# Experimental Evaluation of Sketching Techniques for Big Spatial Data

## A. B. Siddique, Ahmed Eldawy

[msidd005,eldawy]@ucr.edu.

*Department of Computer Science and Engineering, University of California, Riverside.*
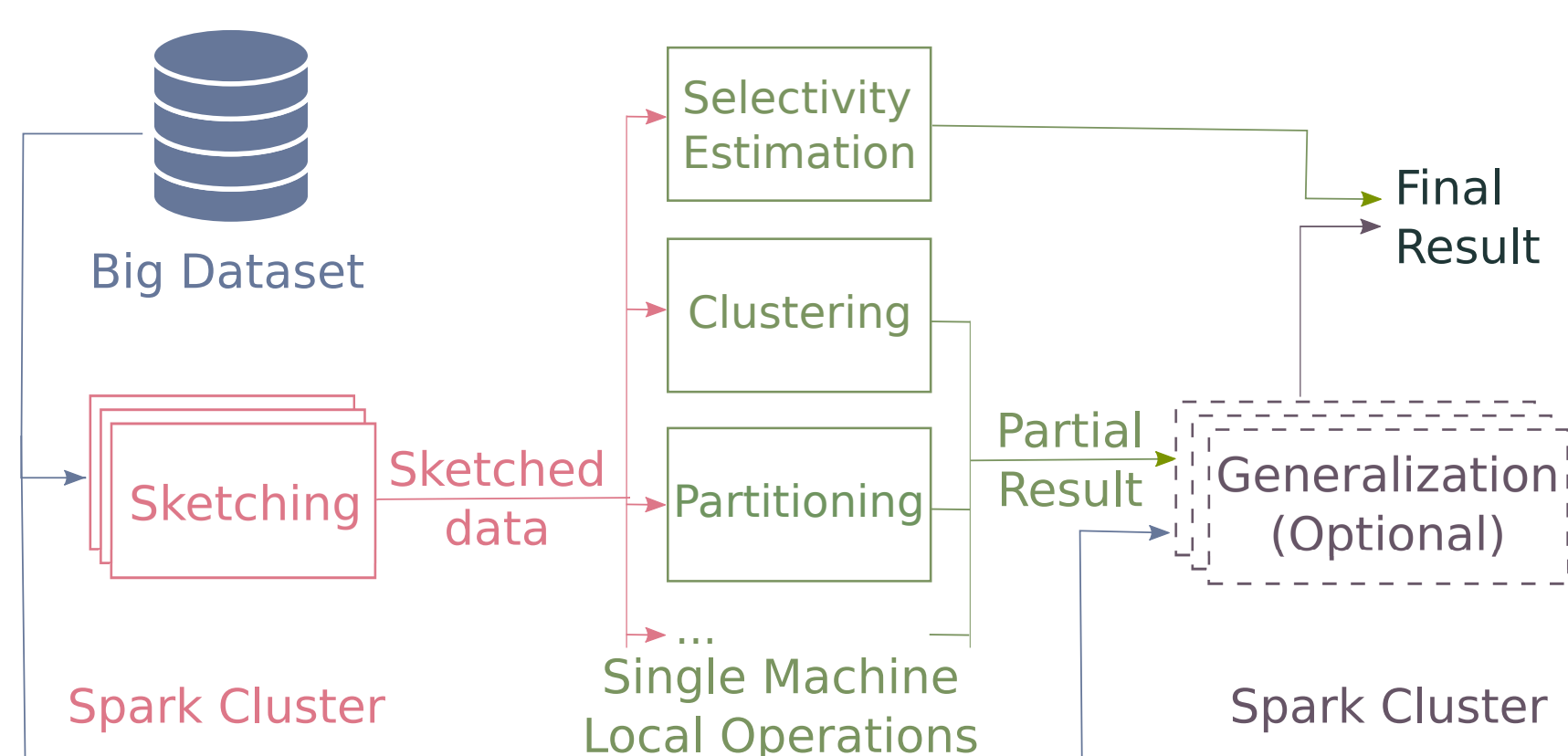
## Motivation

- Swift growth of the data
  - 2.5 exabytes of data is produced daily, of which $60 - 80\%$ is geo-referenced.
  - Space telescopes broadcast about 140 GB data weekly.



- New scalable query processing techniques are need of the hour.

