# Parallel DBSCAN

**Bowen Chen**

# DBSCAN Overview



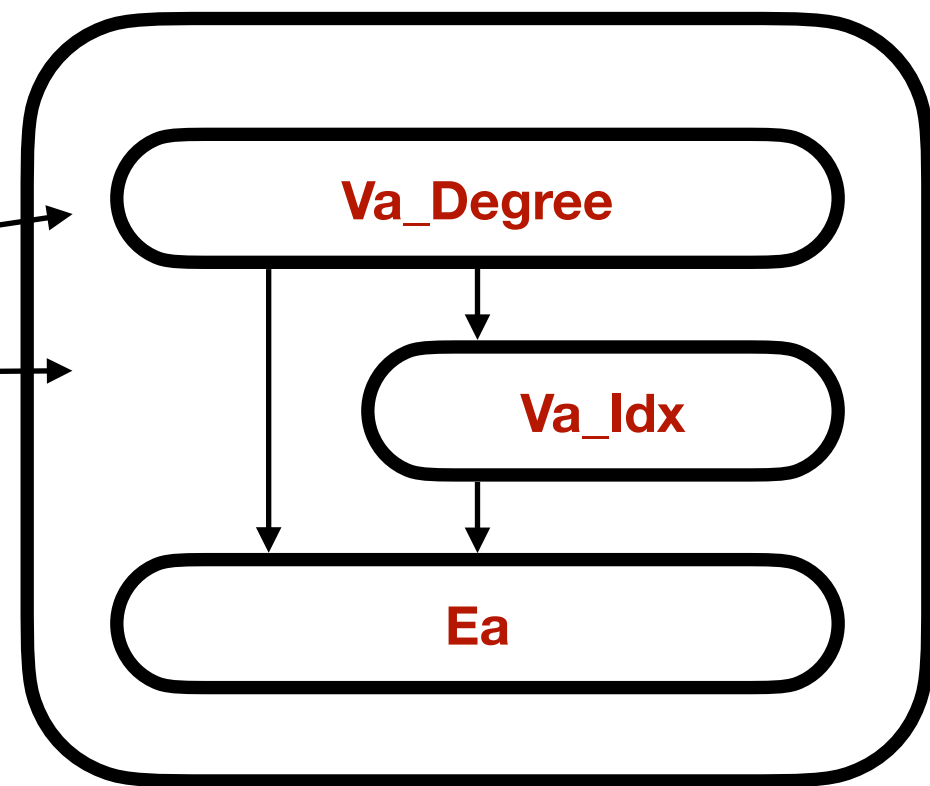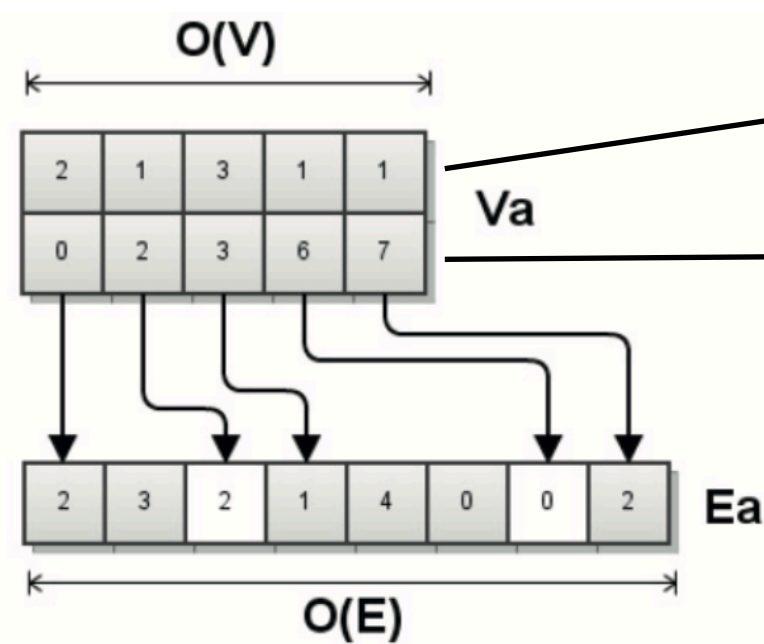min pts = 5
search radius

