

Efficient Evidence Accumulation Clustering for large datasets

Diogo Silva¹, Helena Aidos² and Ana Fred²

¹Portuguese Air Force Academy, Sintra, Portugal

²Instituto de Telecomunicações, Instituto Superior Técnico, Lisbon, Portugal
dasilva@academiafa.edu.pt, {haidos, afred}@lx.it.pt

INTRODUCTION

- EAC is a robust ensemble method but its computational complexity restricts its use to small datasets.
- We propose an optimized implementation of the different EAC steps for faster execution and decreased memory usage.

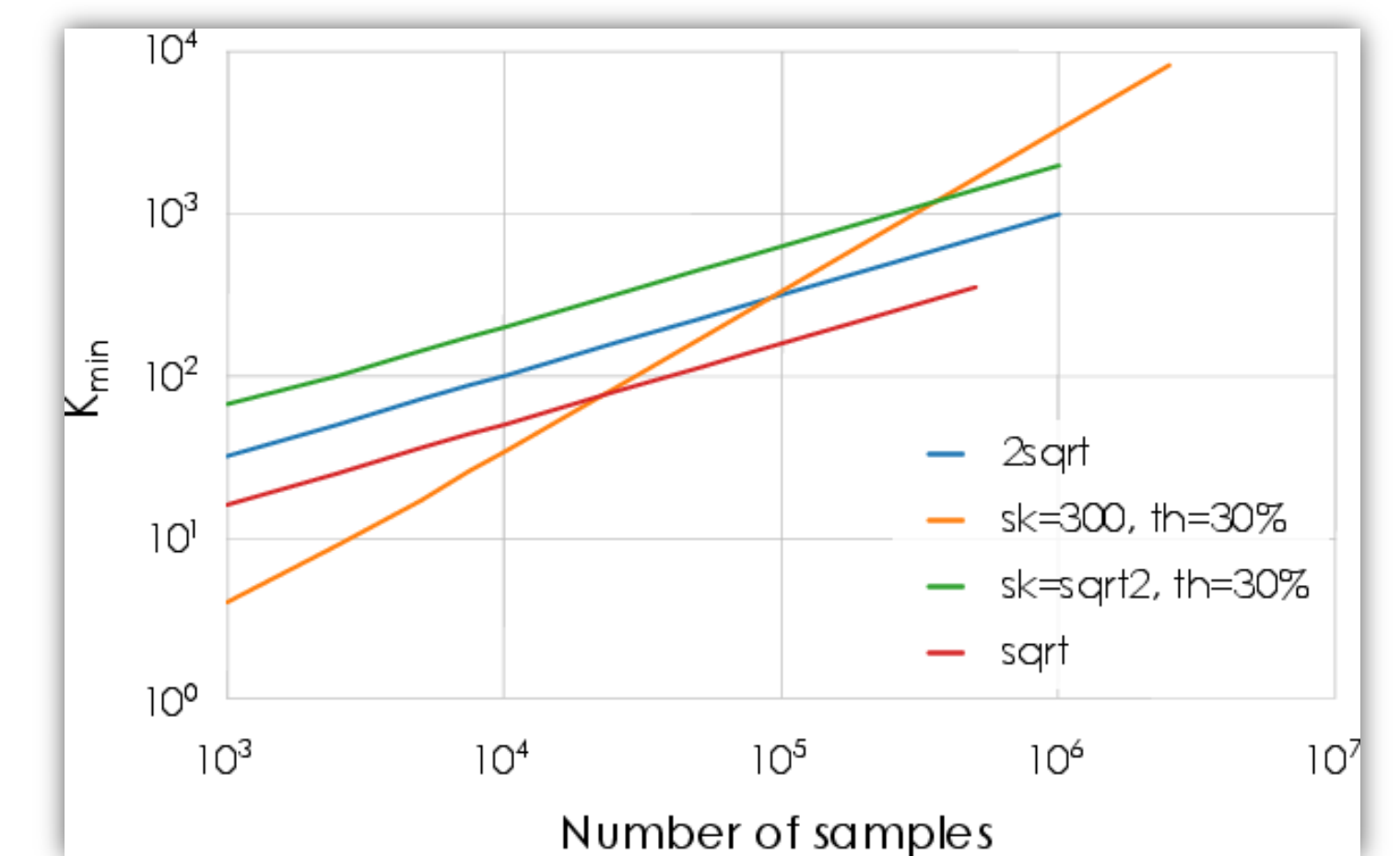
VALIDATION AND SPEED-UP

- The clustering accuracy of the optimized version relative to the original on several small benchmark datasets is negligible, validating its use on large datasets.
- Speed-up over the original version on small datasets varied between 6 and 200 on the different EAC phases.

