

Streaming Algorithms for Halo Finders

Zaoxing Liu , Nikita Ivkin , Lin F. Yang , Mark Neyrinck , Gerard Lemson, Alexander S. Szalay,
Vladimir Braverman, Tamas Budavari, Randal Burns, Xin Wang
Johns Hopkins University, Baltimore MD

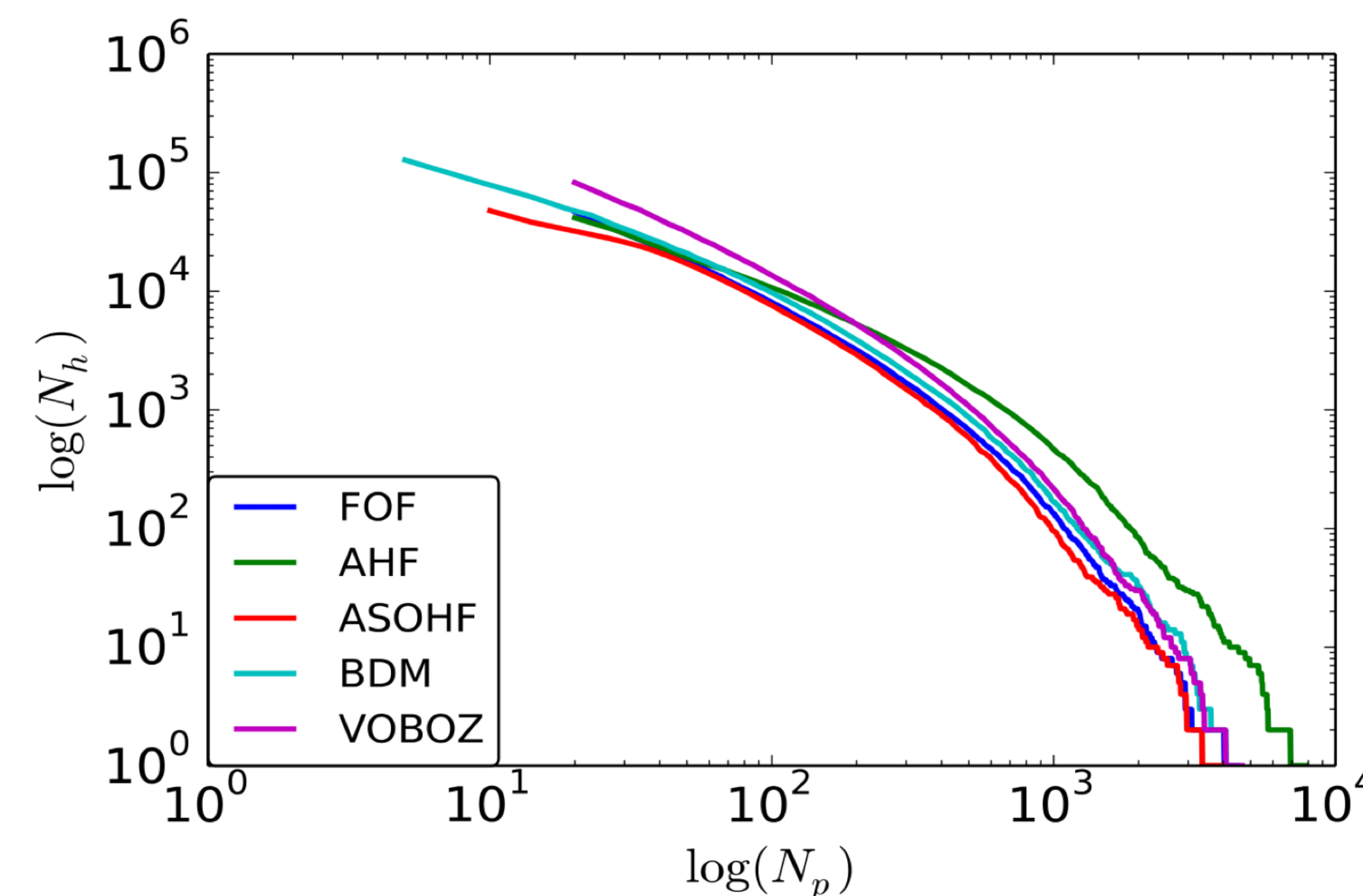
Data

- Cosmological Simulation provides positions of particles.
- Halo – macro structure with high mass concentration.
- Finding haloes is crucial to connect theory to observation.
- Some facts about data:

of particles $\sim 10^{12}$ # of haloes $\sim 10^9$ # of particles in haloes $\sim 10\%$

