

Conquering the Divide: Continuous Clustering of Distributed Data Streams

Graham Cormode
graham@research.att.com

S. Muthukrishnan
muthu@cs.rutgers.edu

Wei Zhuang
weiz@cs.rutgers.edu

Outline

- **Motivation in Data Streams**
- Continuous Distributed Model
- Our CD Clustering Algorithms
- Experiments
- Extensions and conclusion

Data Streams

- Modern data acquisition generates large distributed data coming as distributed streams
 - Environmental sensors
 - Web/blog crawlers
 - Soldier locations on battlefield
 - Different mail servers collate their mails over time
 - Repositories on book, audio, video evolve on their own
- Need support not provided in traditional DBMS

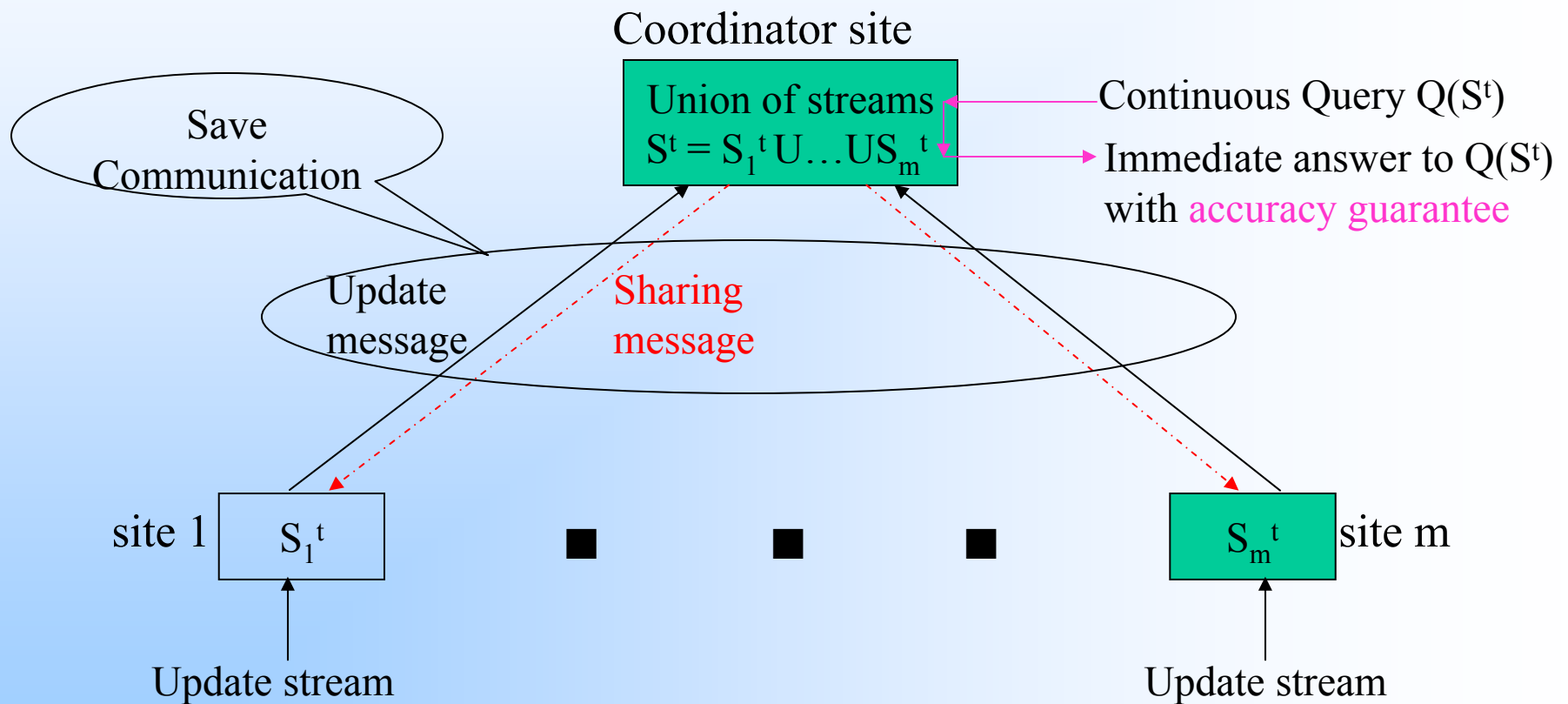
Monitoring Challenges

- Need real-time monitoring of data properties
