# SocialTrove: A Self-summarizing Storage Service for Social Sensing

Tanvir Amin, Shen Li, Muntasir Rahman, Panindra Seetharamu, Shiguang Wang, Tarek Abdelzaher, Indranil Gupta
University of Illinois at Urbana Champaign

Mudhakar Srivatsa, Raghu Ganti
IBM Research

Reaz Ahmed
University of Waterloo

Hieu Le
TCG Corp

**ICAC 2015**

1

# Social Sensing



Egypt Unrest



Fukushima Disaster



Sandy Gas Outage



Crimea Annexation



Syria Chemical Attack



Twitter   Facebook   Google+   Instagram   Flickr

Credibility Estimation   Anomaly Detection   Timeline Reconstruction   …

# 60 Seconds of Social Media*

- 2.66 Million Google searches
- 433 Thousand Tweets
- 67 Thousand Instagram photo uploads
- 293 Thousand Facebook Status updates
- 277 Thousand Snaps over SnapChat
- 20 Million Photo views on Flickr
- 11 Thousand LinkedIn searches
- 17 Thousand Walmart transactions
- 20 Thousand Tumblr new photos
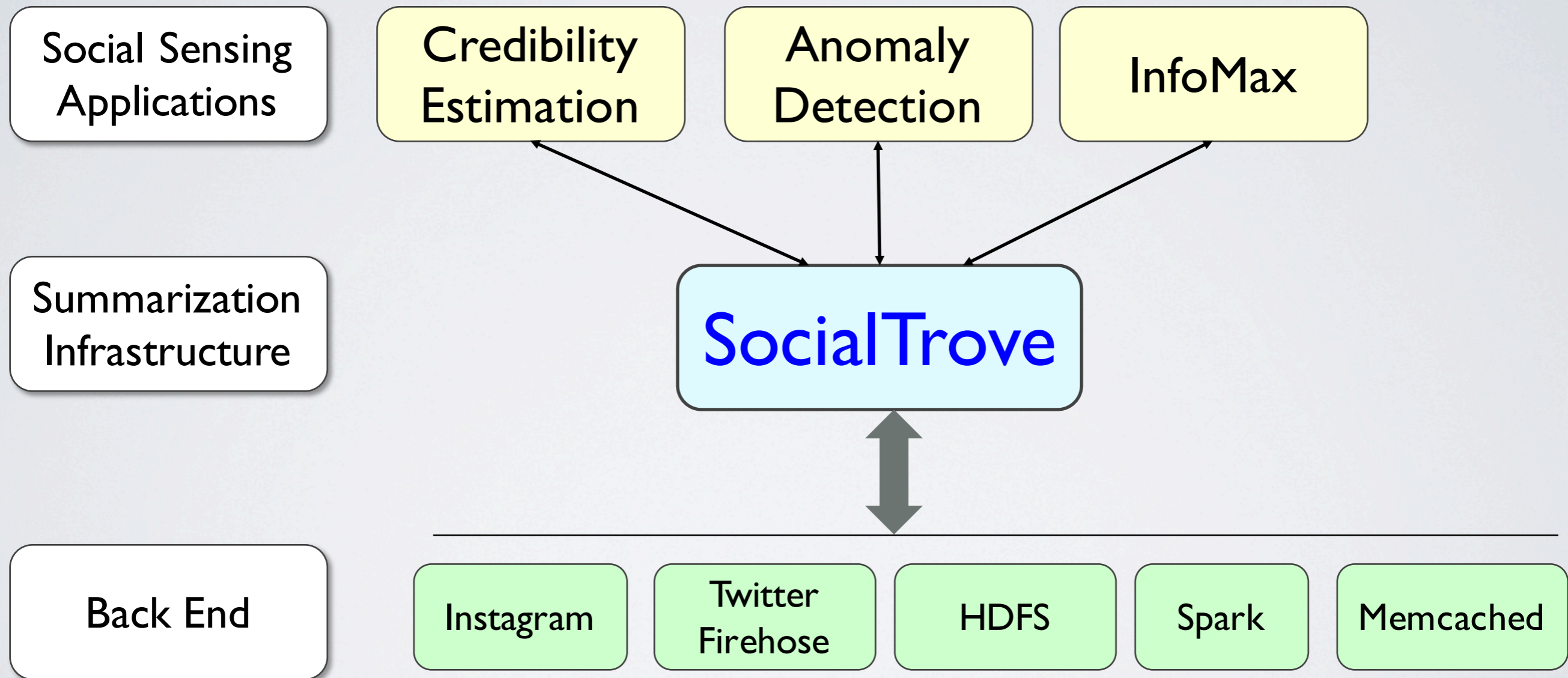- 139 Million emails sent