Distribution of Halo sizes

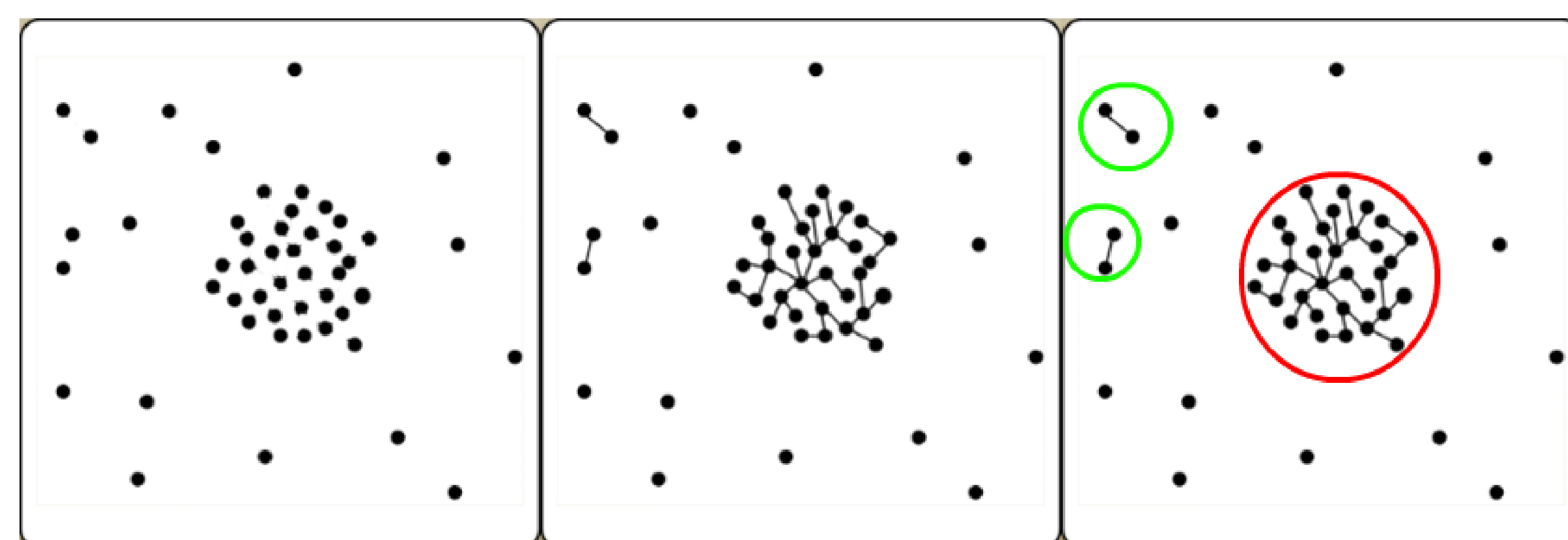
N_h – number of haloes
 N_p – number of particles in halo



Previous Methods to find Haloes

- Current solutions require to load all the data into memory (**$\sim 12\text{TB}$**).

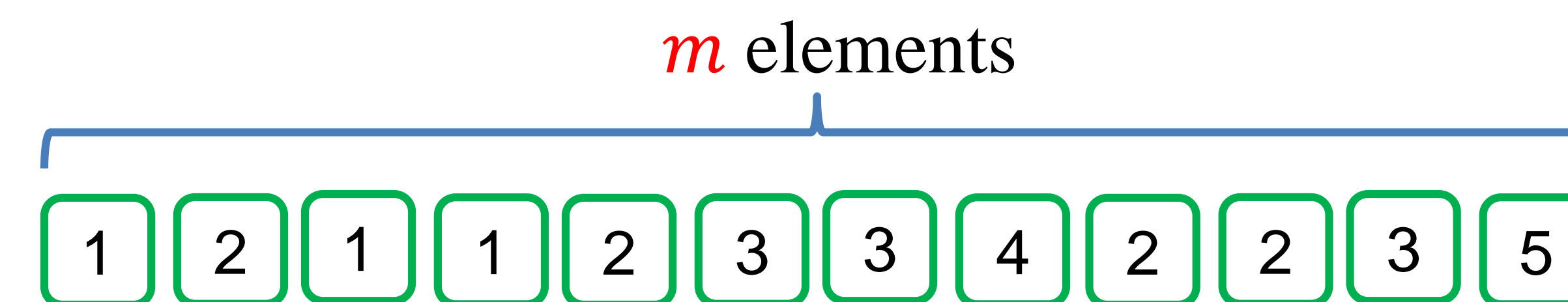
Huge memory!



