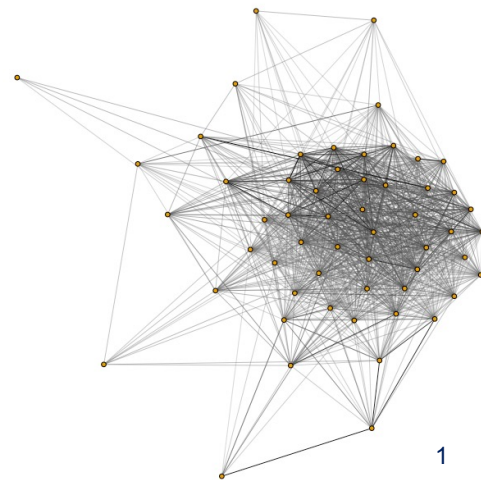
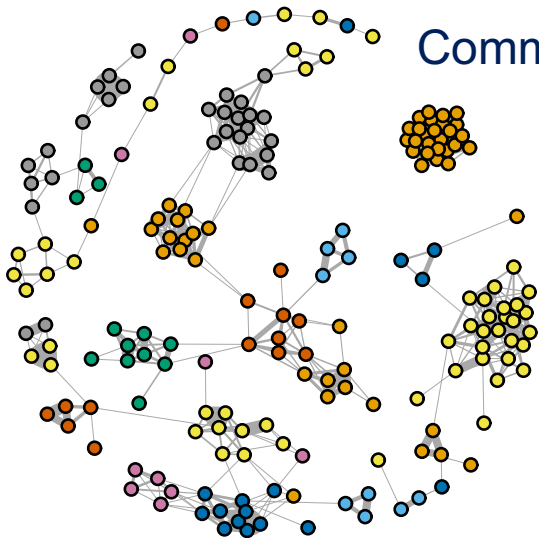
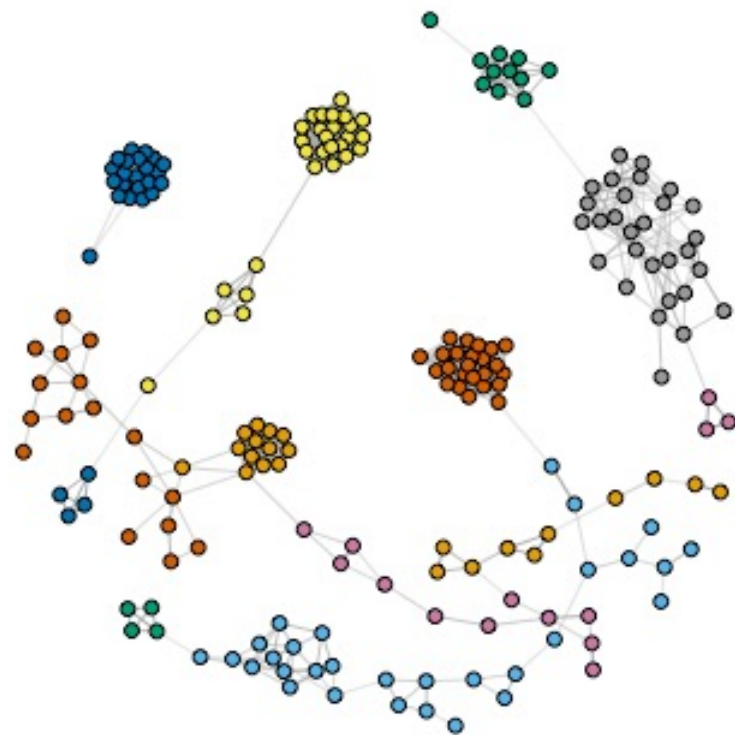
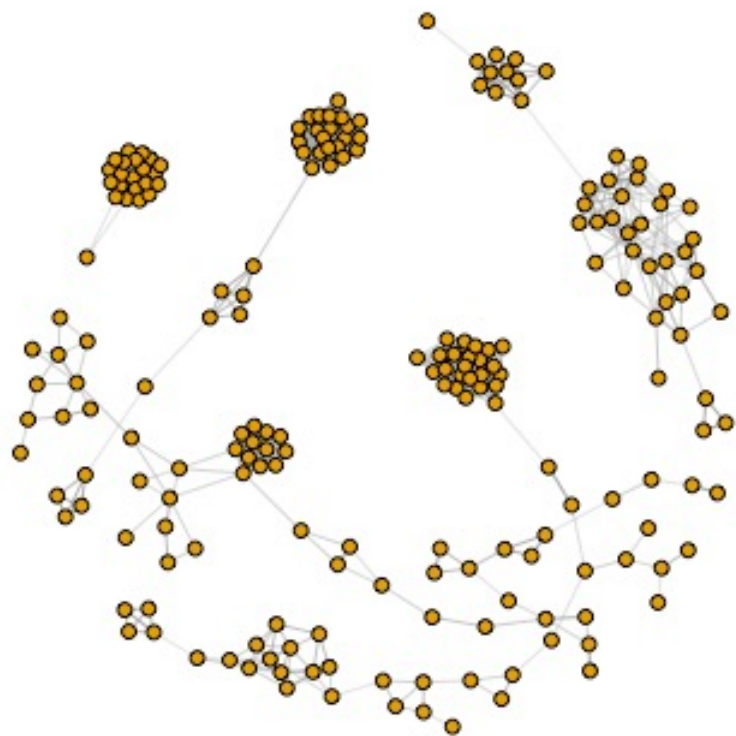


