

# Cost Model Based Approach for Graph Partitioning in Spark GraphX

Rohit Kumar, Alberto Abelló & Toon Calders

*Université Libre de Bruxelles & Universiteit Antwerpen, Belgium  
Universitat Politècnica de Catalunya (Barcelona Tech), Spain*

Large graphs with millions of nodes and billions of edges are becoming quite common today. Social media graphs, road network graphs, and relationship graphs between buyers and products are examples of large graphs generated and processed regularly. With the increase in the size of these graphs distributed graph processing is becoming more popular. One of the challenges in distributed graph processing is the problem of effective graph partitioning. There are many partitioning strategies proposed in the literature for performing efficient graph computations on distributed graph computing (DGC) systems. Despite the abundance of partitioning strategies, however, there exist relatively little guidelines for selecting the best one depending on the algorithm to run and the characteristics of the graph on which the algorithm will be run. Verma et al. in [2] attempt to address this question with an experimental comparison of different partitioning strategies on three different DGC systems. This comparison leads to many interesting insights but unfortunately, lacks theoretical justification for why one partitioning strategy outperforms another for some specific combination of graph characteristics and algorithm. Our first contribution in the presentation will be to exactly tackle this absence of a good theoretical justification by looking into a cost model for the Pregel implementation in GraphX proposed by Rohit et al. [1]. Cost models have been used for decades in the database community for query plan evaluation. If we consider graph algorithms as queries and the choice of partitioning strategy as a query plan, we contend that DGC systems should be able to choose the best partitioning strategy for a given graph and algorithm using the proposed cost model. Another challenge is that some of these large graphs are time-dependent graphs. In such situation committing to the same partitioning strategy before the start of the algorithm may not result in an optimal execution. Therefore, we look into how the cost model can be leveraged to make an adaptive system which can change the partitioning strategy by monitoring the graph evolution.

## References

- [1] Kumar, R., Abelló, A., Calders, T.: Cost model for pregel on graphx. In: Advances in Databases and Information Systems. pp. 153–166. Springer (2017)
- [2] Verma, S., Leslie, L.M., Shin, Y., Gupta, I.: An experimental comparison of partitioning strategies in distributed graph processing. VLDB (2017)