 - To help identify any problem of the data (collection): network attack detection, performance debugging
 - To help understand the nature of the data distribution: classification, distinct count, etc.
- Clustering is a foundational monitoring problem
 - Web crawlers cluster pages into groups to monitor traffic load
 - Underwater fish tracking monitors the clusters of schools to quickly deploy attracting or dispelling devices

Outline

- Motivation
- **Continuous Distributed Model**
- Our CD Clustering Algorithms
- Experiments
- Extensions and conclusion

Continuous Distributed (CD) Model

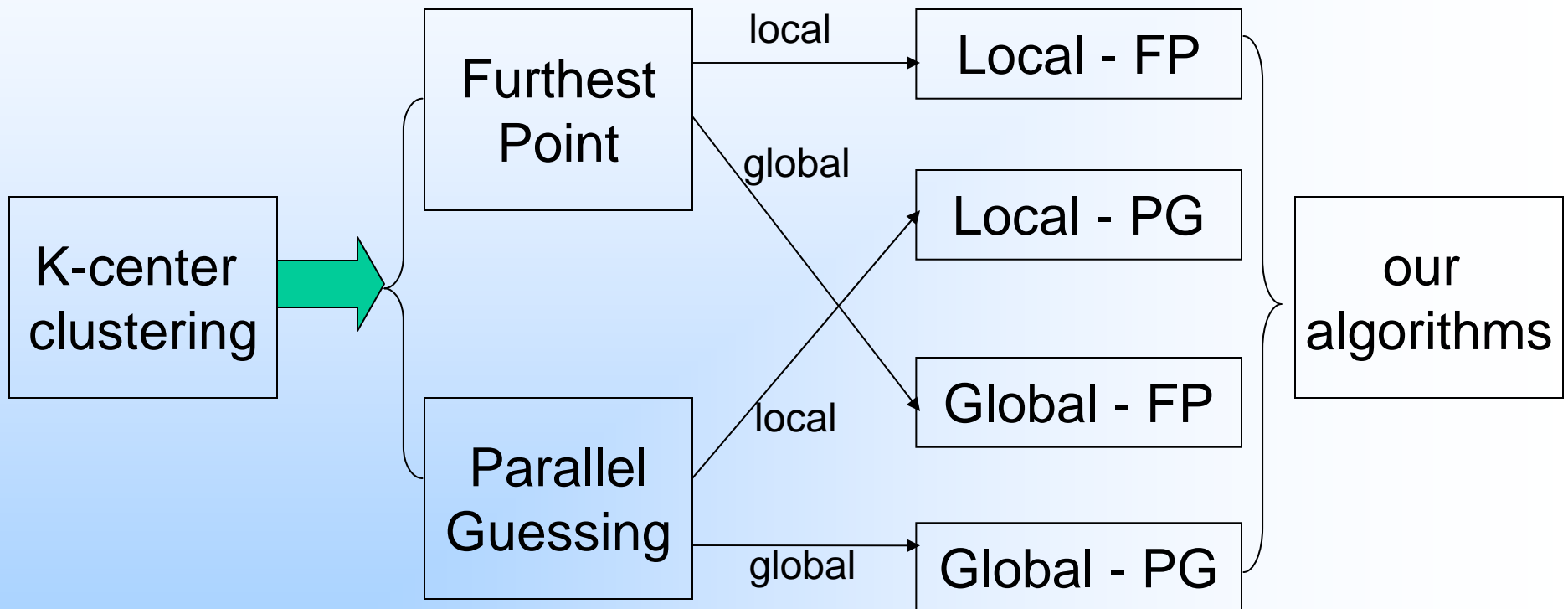


Outline

- Motivation
- Continuous Distributed Model
- **Our CD Clustering Algorithms**
 - Local algorithms
 - Global algorithms
- Experiments
- Extensions and conclusion