Close particles are friends.
A halo is a group of friends.

Friends-of-Friends Algorithm (FOF)

Streaming Model



Stream:

- m elements from dictionary of size n
e.g. $D = \{x_1, x_2, \dots, x_m\} = 3 \ 5 \ 3 \ 7 \ 5 \ 4 \ \dots$

Goal:

- Compute the function of stream (e.g. k most frequent items) in sublinear memory.
- Approximate answer with **high probability** is OK.

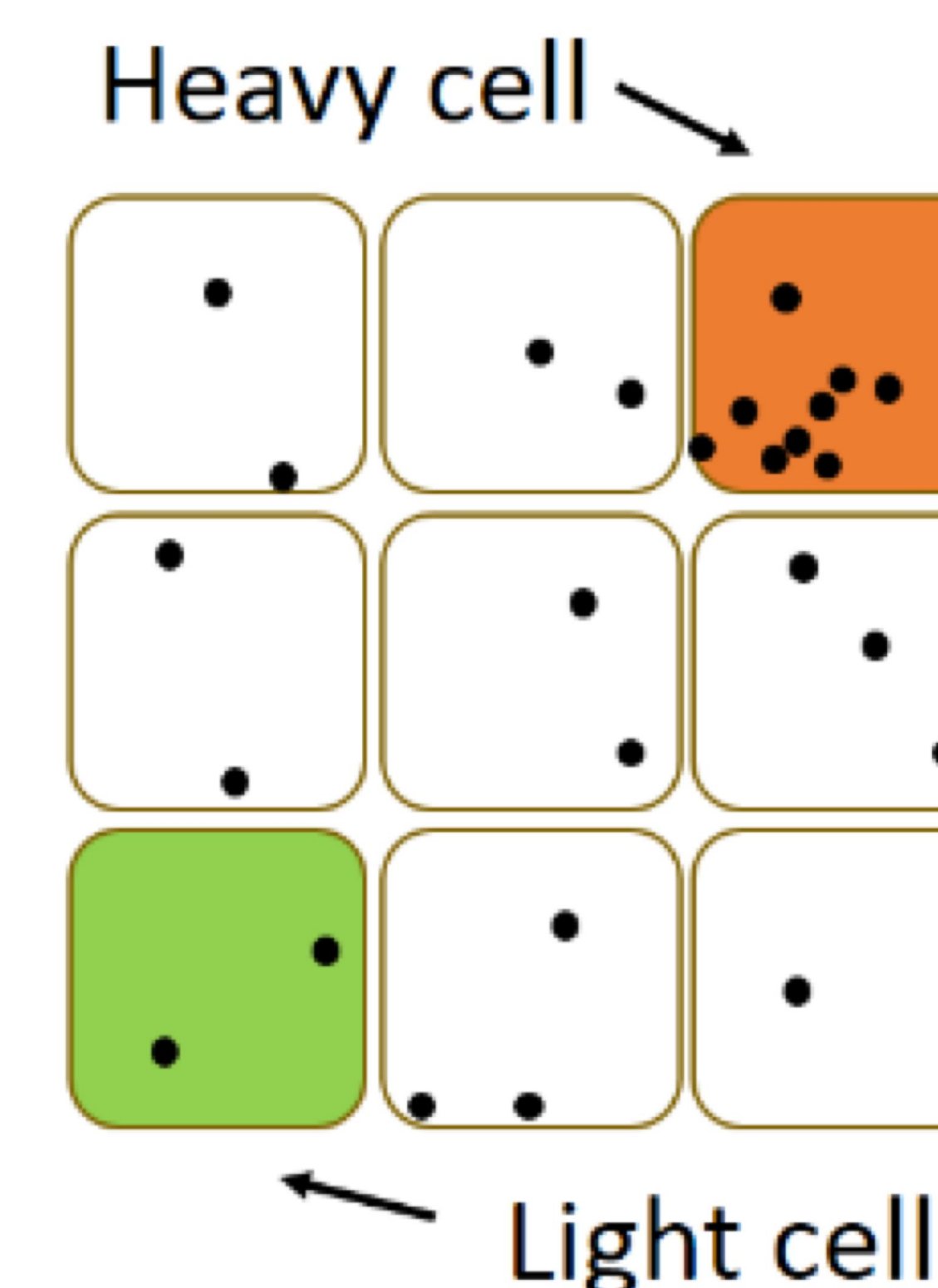
Streaming Solutions

haloes \approx heavy hitters?