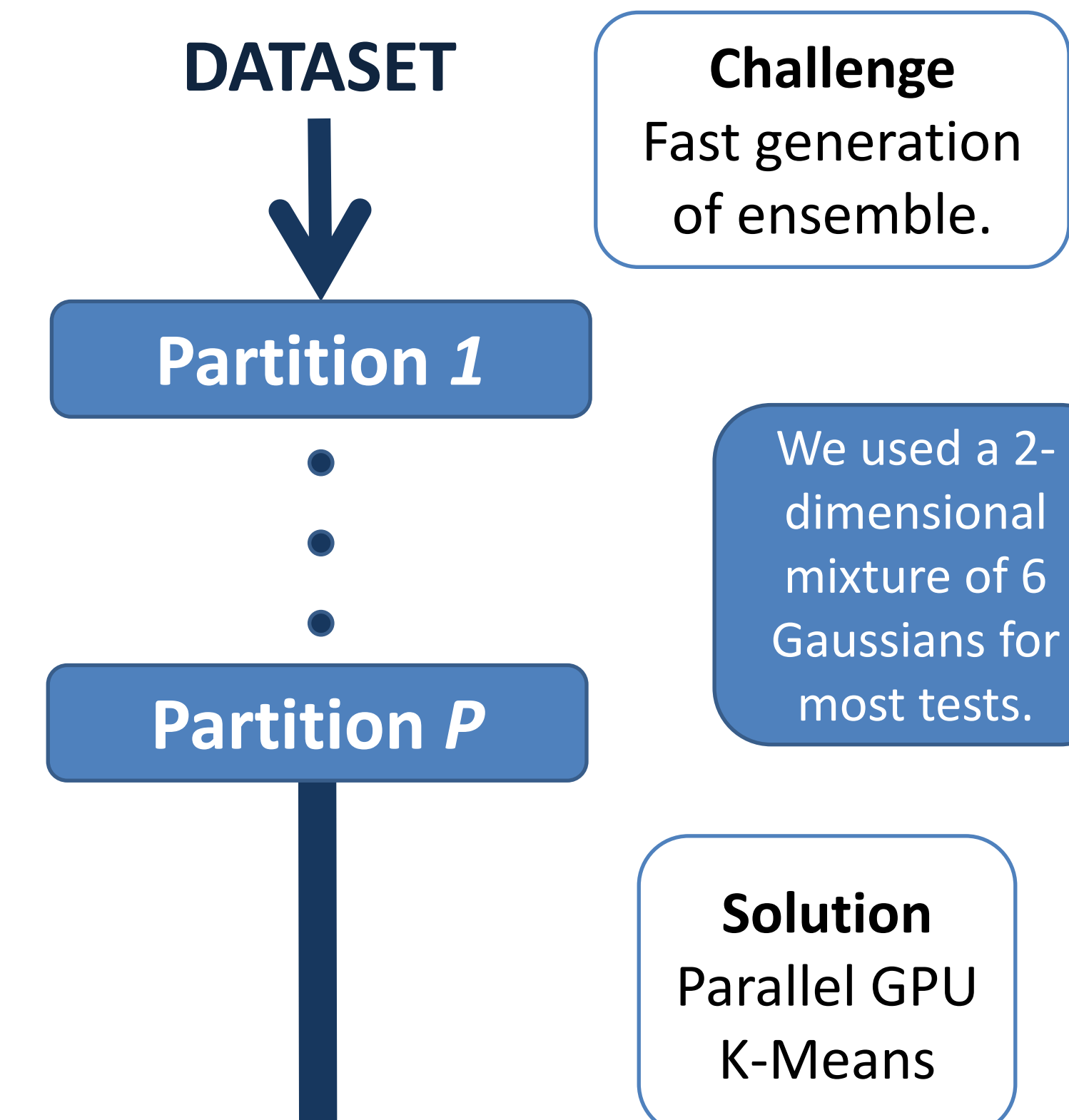
RULES FOR ENSEMBLE

The rules for the minimum and maximum number of clusters of the ensemble have a big impact on performance and memory usage. Four rules were tested.

