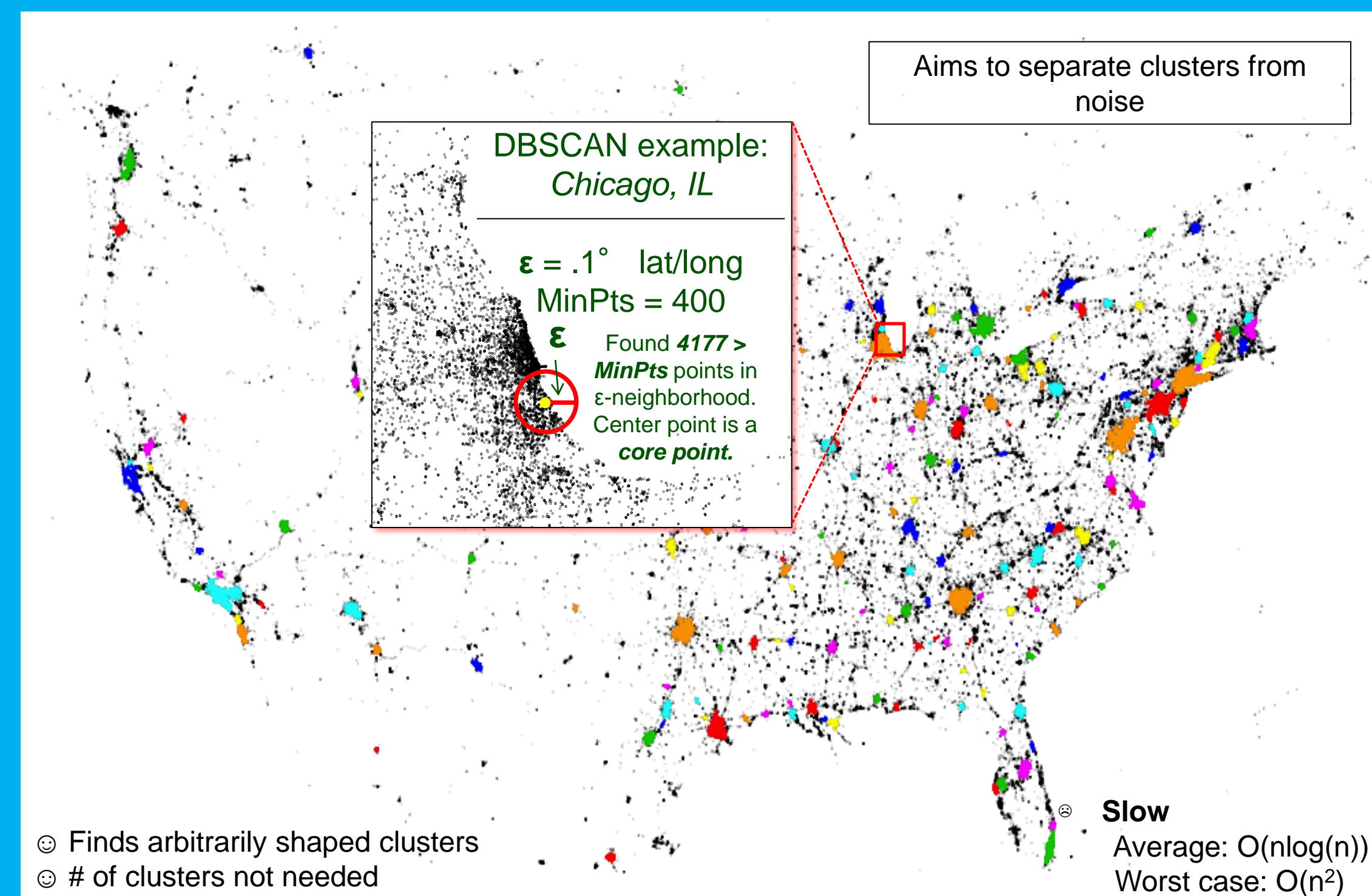


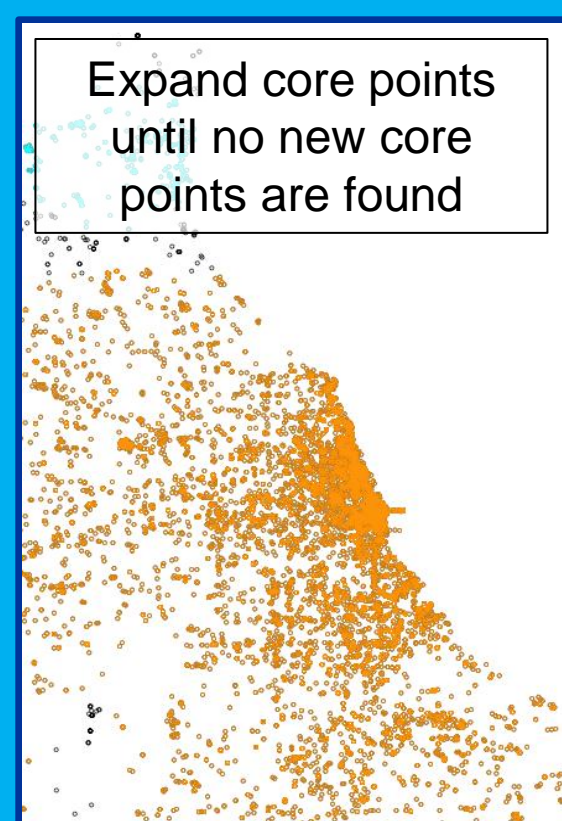
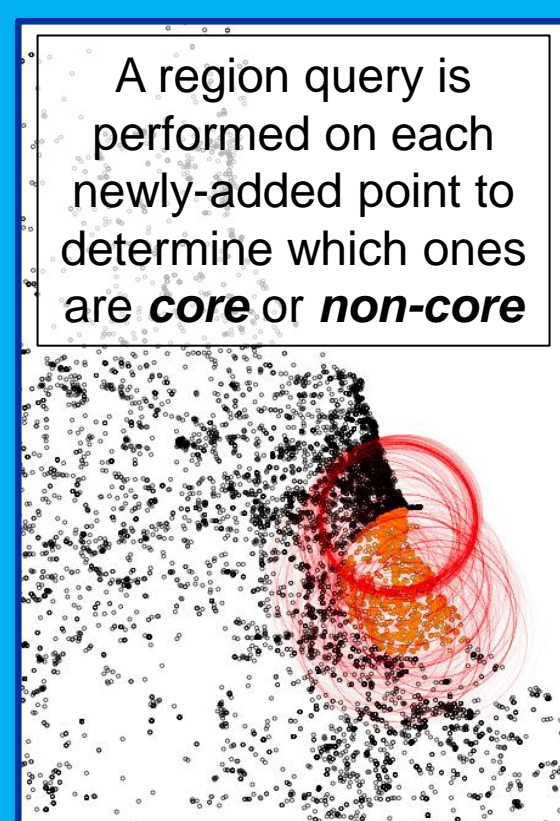
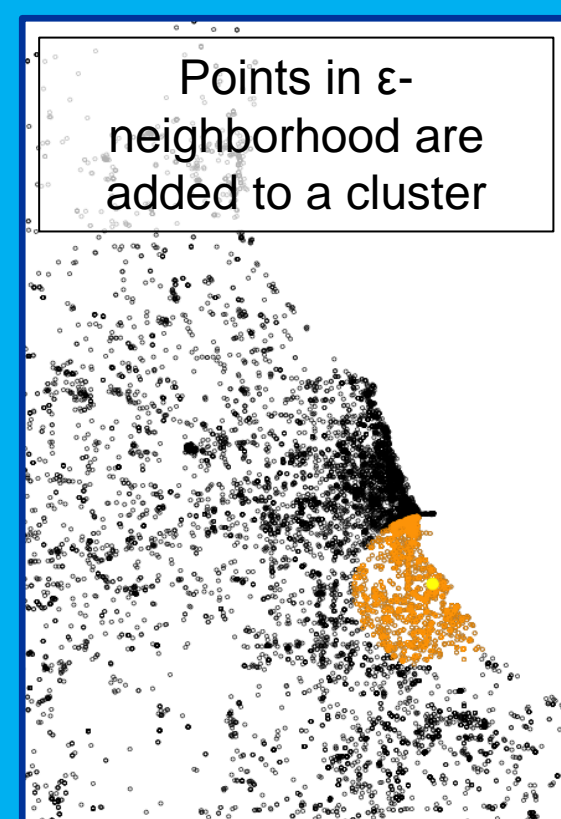
DBSCAN



Geo-located tweets from Twitter gathered on Nov. 10, 2011

DBSCAN Scaling

- DBSCAN's serial performance limits feasible calculations to tens of millions of points
- Need parallelism to cluster datasets with billions of points