Heavy hitters – most frequent items in the stream.
Naïve solution is to use 3D mesh:

1. Particle \rightarrow Cell ID.
2. Heavy cells \approx Haloes.

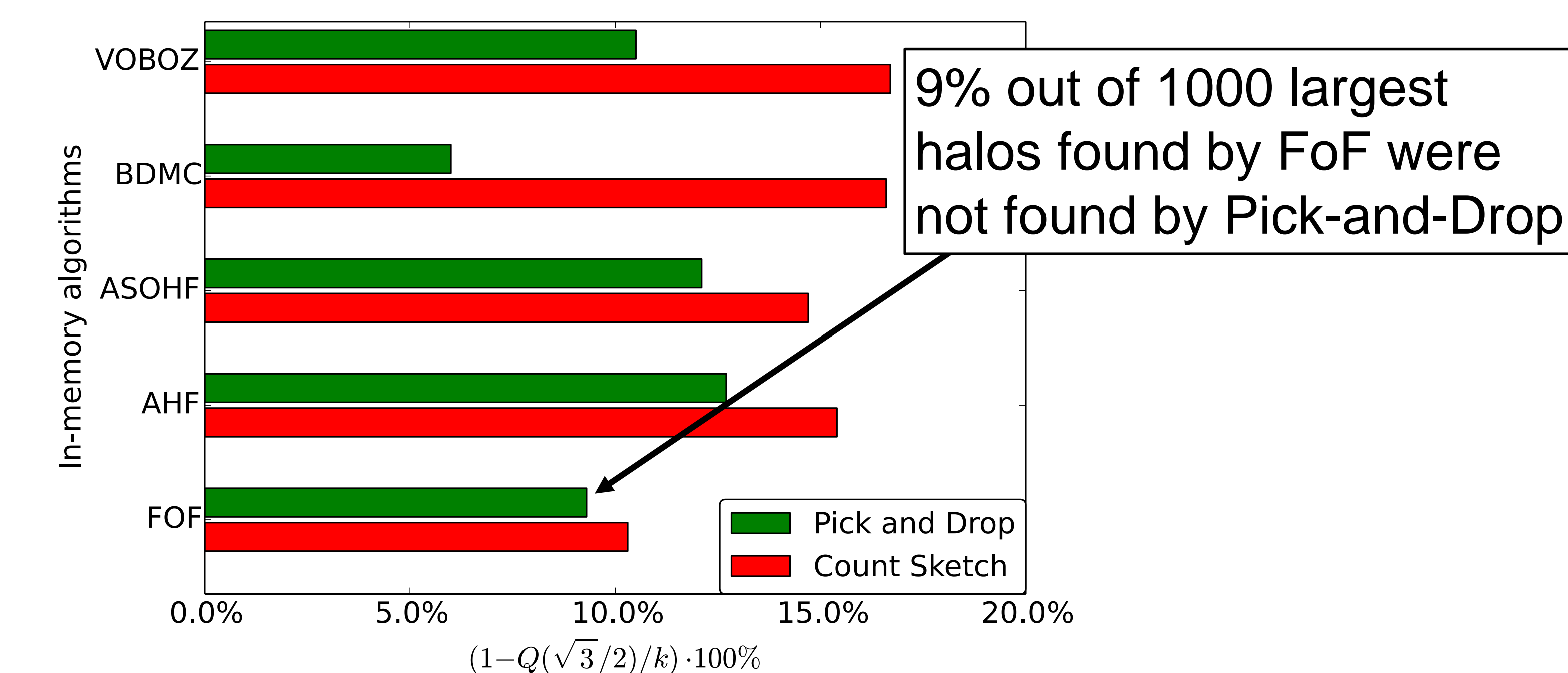
Use any heavy hitter algorithm as black box



Heavy Hitter Algorithms

Algorithms implemented:

- Count Sketch algorithm: $\sim 1\text{GB}$
- Pick-and-Drop algorithm: $\sim 30\text{MB}$
- Friends-of-Friends: $\sim 12\text{GB}$



Results

- Connection between haloes and heavy hitters.
- Two streaming algorithms for finding top- k haloes with 90% accuracy.
- Sublinear memory provides scalability.

Future Directions

- Extend result for large k .
- Consider 6-dimensional space (position + velocity).