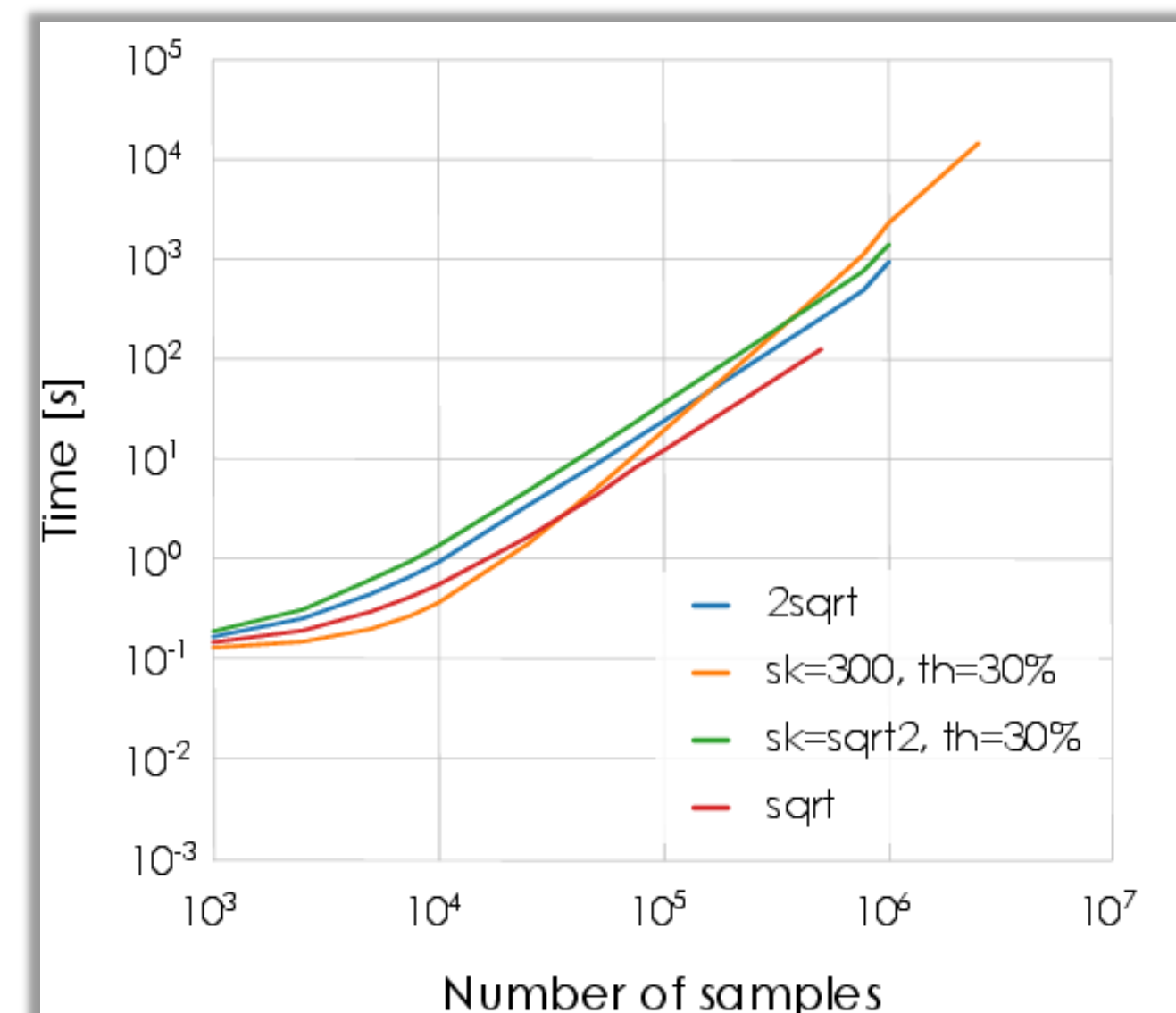
Evolution of K_{min} with different rules