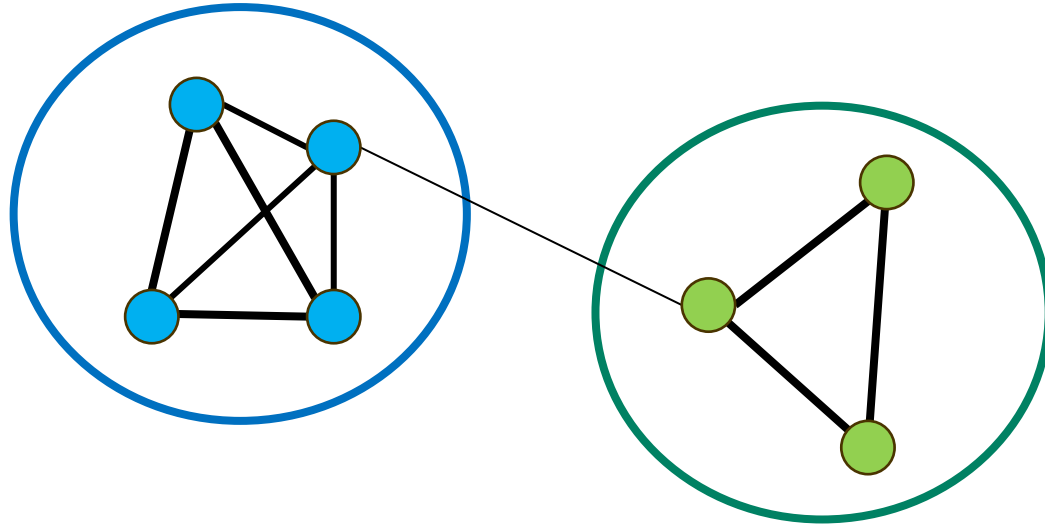
ExoLabel: Community Detection for Big Biological Networks