Our Contribution

- We provide several k-center clustering algorithms that
 - Monitor **distributed dynamic** streams
 - Answer **continuous** query
 - Achieve **constant factor** approximation as good as centralized algorithm
 - **Minimize the communication** between sites



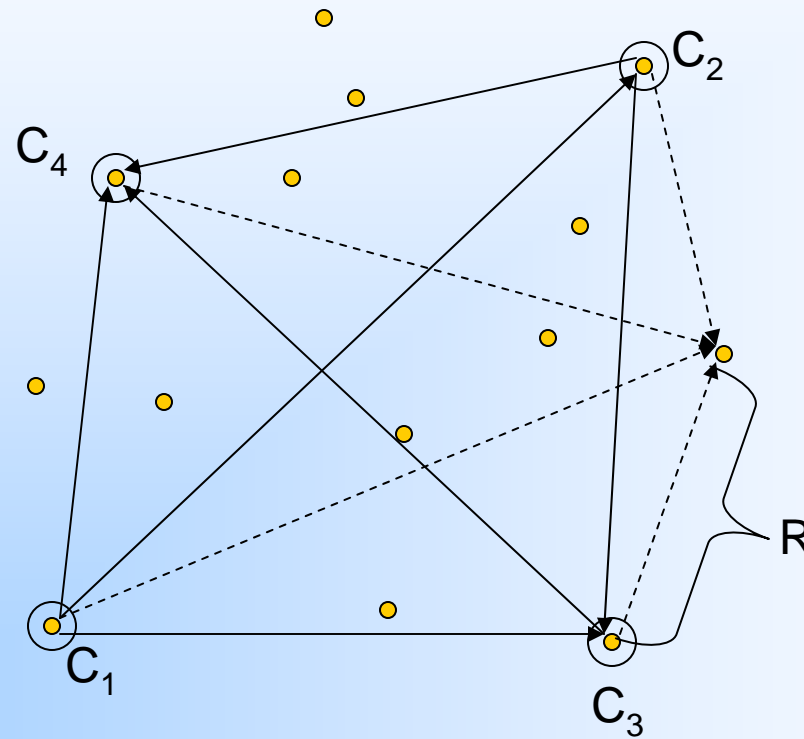
K-center clustering

- Given
 - a set P of n points
 - a distance function $d(p, q)$ between pairs of points
- Find a set $C \subset P$ of $k < n$ points as centers
- Minimize the radius $r = \max_{p \in P} d(p, C(p))$
 - $C(p) = \operatorname{argmin}_{c \in C} d(c, p)$, distance of p to closest center

“Furthest Point” (FP) Algorithm

- Furthest Point (FP) (Gonzalez '85)
 - Pick an arbitrary point as the first center, $C_1 = \{c_1\}$
 - Given i centers in C_i , find the point p to maximize $d(p, C_i(p))$, set $C_{i+1} = C_i \cup \{p\}$.
 - Stop after k iterations with k centers C_k .
- Analysis:
 - Factor 2 approximation to the optimal clustering
 - k passes, $O(kn)$ distance computations