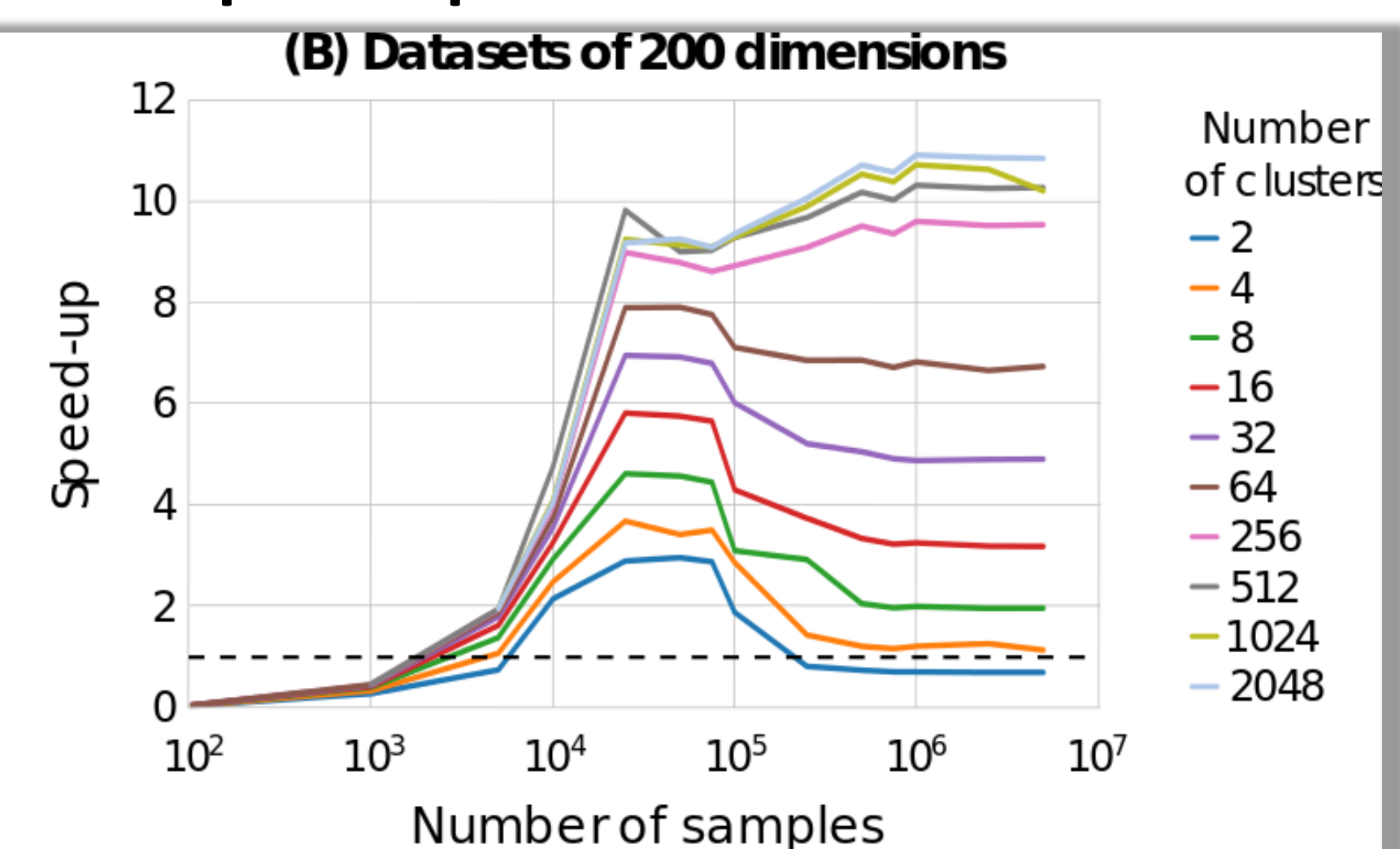
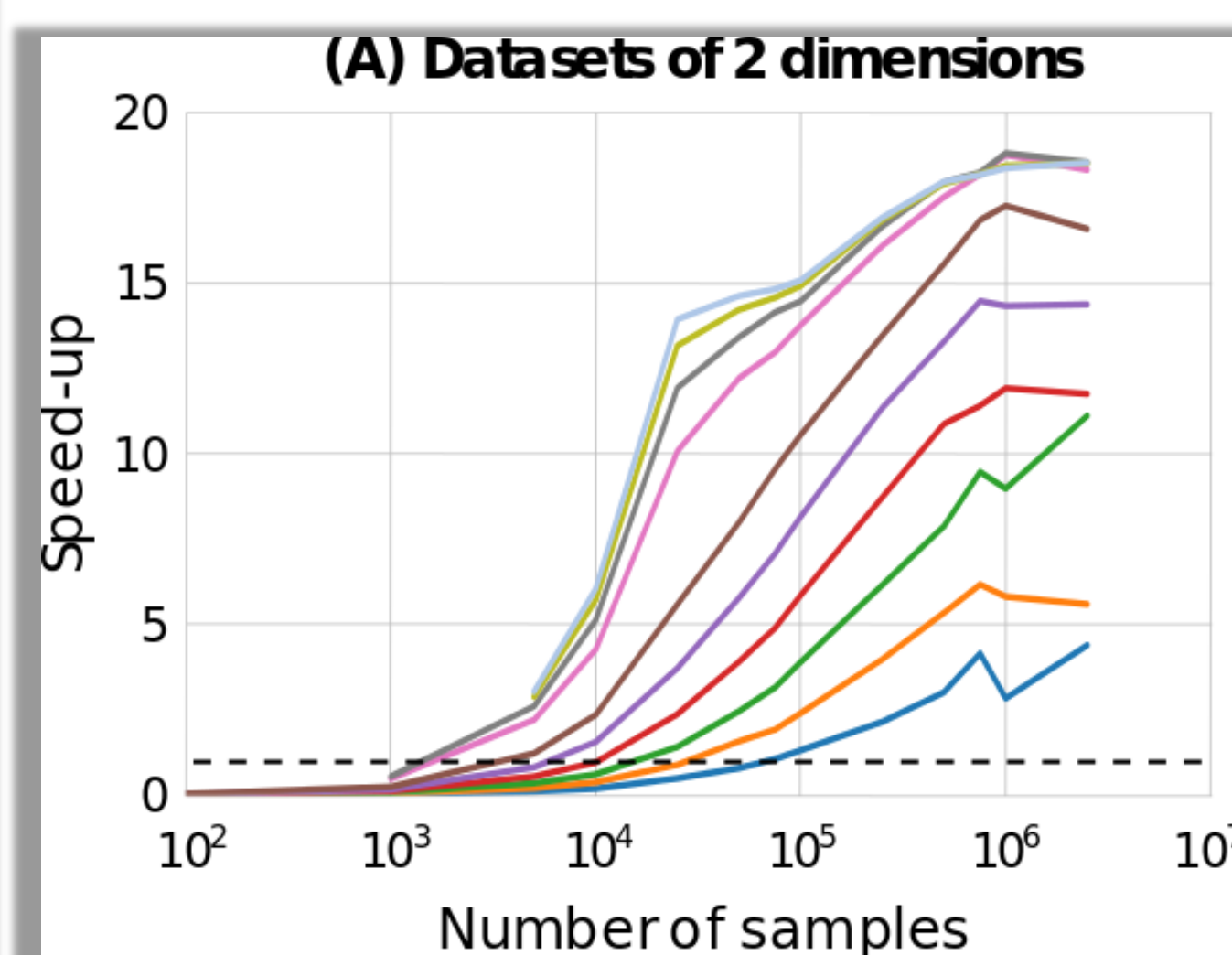
PRODUCTION OF ENSEMBLE



