identifying and correlating objects belonging to different observations in order to make new scientific achievements by studying the temporal evolution of the sources or combining physical properties.
- Histogram queries distribute the dataset into a specified number of groups and summarize astronomical information about each group.

Cross-Match Definition



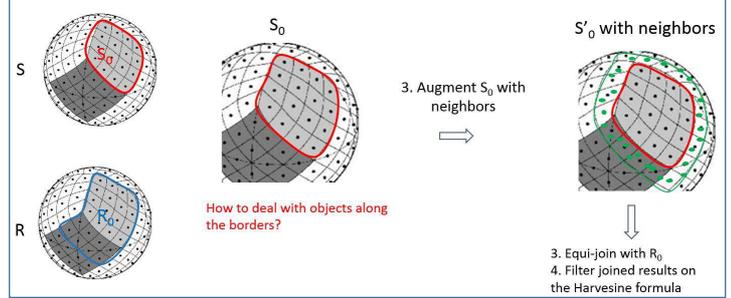
Cross-Match Query



HEALPIX [4]:

- Stands for *Hierarchical Equal Area isoLatitude Pixelization*.
- A hierarchical sky partitioning developed at NASA.
- Used as a linearization technique to transform a multi-dimensional space into a single dimension.

HX-MATCH Functioning [2]



HX-Match Algorithm [2]

1. Partitioning the two input datasets R and S using Healpix and range partitioning ;
2. Duplicate all objects in S and assign these duplicates the Healpix index of each neighbor cell -> Let's call it S' ;
3. Equi-join (R, S') on Healpix indices ;
4. Filter joined results on the Harvesine formula.

Experimental Evaluation

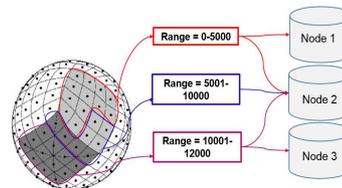
AstroSpark partitions (Our approach)



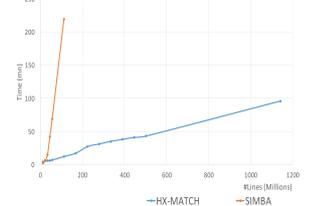
SIMBA [4] partitions (State-of-the-art)



Healpix-based Range Partitioning in AstroSpark



Cross-Matching Performance (Execution Time in AstroSpark & SIMBA[5])



Conclusion

- ✓ Design of AstroSpark, a distributed system based on Spark to process astronomical data.
- ✓ Data partitioning with Healpix to speed up query processing.
- ✓ Implementation of a cross-matching algorithm based on Spark capabilities and Healpix library.
- ✓ Extension of the Spark 2.0 Catalyst optimizer to implement the query optimizer.
- Future work:
 - Propose other algorithms for NN queries, NN join, histograms, ... with ADQL.
 - Explore other techniques of optimization: cost model assessment, multi-query optimization, caching, ...

References

- [1] Brahem, Mariem et al. *AstroSpark: towards a distributed data server for big data in astronomy*. SIGSPATIAL PhD Symposium (2016).
- [2] Brahem, Mariem et al. *HX-MATCH: In-Memory Cross-Matching Algorithm for Astronomical Big Data*, International Symposium on Spatial and Temporal Databases (SSTD17).
- [3] ADQL. <http://www.ivoa.net/documents/latest/ADQL.html>.
- [4] Gorski, Krzysztof M., et al. *HEALPix: a framework for high-resolution discretization and fast analysis of data distributed on the sphere*. The Astrophysical Journal 622.2 (2005): 759.
- [5] Xie, Dong, et al. *Simba: Efficient in-memory spatial analytics*. Proceedings of the 2016 International Conference on Management of Data. ACM, 2016.