Aidan Lakshman, Erik S. Wright

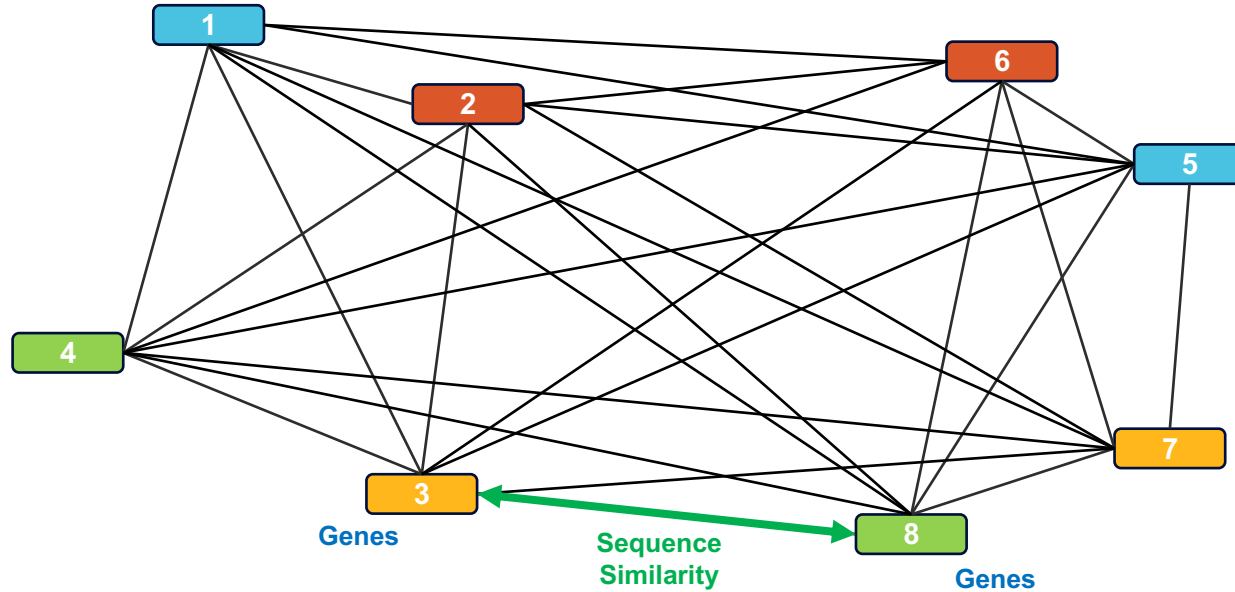




Network Community Detection Communities



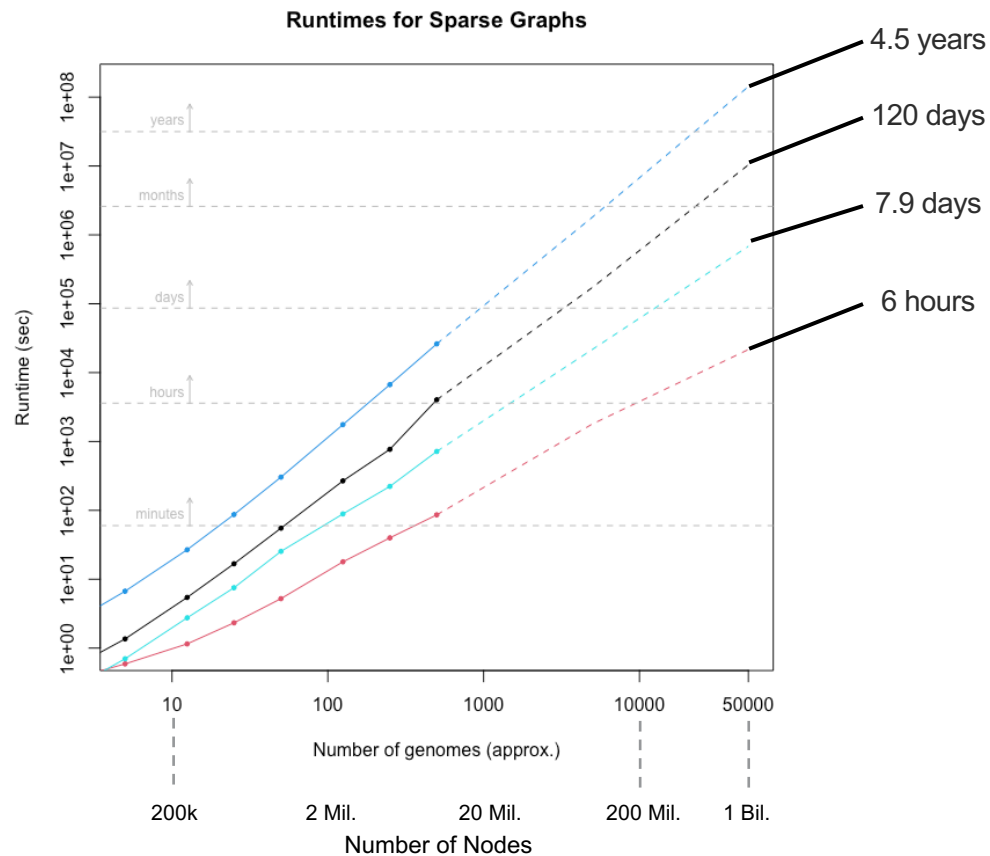
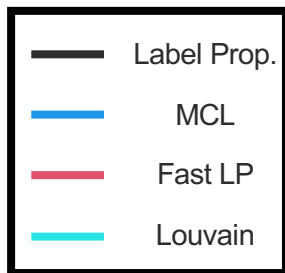
Comparative genomics depends on detecting community detection in Sequence Similarity Networks