Production of the ensemble

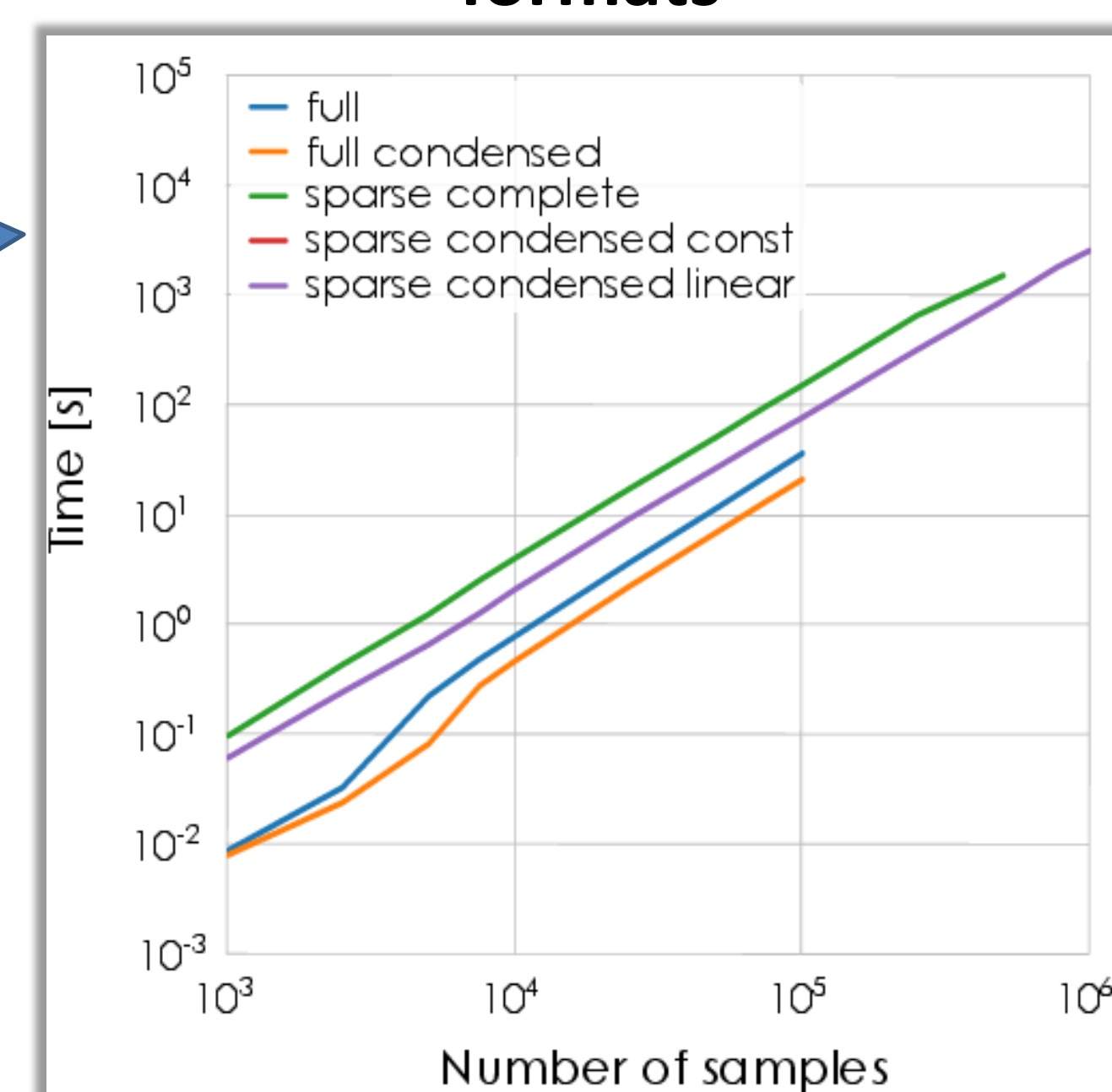


GPU K-Means speed-up