An example



- A 4-center clustering using FP algorithm

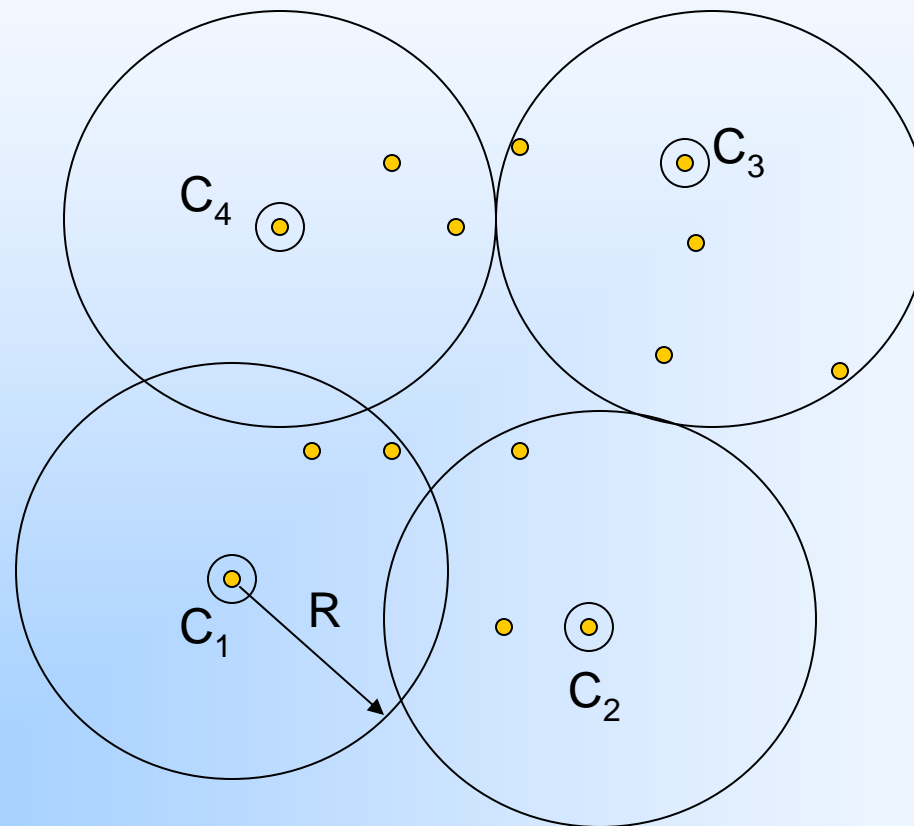
“Parallel Guessing” (PG) Algorithm

- Assume R is a lower bound on the optimal radius
- Pick an arbitrary first center, $C = \{c_1\}$
- For each point p , if $d(p, C(p)) > R$, then $C = C \cup \{p\}$
- In parallel, guess $R = (1 + \epsilon/2), (1 + \epsilon/2)^2, (1 + \epsilon/2)^3, \dots$
- Use result for smallest R that generates k or fewer centers

Analysis of PG Algorithm

- Makes one pass over the input
- Gives $(2 + \epsilon)$ approximation
- If the ratio $\max_{p,q} d(p,q) / \min_{p,q} d(p,q) = \Delta$, the number of parallel guesses is $O(1/\epsilon \log \Delta)$
- Stores $O(k/\epsilon \log \Delta)$ points, and each update needs $O(k/\epsilon \log \Delta)$ distance computations

An example



- A 4-center clustering using PG algorithm

Outline

- Motivation
- Continuous Distributed Model
- **Our CD Clustering Algorithms**
 - Local algorithms
 - Global algorithms
- Experiments
- Extensions and conclusion

Merging Local Clusterings

- **Local-FP**: each site runs FP on its points, send k-centers to coordinator on recluster.
- **Local-PG**: each site runs PG on its points, shares centers for current good guess R with coordinator
- Both have variants which share more or less data:
 - Coordinator does not share cluster information
 - Coordinator broadcasts the global cluster radius
 - Coordinator lazily shares radius (when contacted by site)
- All versions guarantee a $(4+\epsilon)$ approximation

Space and Communication Cost

— Local Algorithms

- **Local-FP**: each local reclustering sends $O(k)$ points
 - No tight bounds on worst case number of reclusterings.
- **Local-PG**: one pass over the input, each site sends up to k points for each guess.
 - Total communication is at most $O(km/\epsilon \log \Delta)$.
 - Space required at each site is $O(k/\epsilon \log \Delta)$.