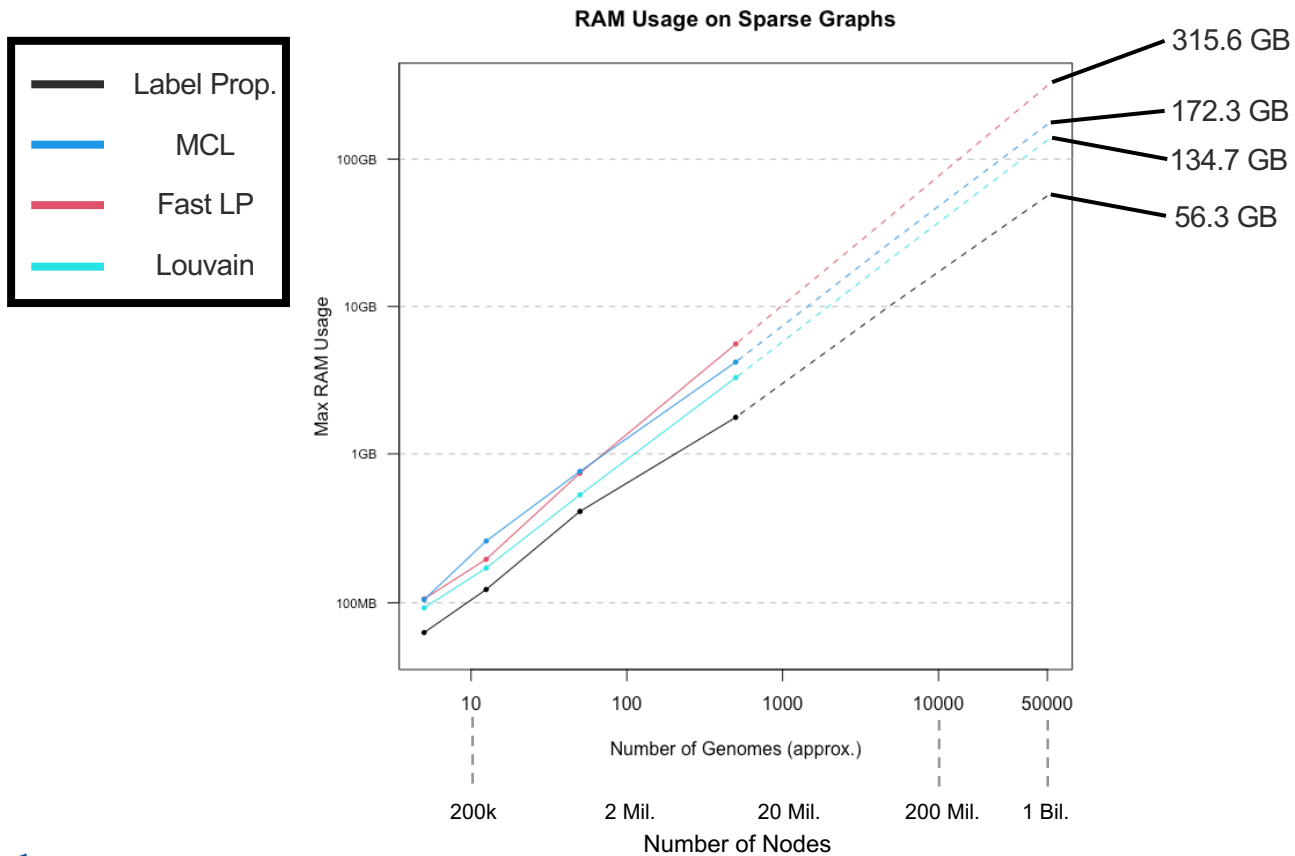




Each additional
gene → 1 edge per
existing genome

Network size is quadratic in number of genomes



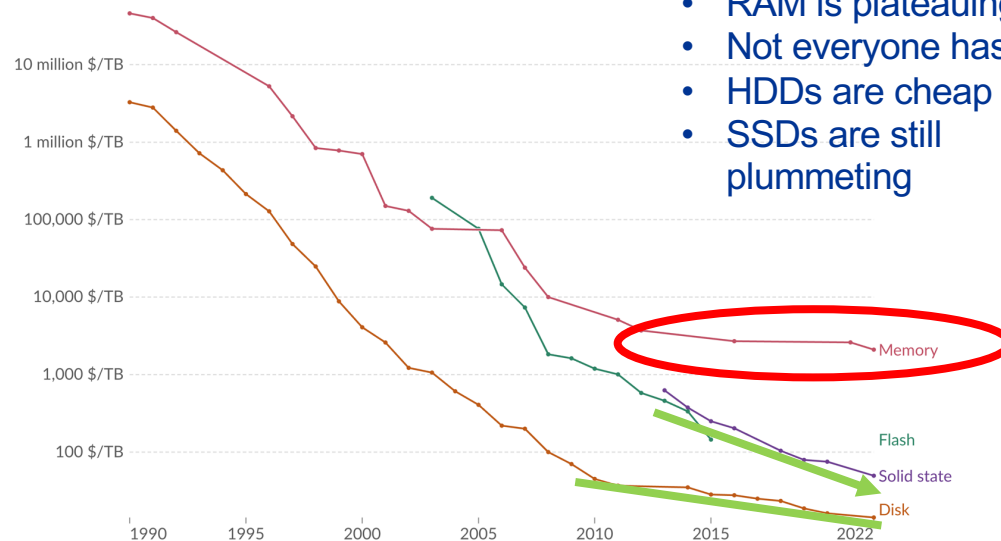


Historical cost of computer memory and storage

This data is expressed in US dollars per terabyte (TB). It is not adjusted for inflation.



- RAM is plateauing
- Not everyone has HPCs
- HDDs are cheap
- SSDs are still plummeting

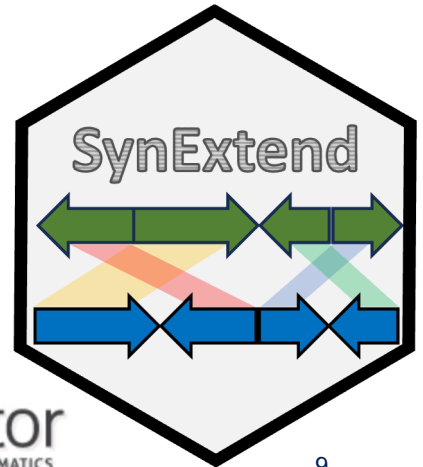


Data source: John C. McCallum (2022)

OurWorldInData.org/technological-change | CC BY

Note: For each year, the time series shows the cheapest historical price recorded until that year.