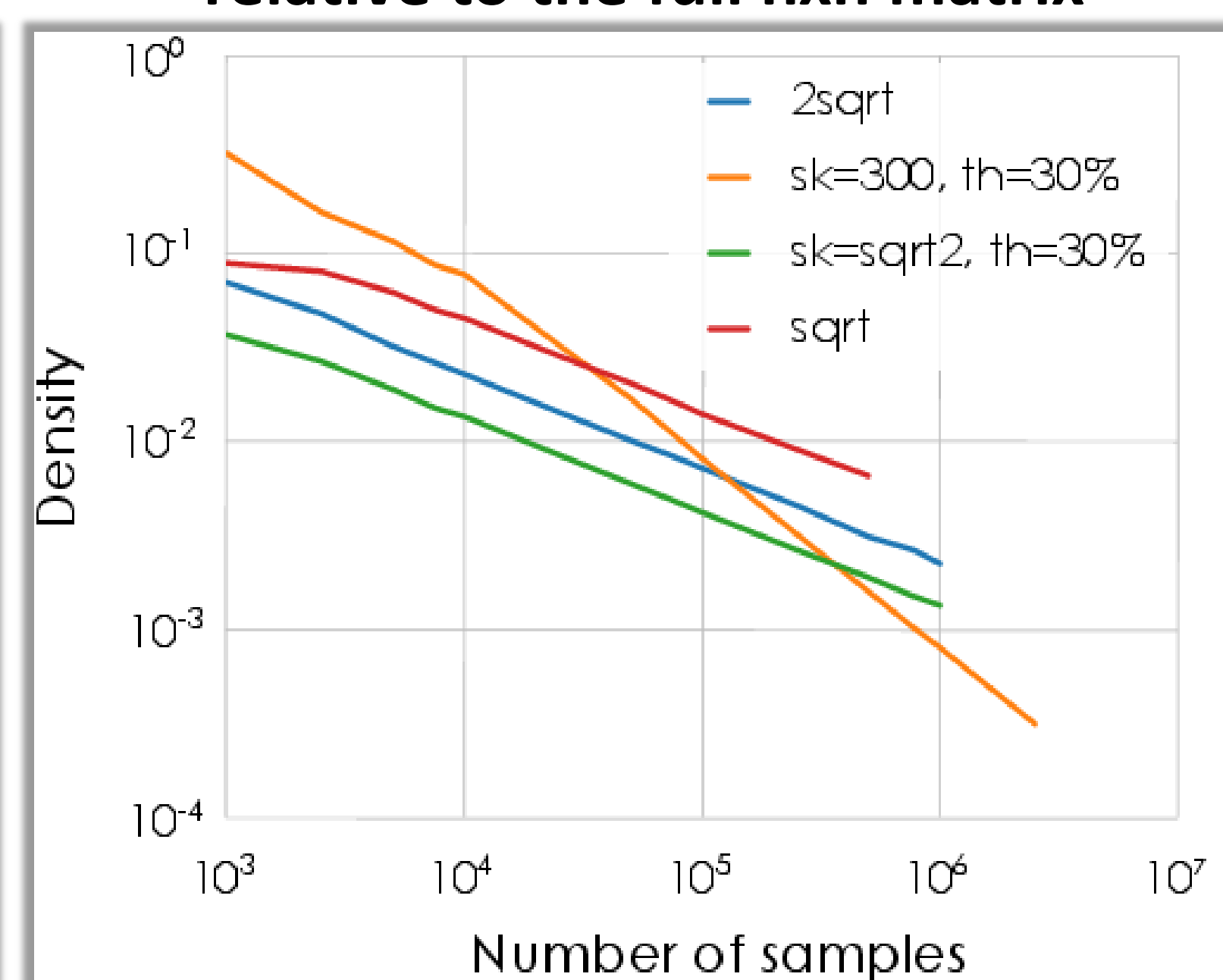


COMBINATION OF PARTITIONS

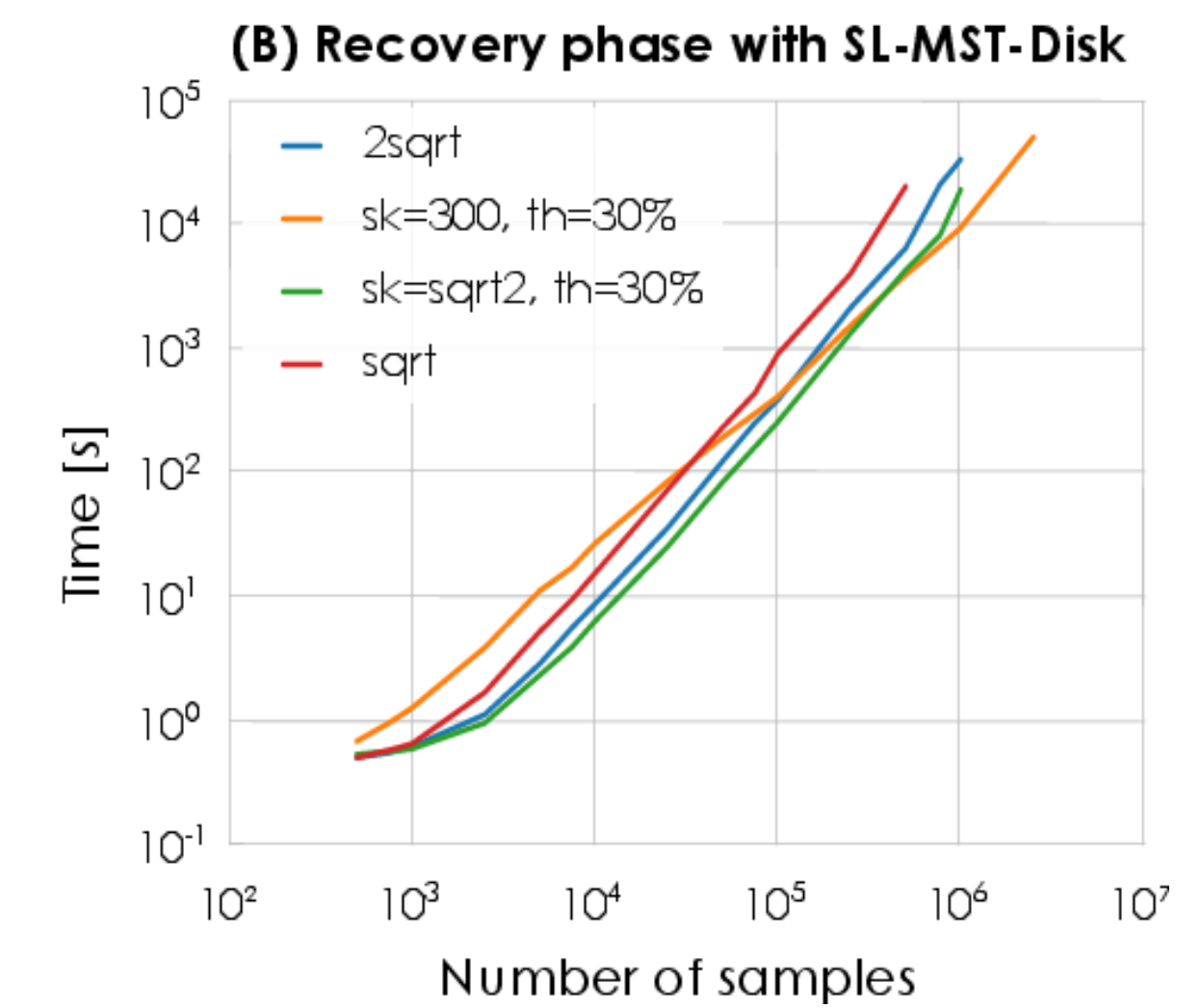