Cluster Merging

- Clusters can be merged by clustering the cluster centers
- Given point set P_i , $i=1,\dots,m$, let $C_i \subseteq P_i$ be the k centers of α -approximate k -center clustering of P_i
- Run a β -approximate k -center clustering on union of C_i 's
- **Result:** the resulting k centers C form an $(\alpha+\beta)$ clustering of the point set $P = P_1 \cup P_2 \cup \dots \cup P_m$

Outline

- Motivation
- Continuous Distributed Model
- **Our CD Clustering Algorithms**
 - Local algorithms
 - **Global algorithms**
- Experiments
- Extensions and conclusion

Global-FP Clustering

- Distributed **collaboration** to find k global centers using FP:
 - In round i , each site sends its furthest point from current centers C_i
 - Coordinator picks furthest point as a center and sends it to all sites
- Each remote client **monitors** its local points, recluster if some point is more than $(1 + \epsilon/2)R$ from the global centers
- Alternate **collaborating** and **monitoring** phrases
- Each reclustering sends $O(mk)$ points

Global-PG Clustering

- Coordinator maintains a global PG clustering, shares the current good guess and centers with all clients
 - Client sends new point p if $\exists i, d(p, C_i(p)) > R_i$
 - Coordinator shares new centers with all clients
- Total communication bounded by $O(km/\epsilon \log \Delta)$
- Coordinator always has a $(2+\epsilon)$ -approximation

Comparison

- A summary of the algorithms:

Algorithm	Guarantee	Space	Communication
Local-FP	$4 + \epsilon$	$O(n)$	—
Local-PG	$4 + \epsilon$	$O(k/\epsilon \log \Delta)$	$O(km/\epsilon \log \Delta)$
Global-FP	$2 + \epsilon$	$O(n)$	—
Global-PG	$2 + \epsilon$	$O(k/\epsilon \log \Delta)$	$O(km/\epsilon \log \Delta)$

Outline

- Motivation
- Continuous Distributed Model
- Our CD Clustering Algorithms
- **Experiments**
- Extensions and conclusion