ExoLabel

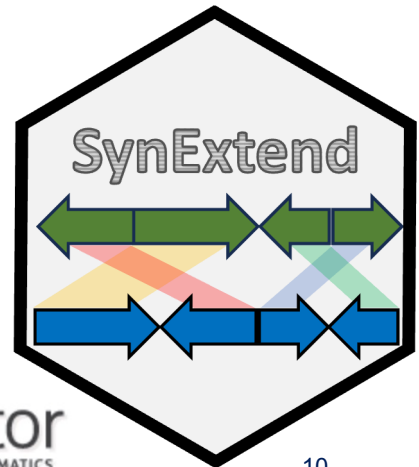


Core algorithm:
Fast Label Propagation

Label

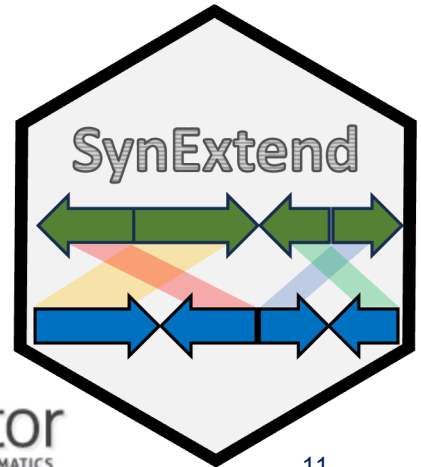
Exo

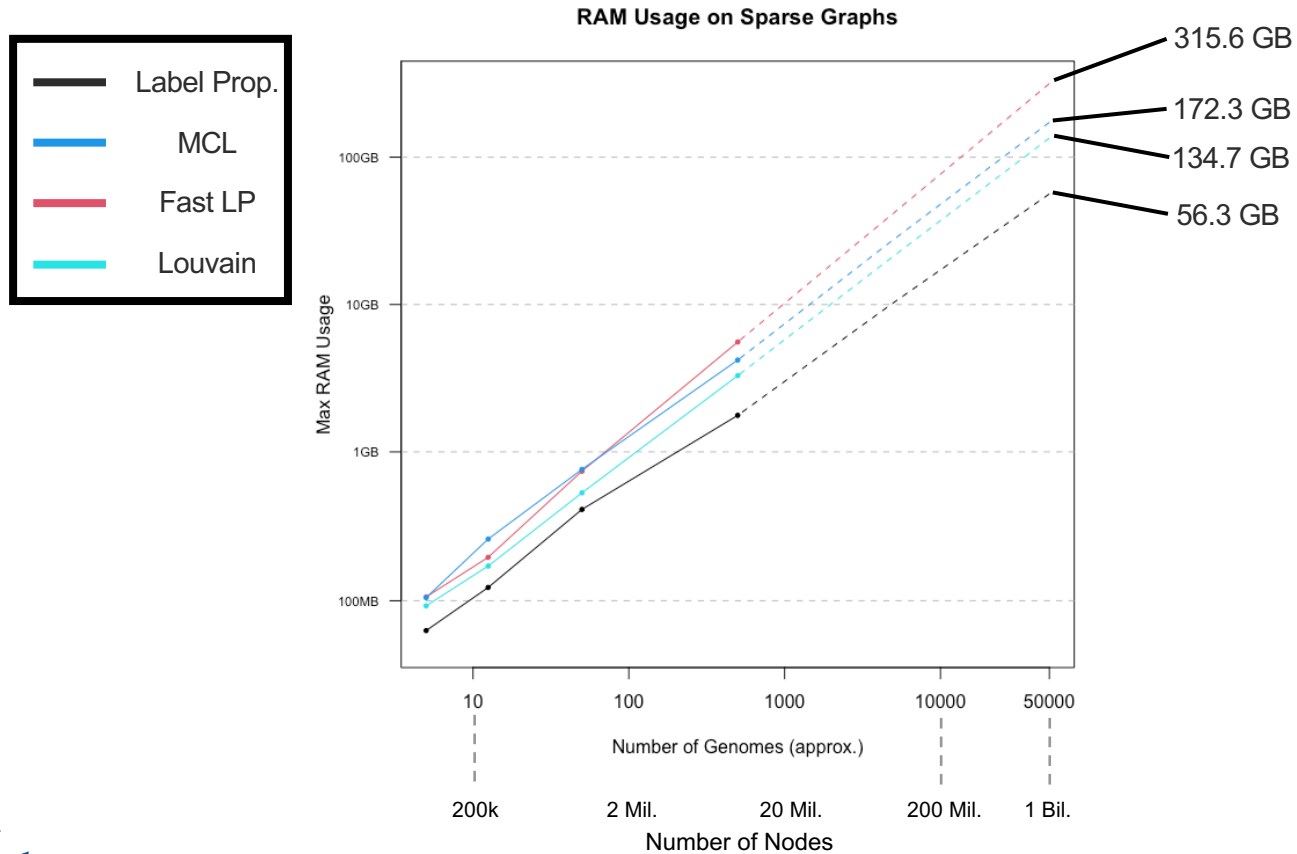