Cluster 1

Cluster 2
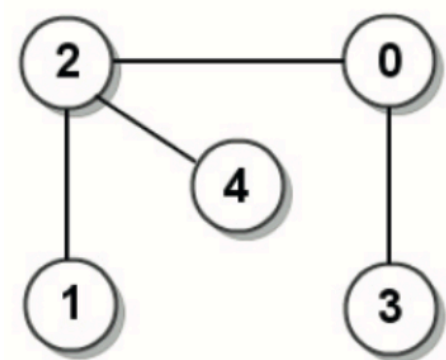
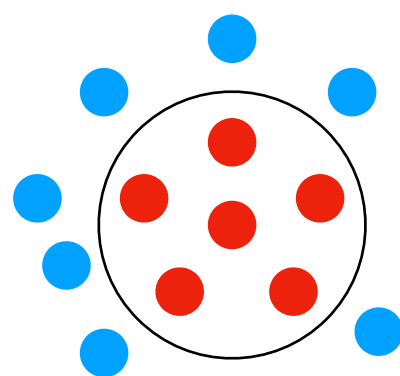Break down:   **1) find core points          2)  growing neighbours**
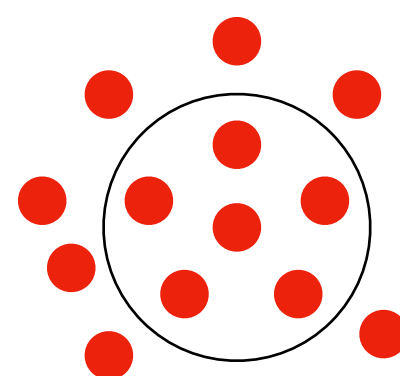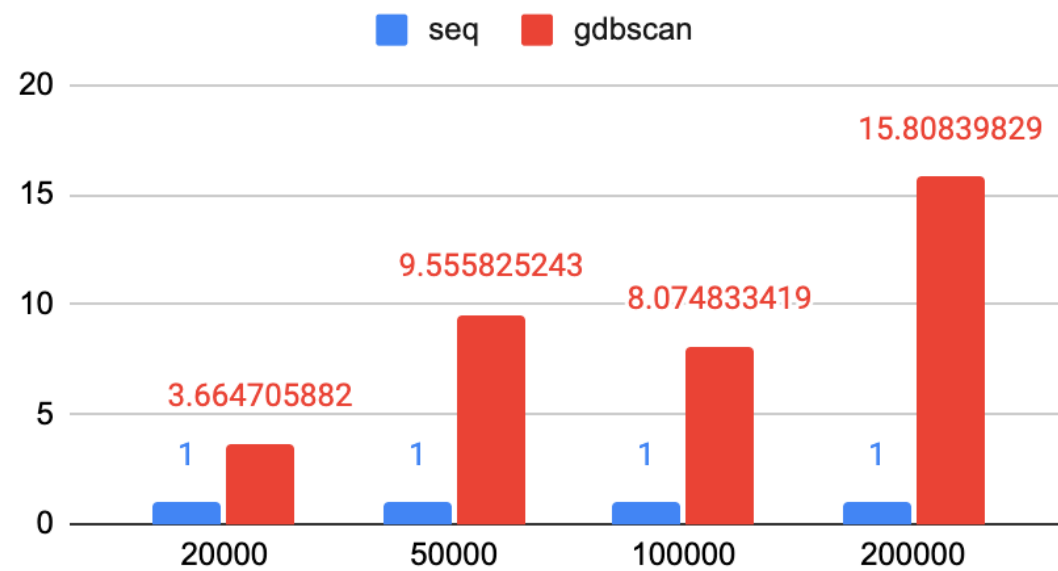
# G-DBSCAN

**Graph Construction**



**BFS SCAN**



Level 0 → Level 1 → Level 2

# G-DBSCAN

## speedup (x)

Legend: seq (blue), gdbscan (red)

- 20000: seq 1, gdbscan 3.664705882
- 50000: seq 1, gdbscan 9.555825243
- 100000: seq 1, gdbscan 8.074833419
- 200000: seq 1, gdbscan 15.80839829

## gdbscan execution break down

Legend: bfs_scan (red), graph_construct (blue)

- 20000: ~27% graph_construct, ~73% bfs_scan
- 50000: ~36% graph_construct, ~64% bfs_scan
- 100000: ~20% graph_construct, ~80% bfs_scan
- 200000: ~30% graph_construct, ~70% bfs_scan

## Pie chart

- bfsKernel 68.5%
- collectDegreeKernel 13.5%
- constructEaKernel 13.0%
- exclusive_scan 5.0%