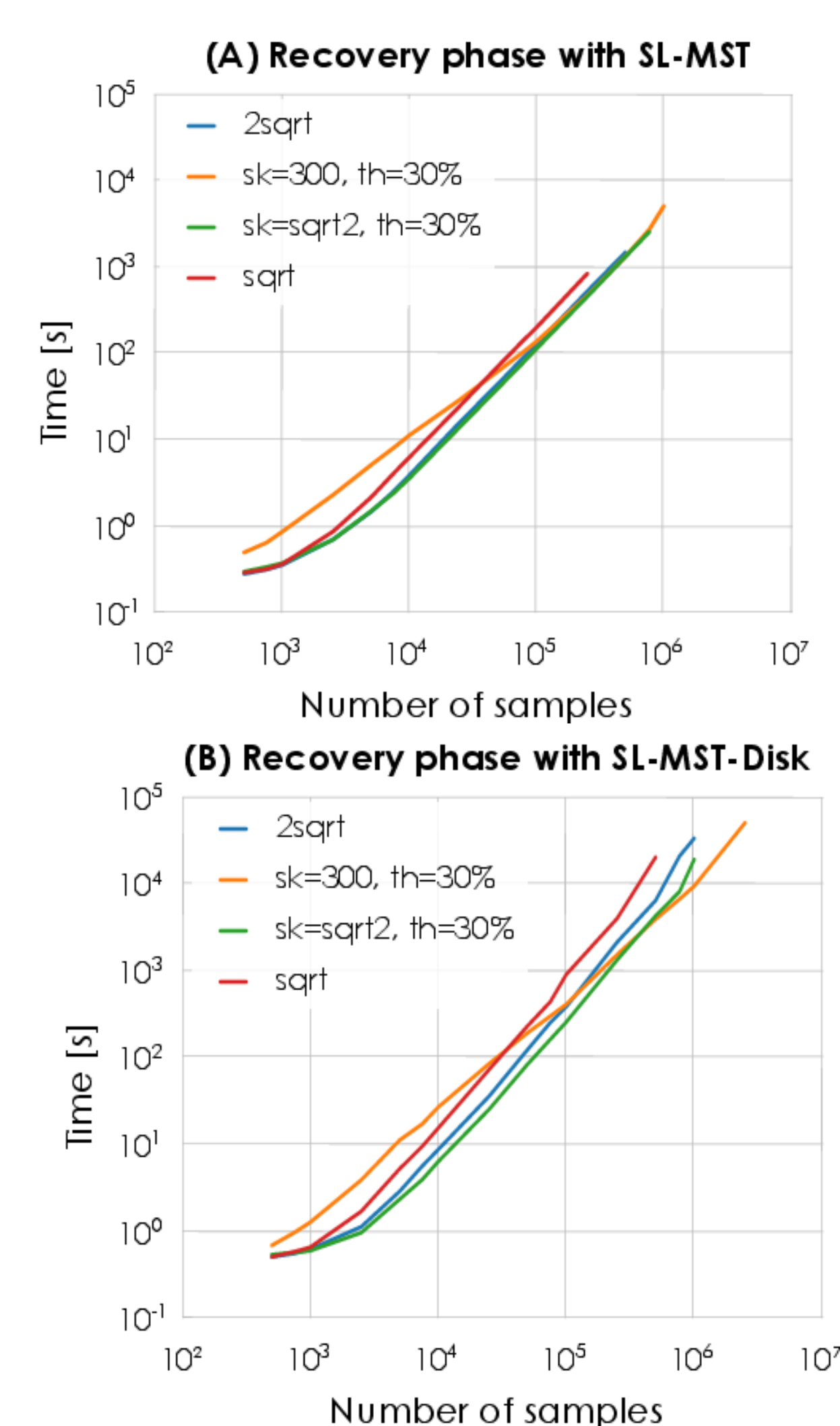
Building with different matrix formats