Data stored out of memory

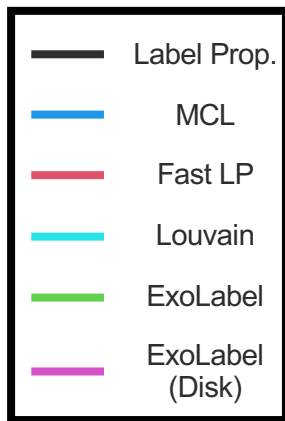


ExoLabel

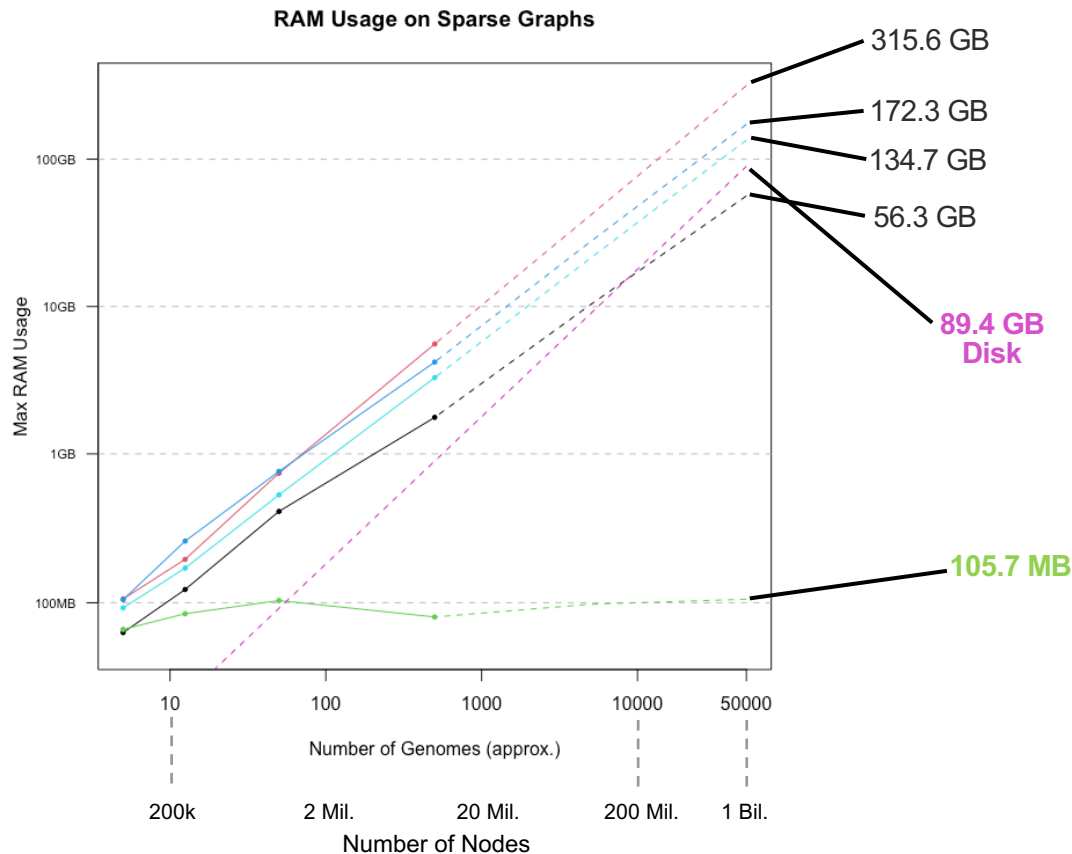
- Constant memory scaling