- Sketching techniques excluding sampling, are not well-studied due to two challenges.
  - Hard to compare their performance.
  - Might require some tweaks to the algorithms to work.

- A comprehensive evaluation to understand the trade-offs in the different sketching techniques for big spatial data.

## Overview

- Three-phase sketching-based framework for big data processing.



- Data is sketched only once for all future local operations.
- To make the sketching methods comparable, a parameter $B$ is used.
- Local operations phase allows to reuse existing algorithm(s) with minimal changes.
- Optional generalization phase is merely a scan of the whole dataset in parallel.
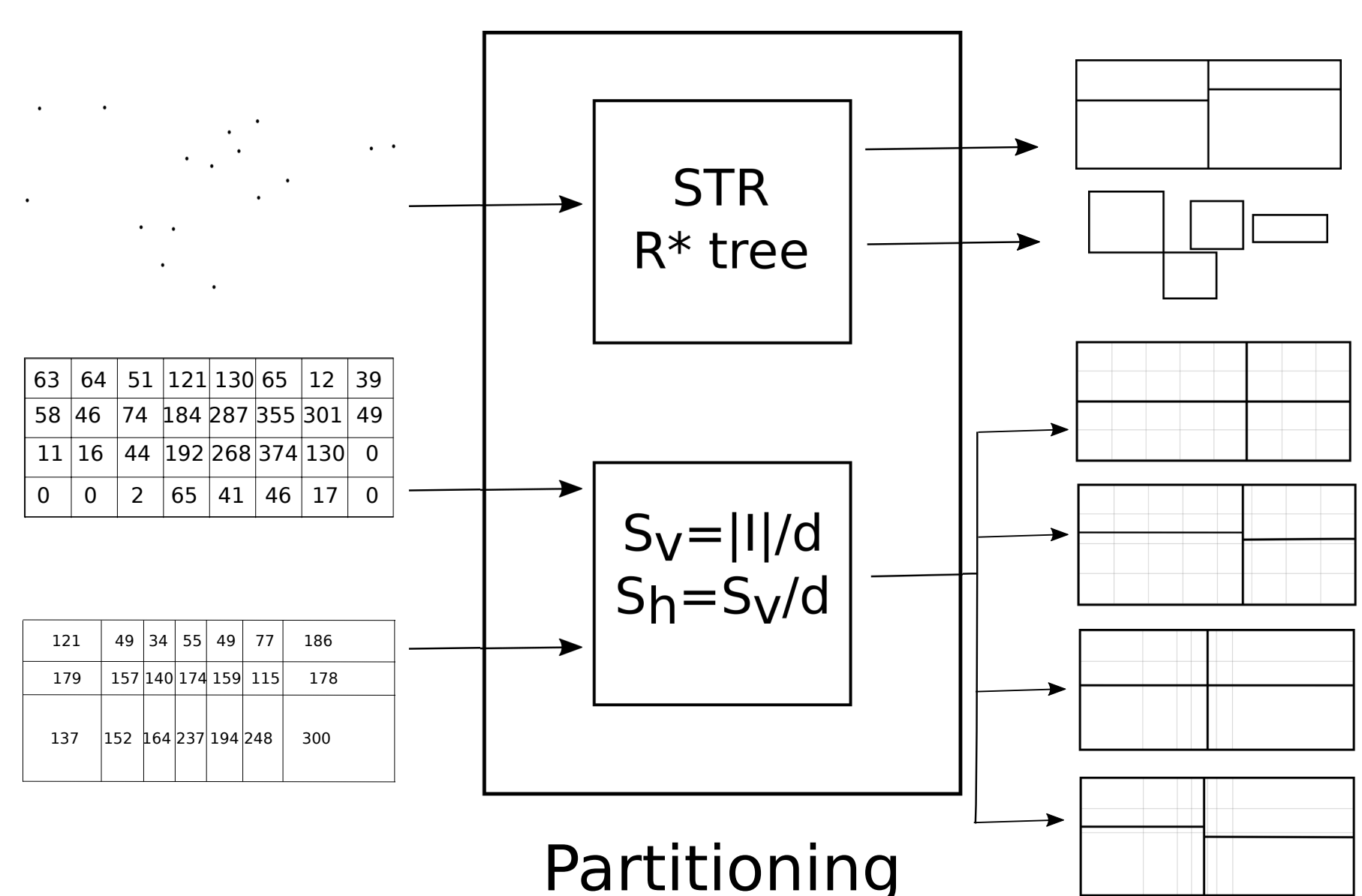
## Selectivity Estimation