## BFS SCAN is the bottleneck!!!

# PDSDBSCAN (Disjoint-Set)

## Iterate over points

**Iteration 3**

**Iteration 7**

**With bfs**

## Key Idea: postpone union



(a) Data points

(b) Singleton trees

(c) Intermediate trees

(d) Final trees

## Shared Memory Parallel with omp
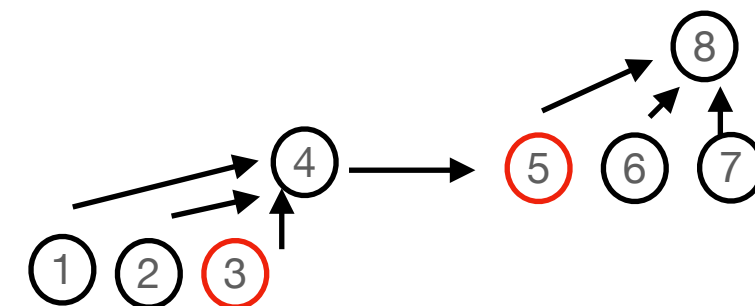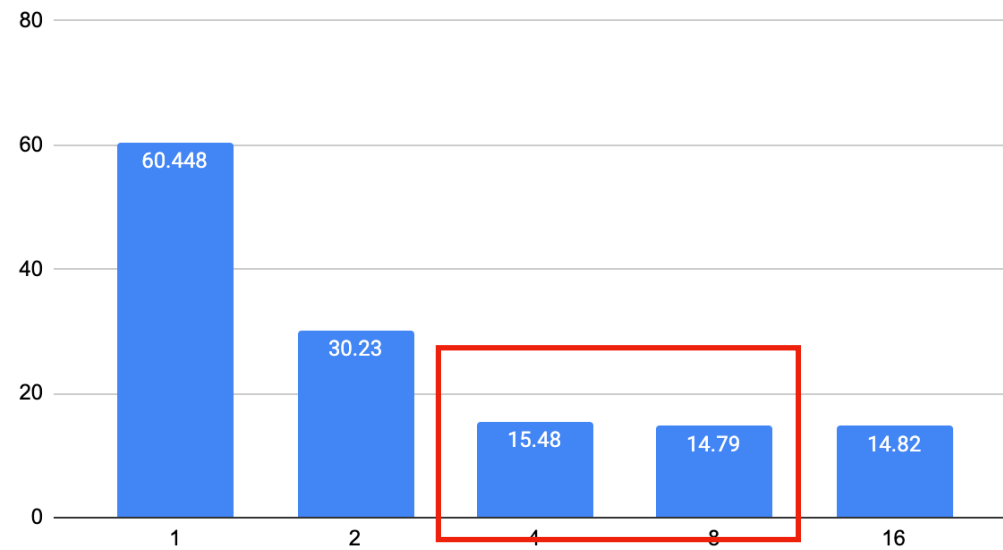
**Thread 0**    local computation    **Thread 1**

**crossThread**

**cross thread union with lock**

# PDSDBSCAN (Disjoint-Set)

**Execution time with different threads**
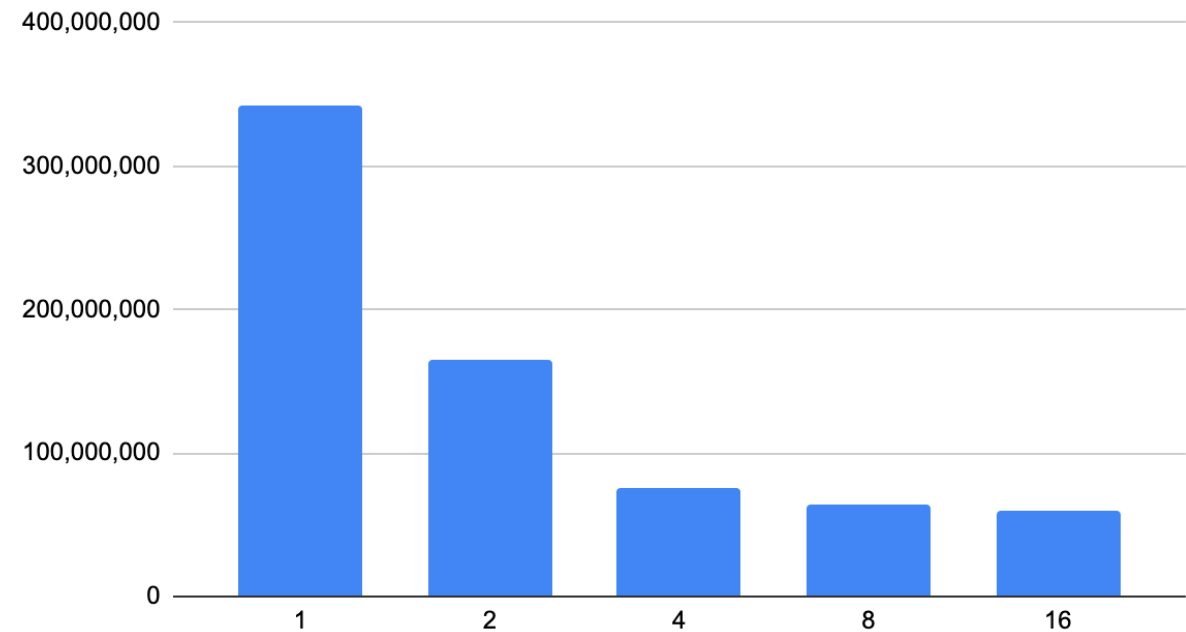
ds-shm twitter_20000 (s)



**Cache Misses**

ds-shm cache-misses



**Workload Distribution**

per thread execution time (s)