# Observations and Goals

- Social Sensing Applications
  - User facing – real-time updates
  - Many users – many queries
- Social Sensing Data Streams
  - Real-time *information overload*
  - Large amounts of *redundant data*
- Obtain a *representative sampling*
  - in a *content agnostic* manner
  - based on *application specified distance metric*
  - at a *configurable granularity*
  - with high query *throughput*
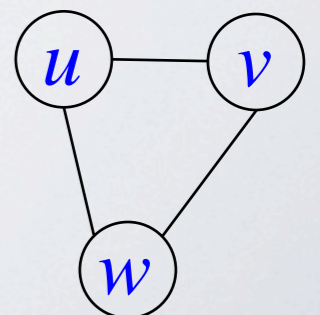  - and low access *latency*

4

# Real-time Social Sensing Stack



Social Sensing Applications

Credibility Estimation

Anomaly Detection

InfoMax

Summarization Infrastructure

SocialTrove

Back End

Instagram

Twitter Firehose

HDFS

Spark

Memcached

# Customization Interface

- SocialTrove agnostic to stream content
- Application provides two callbacks
- **`Vectorize(`** $u$ **`)`**
  - Convert stream content $u$ to a feature vector
- **`Distance(`** $u$ **`.FeatureVector,`**
  **$v$`.FeatureVector)`**
  - Compute distance between two feature vectors
  - Should be a metric, obey Triangle Inequality

$$\text{distance}(u,v) + \text{distance}(v,w) \geq \text{distance}(u,w)$$
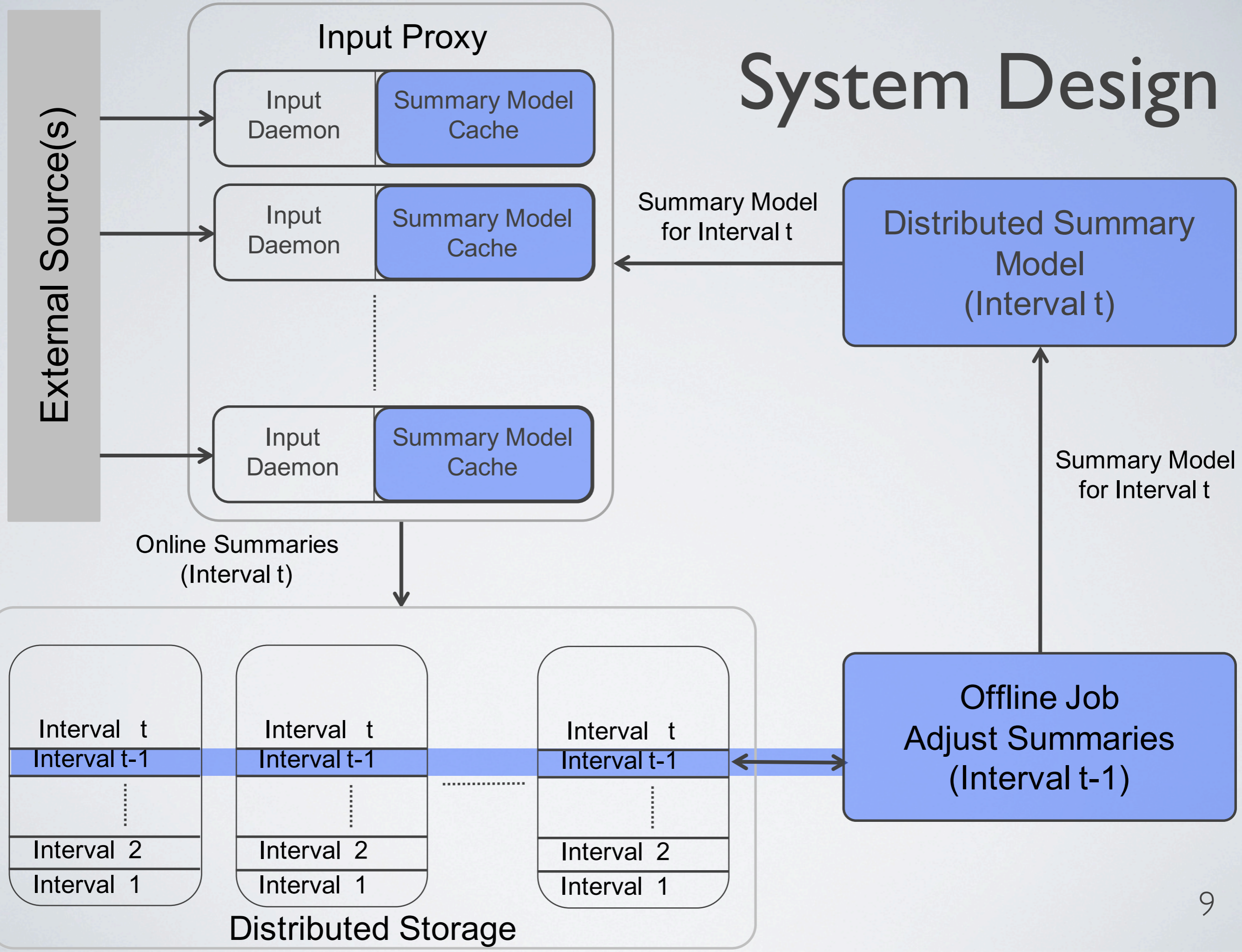
6

# Summary Query Interface

- Transform application specific query.
  - Degree of granularity for summary retrieval.
- Transform SocialTrove query response to application specific results.
- Example:
  - Rank the matching summaries and retrieve top 10.
  - Retrieve matching summaries as a list.
  - Retrieve matching objects and additional metadata.
    - Source ID, Source relations etc.
  - Retrieve summary clusters and member counts.