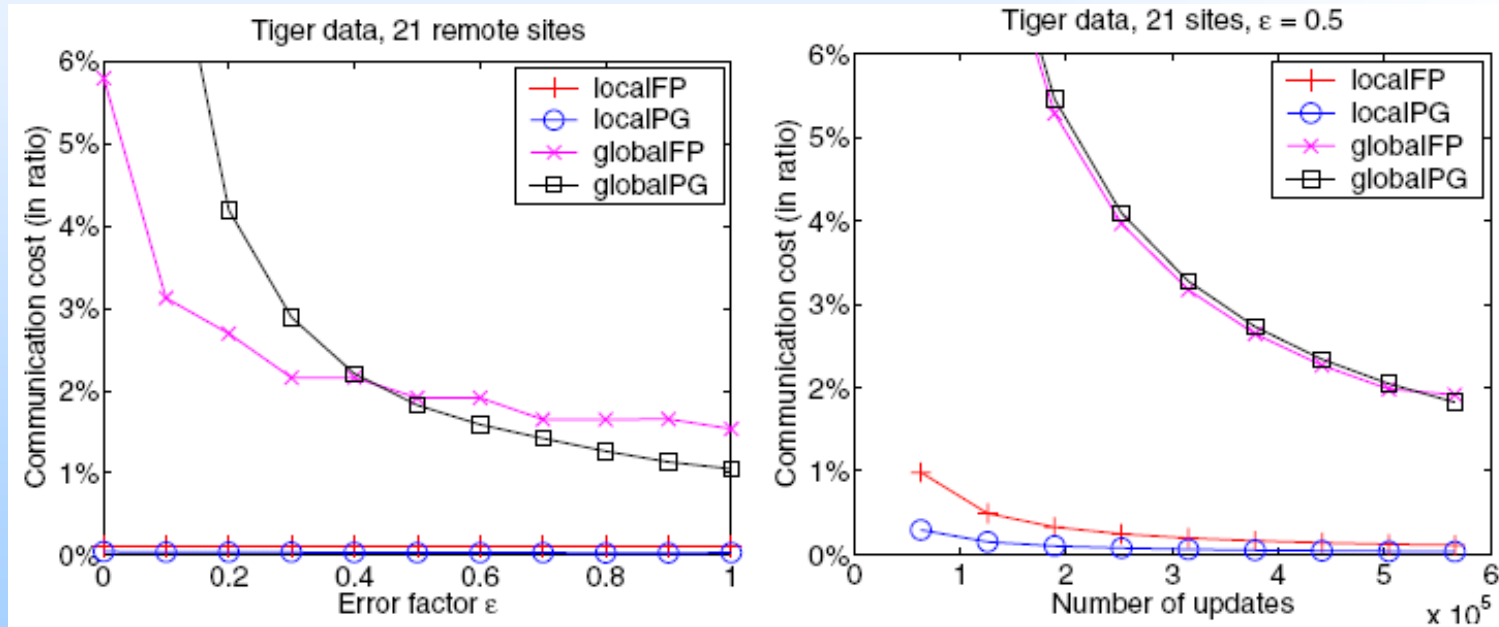
Experimental Data Sets

- **US Census 2005 (TIGER) data**
 - $N = 600k$ 2-D points
 - State of New Jersey, $m = 21$
- **Stock price series**
 - Price series of a stock with $n = 330k$ values
 - Simulate 100-D vectors
- Measure communication cost as ratio compared to sending complete data

Experimental Results

Communication cost

- Fix sites m and clusters k , vary ϵ : compare communication

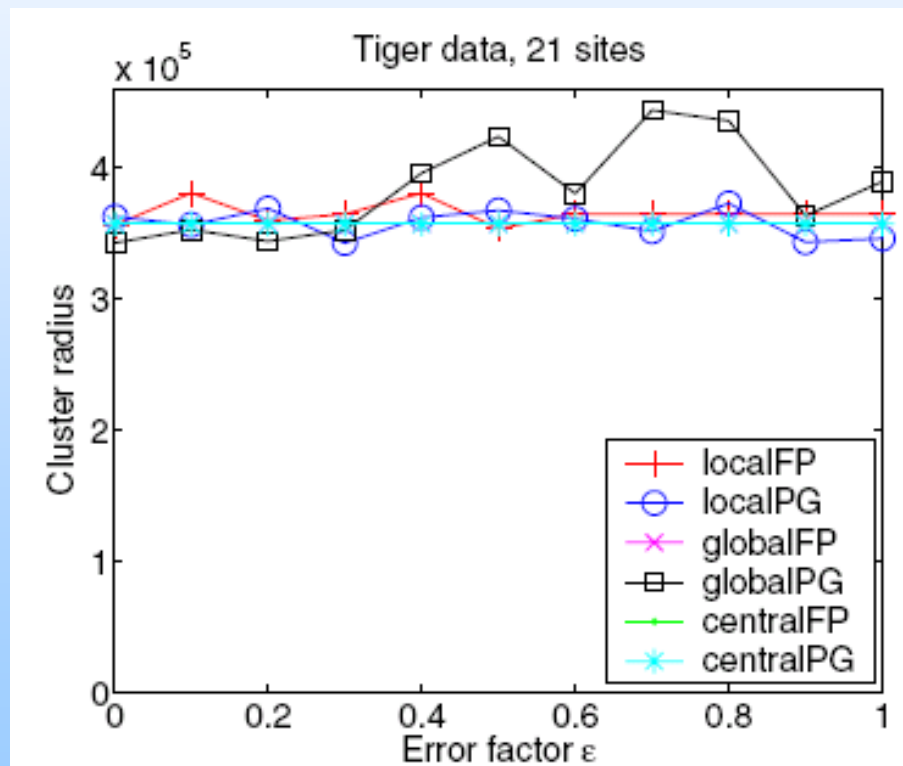


- The local algorithms are more cost effective than naive global implementations of distributed algorithms