Density of associations relative to the full nxn matrix

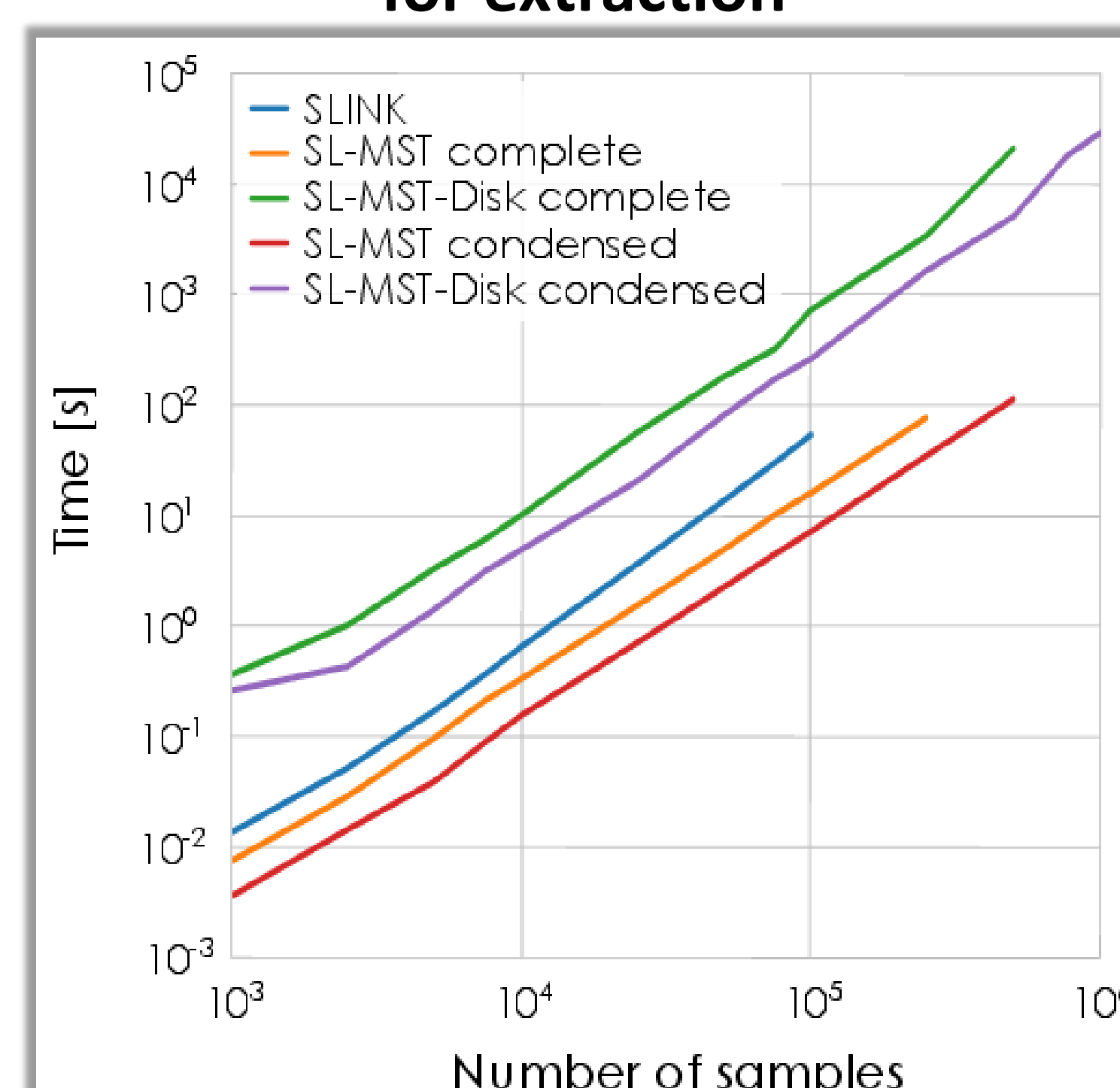


Total time



RECOVERY OF FINAL PARTITION

Comparison of three methods for extraction



CONCLUSIONS

- EAC is now applicable to a wider spectrum of datasets – we clustered datasets of up to 10 times bigger what was before possible, but the implementation supports bigger.
- Speed-up from 6 to 200 compared to original implementation on the different phases for small datasets.
- Better understanding of how ensemble rules affect the performance of the overall algorithm.