- Linear disk scaling
- Log-linear runtime
- Comparable accuracy to other standard methods

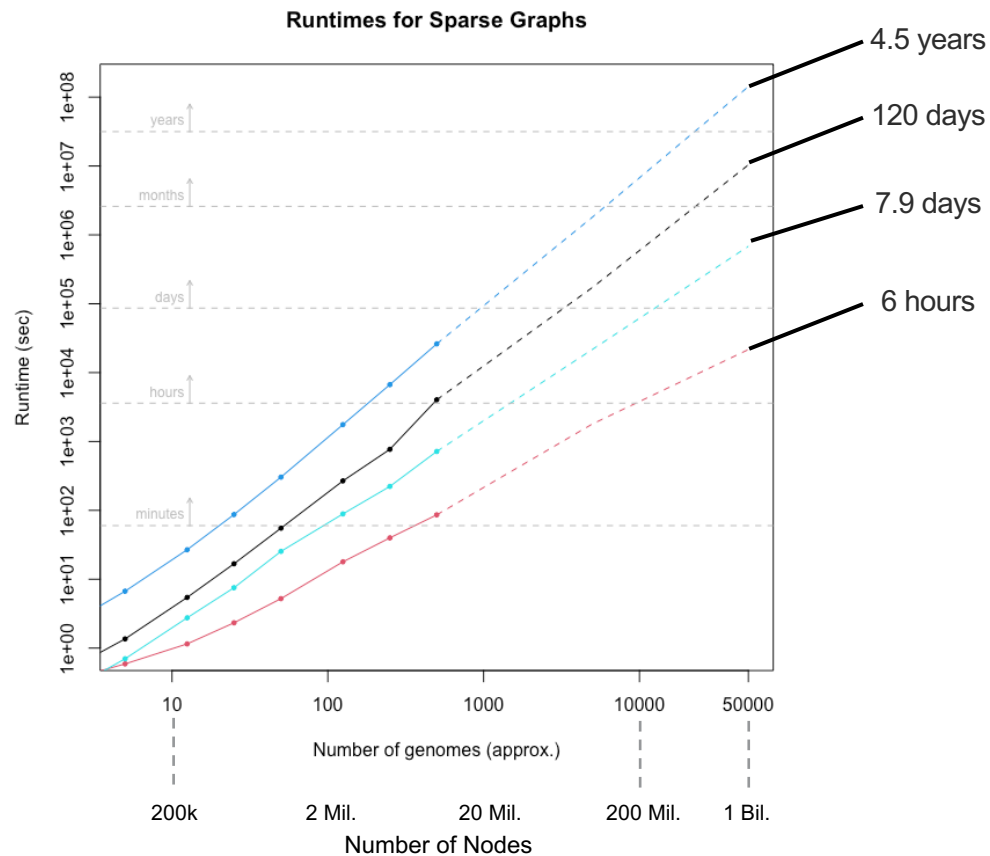
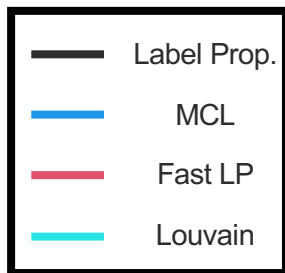