### Prefix Sum



### Euler Histogram



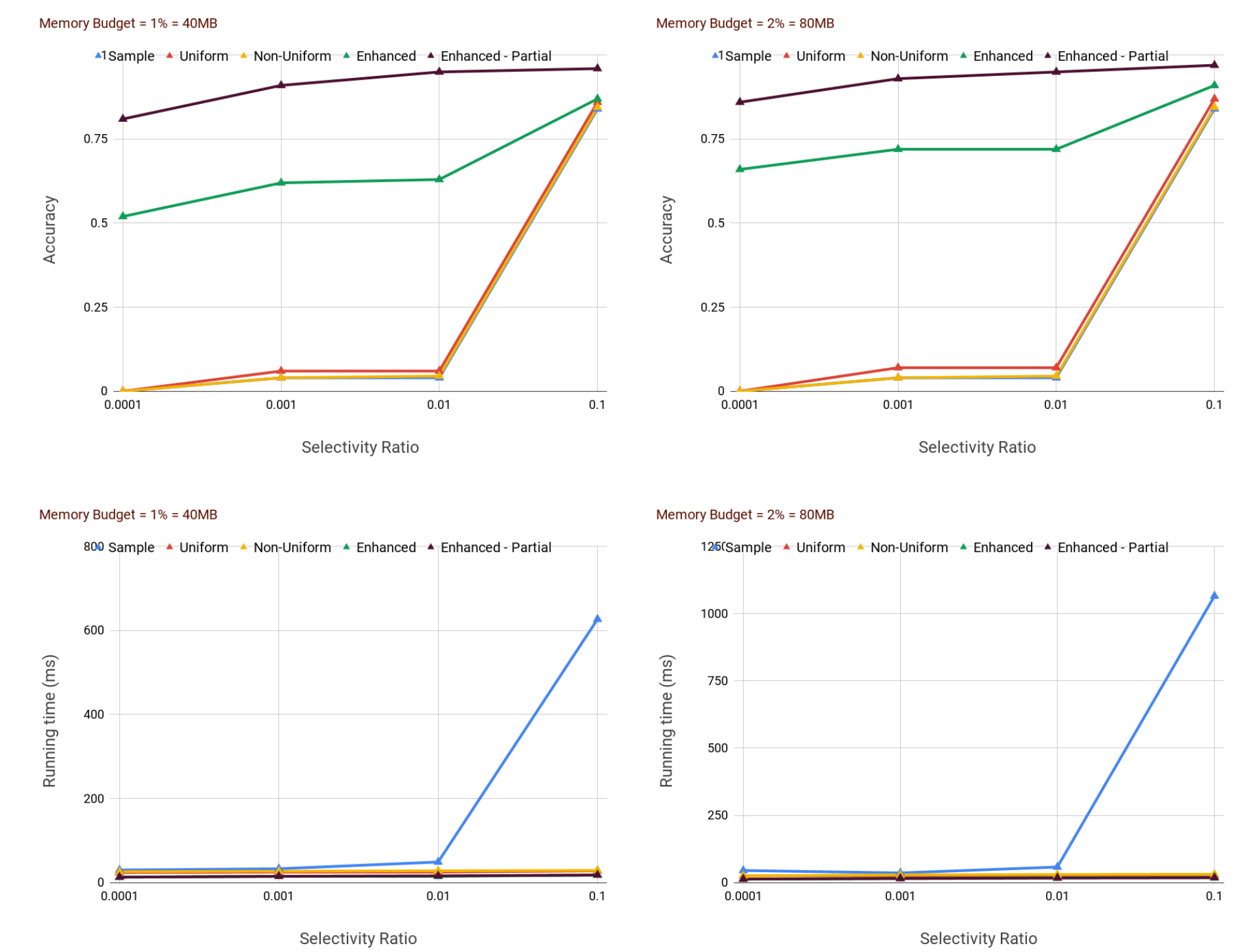