Experimental Results —

Cluster Quality

- Fix m , k , vary ϵ , compare radius: smaller is better

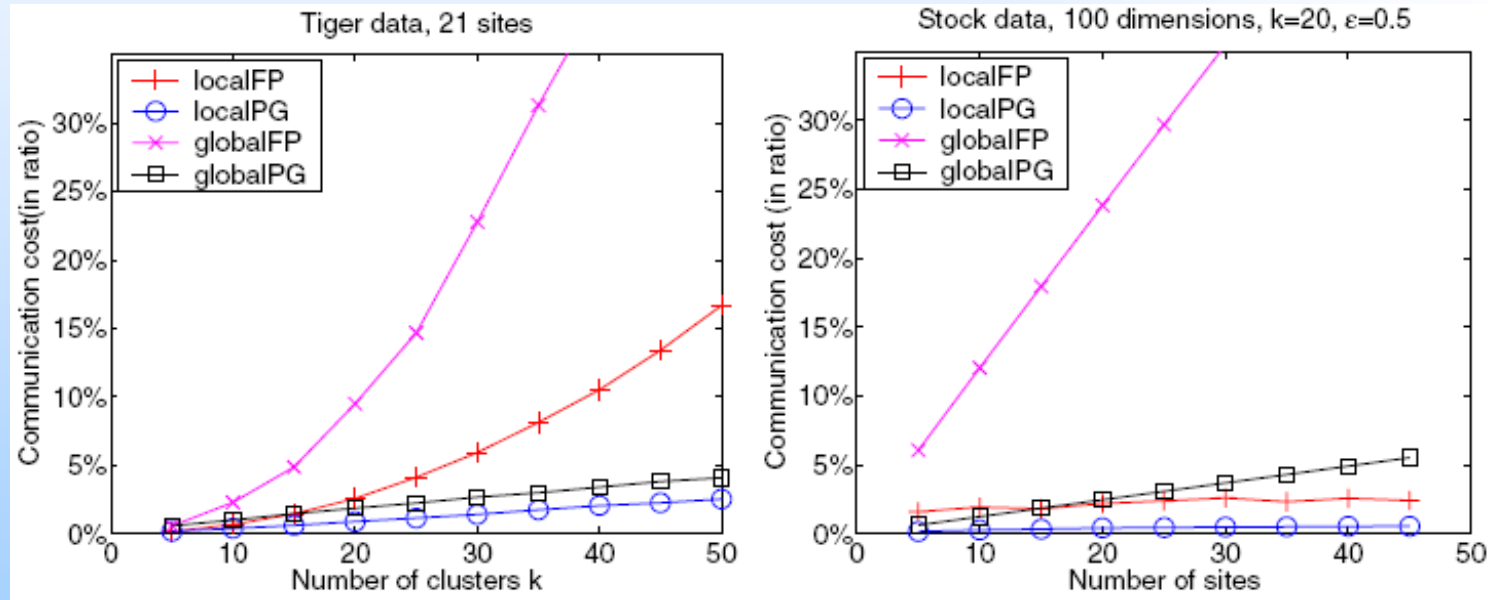


- Not much to choose between algorithms in terms of accuracy
- Choice of ϵ seems not very sensitive – can choose a large value and not suffer

Experimental Results

Other Dependence

- Varying number of clusters k and sites m



- Analysis of cost seems a good predictor: all algorithms very close to linear on m , quadratic on k for FP algorithms 28

Outline

- Motivation
- Continuous Distributed Model
- Related Work
- Centralized K-center Clustering Algorithms
- Our CD Clustering Algorithms
- Experiments
- **Extensions and conclusion**

Extensions

- **Deletions**: run same algorithms as before. Monitor the $(k+1)$ points which are ‘proof’ of current clustering and recluster if any get deleted.
- **Moving Points**: can treat as deletion of old point, reinsertion at new location if possible.
- **Sliding window**: treat the point that is no longer in the window as a deletion of point

Variable number of clusters

- Although we built a k -center clustering, we can get any k' -clustering for $k' \leq k$
 - For FP algorithms, pick the first k' centers
 - For PG algorithms, the k' clustering is the smallest guess that generates at most k' centers
- 1-center problem is a special case, and we can get a stronger $(1+\epsilon)$ -approximation

Conclusion

- Introduced **distributed continuous clustering**
- Shown local and global k-center algorithms
 - Minimize communication, guarantee clustering quality
- Experimental results show
 - local and parallel guessing uses least communication
 - clustering quality is similar to centralized algorithm
 - space needed independent of the input size
- Open to study other popular clustering algorithms, and other settings

- Thank you!