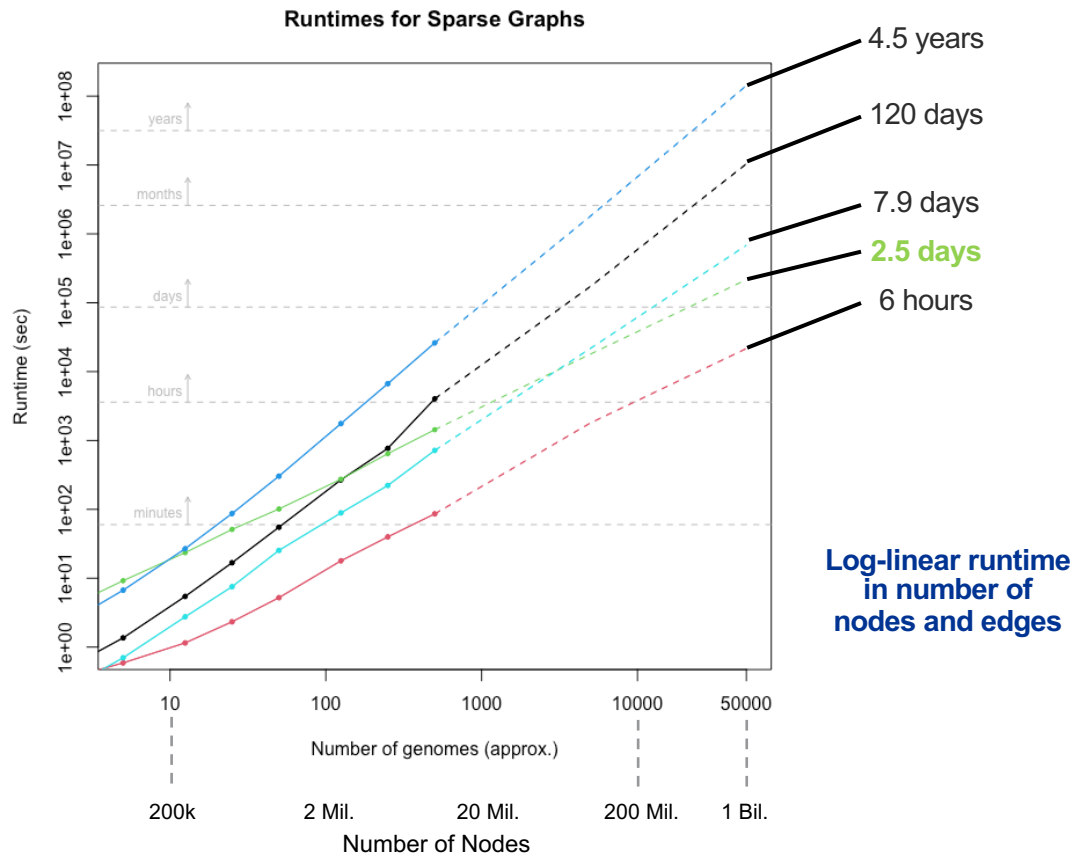
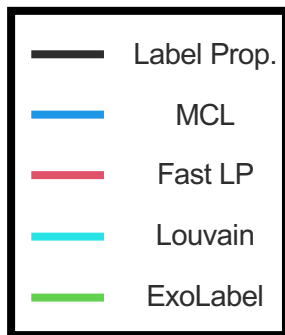




Constant RAM
Linear disk space
in number of
nodes and edges







Acknowledgements

Erik Wright, PhD

Nick Cooley, PhD

Monica Tomaszewski, PhD

Shu-Ting Cho

Sam Blechman

Nishant Panicker

Maria Bond



www.WrightLabScience.com

AHL27.com
AHL27@pitt.edu