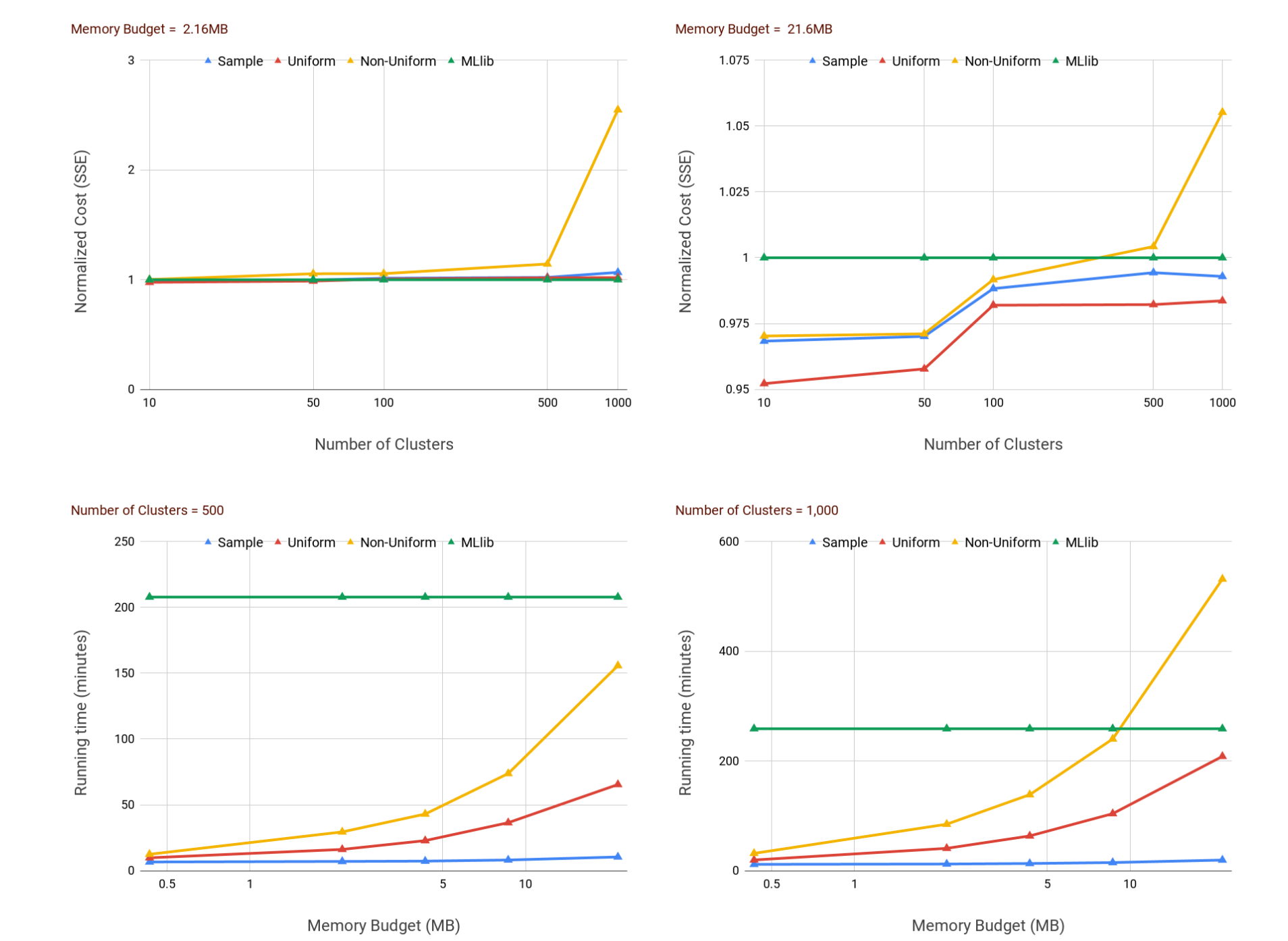## Clustering