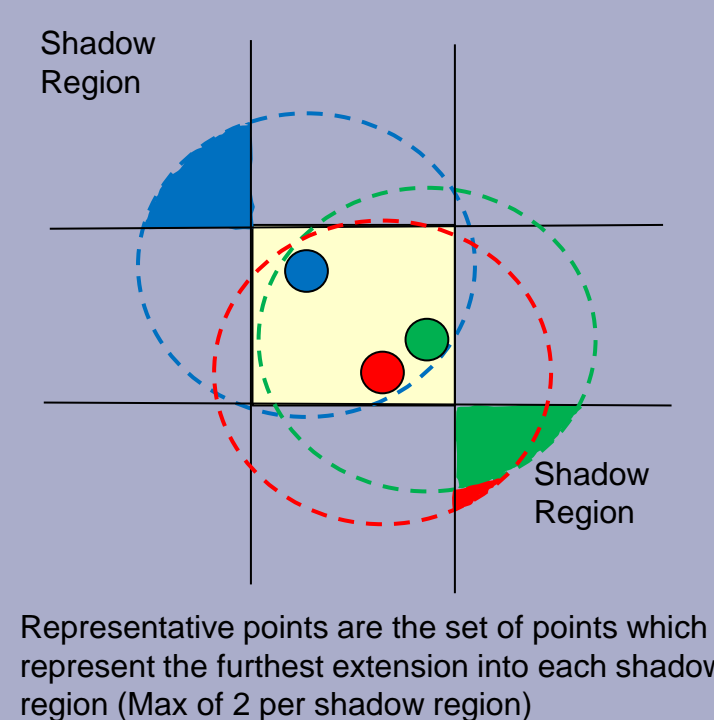
Mr. Scan Execution

3. Merge Overlapping Clusters (CP)

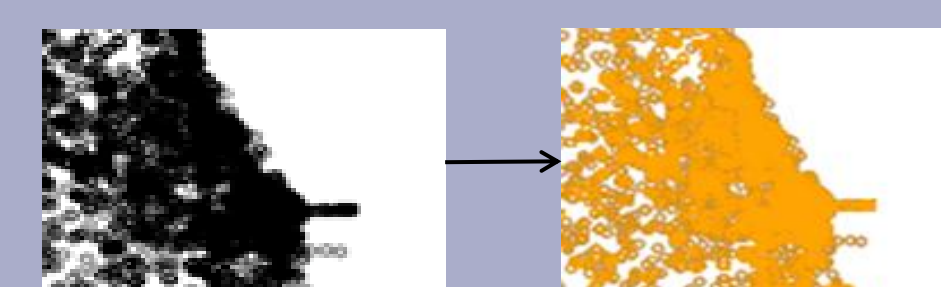
1. Detect possible cluster overlap with grid
 2. Check for representative point overlap
 3. Calculate the diff of the non-core points
 4. Calc distance from representative points. Merge if at least ONE's distance from $\epsilon < \epsilon$
- If $D < \epsilon$ merge
- Orange Non-Core
Blue Non-Core

2. Pick Representative Points (BE)

- Representative points are a small finite set of points that represent all core points in a $\epsilon \times \epsilon$ box
- Used to detect cluster overlap
- A core point that would cause a merge has to fall within ϵ of a representative point



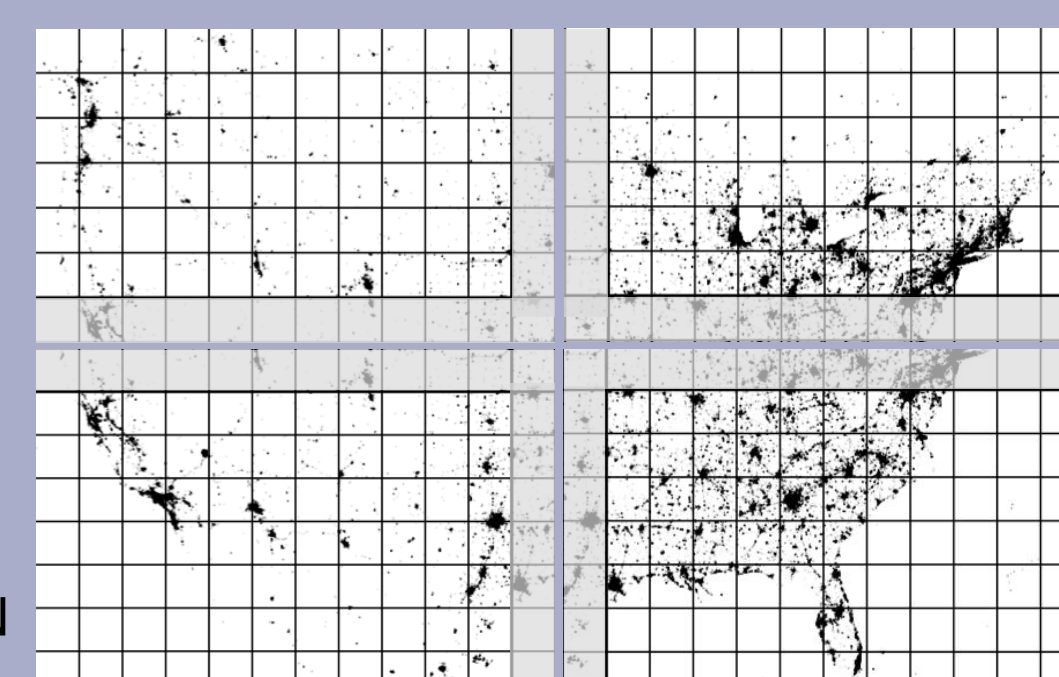
1. Run DBSCAN (BE)



- Run DBSCAN to classify points into clusters on local partition

0. Prep Input Data

- Divide data into $\epsilon \times \epsilon$ boxes
- Form roughly equal partitions from boxes
- Outline each partition with "shadow area"
- Target # points per BE:
800K – GPU DBSCAN
50K – CPU DBSCAN



4. Color the Clusters (FE)

- Color each cluster
- Send coloring to BE for cluster output