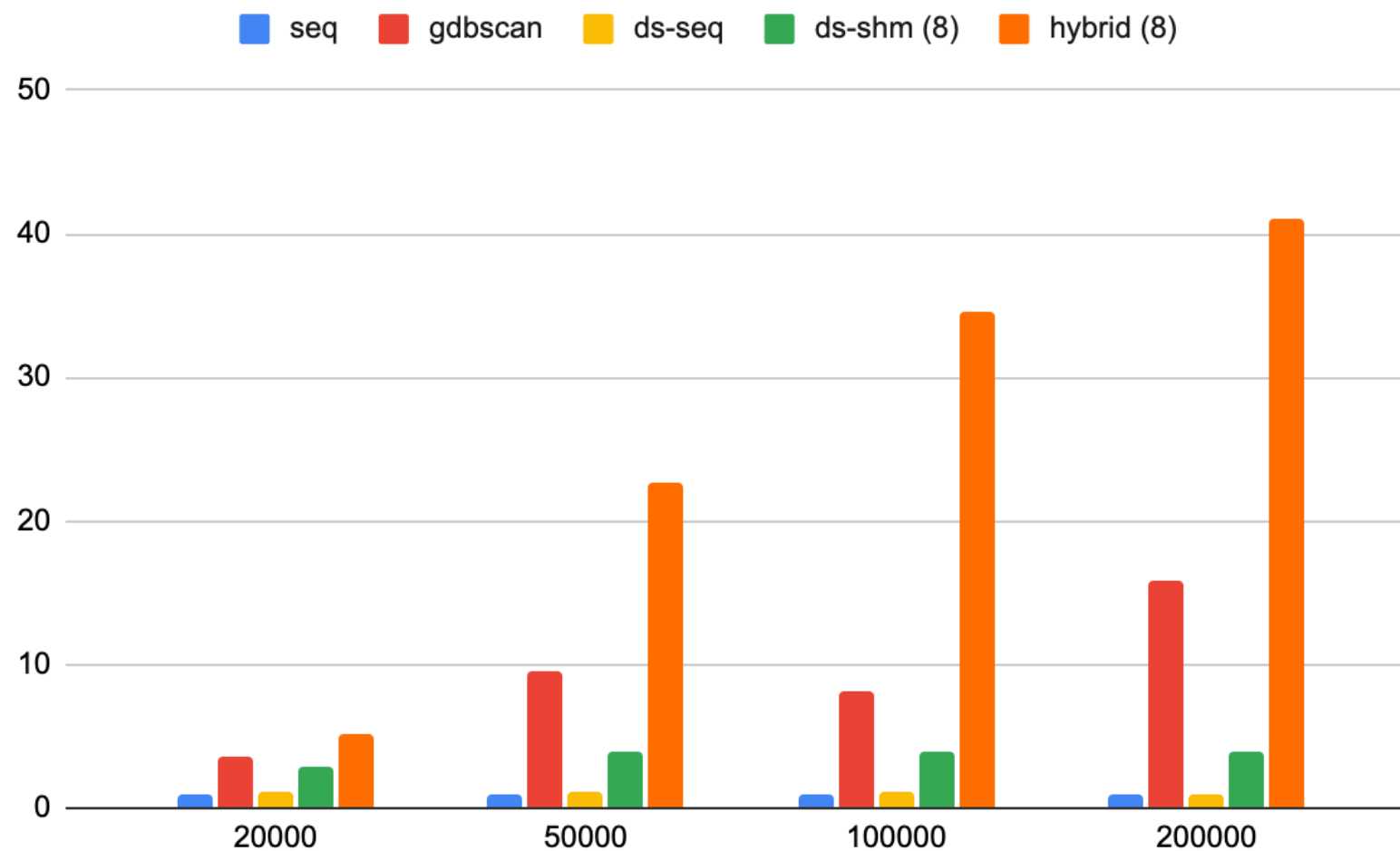**Execution Time Breakdown**

ds-shm execution time break down



**Local Computation is the bottleneck!!!**

# Hybrid

**Key Idea:** dispatch computation intensive workload to cuda,
Then leverage disjoin set for efficient merge (union)

**Graph Construction (cuda)** ⟶ **Disjoint Set (omp)**

# Correctness Validation