## Partitioning

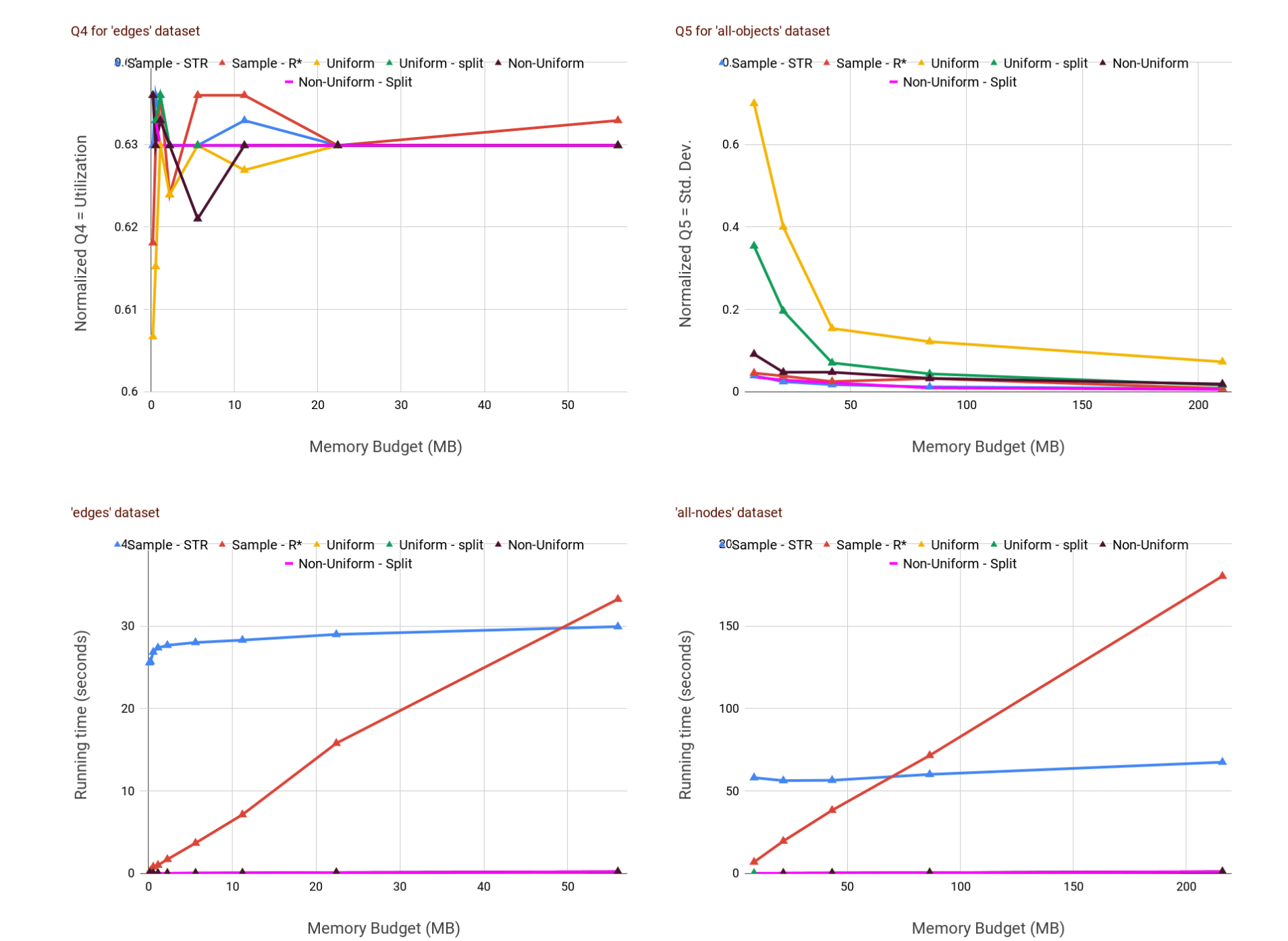

## Experimental Evaluation

- **Selectivity Estimation**



- **Clustering**



- **Partitioning**



## References

[1] Chasparis, Harry, and Ahmed Eldawy., "Experimental evaluation of selectivity estimation on big spatial data" in *Proceedings of the Fourth International ACM Workshop on Managing and Mining Enriched Geo-Spatial Data*, 2017, pp. 8. ACM.

[2] Bahmani, Bahman, et al., "Scalable k-means++" in *Proceedings of the VLDB Endowment*, 2012, pp. 622–633.

[3] Eldawy, Ahmed and Alarabi, Louai and Mokbel, Mohamed F, "Spatial partitioning techniques in SpatialHadoop" in *Proceedings of the VLDB Endowment*, 2015, pp. 1602–1605.