# SocialTrove Runtime

- Maintains a *summary model* in *memory*
  - Nearest neighbor data structure based on application defined distance metric.
  - Encapsulates distance information.
  - Hierarchical organization summarizes data at arbitrary granularity.
  - In-memory model improves throughput and latency.
- Updates limited to once per *batching interval*
  - Local insertions without synchronizing model.
  - Improves throughput and latency.
  - Temporal locality for social sensing content.
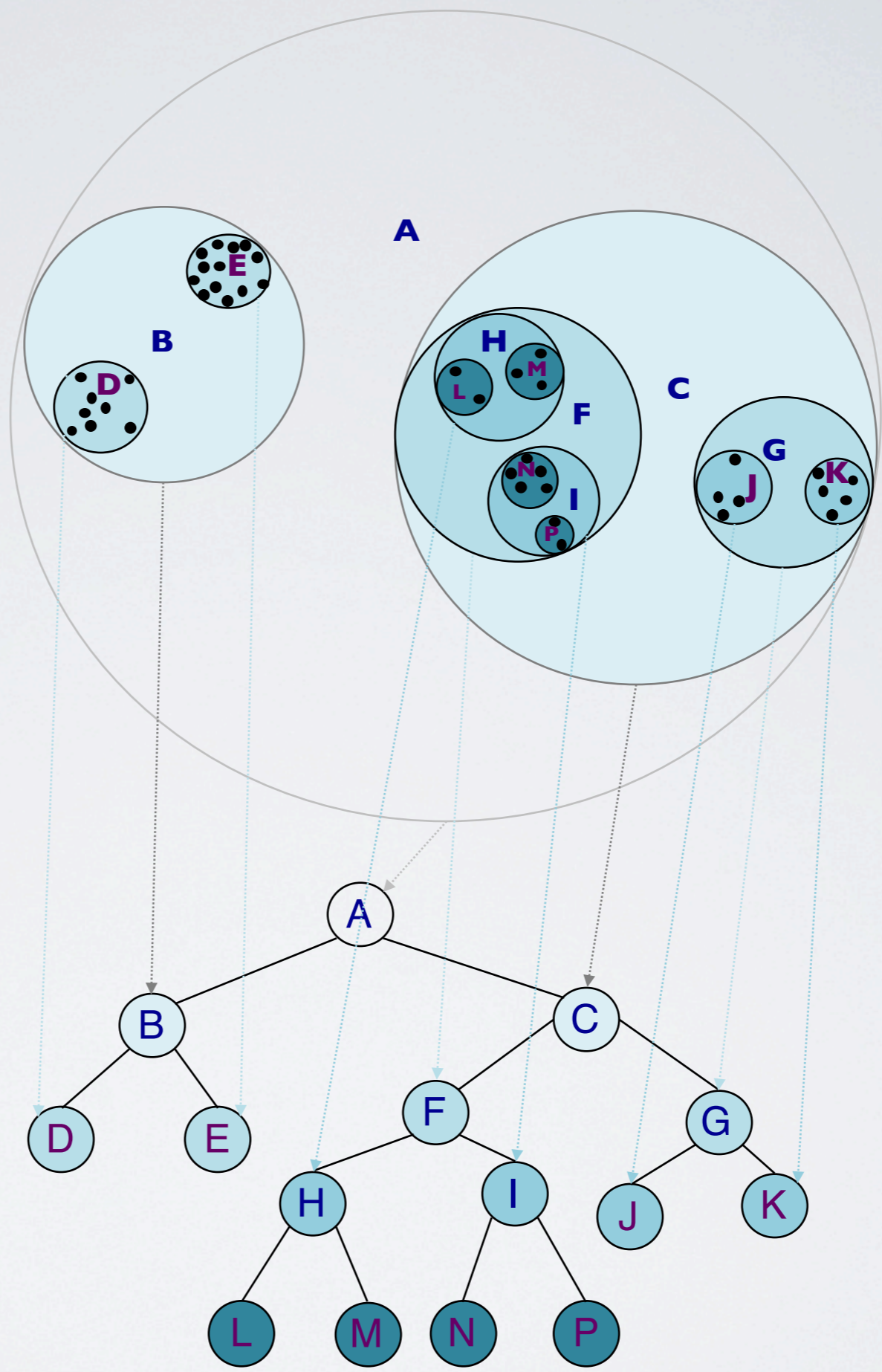  - Eventual consistency through gossip.

# System Design



External Source(s)

## Input Proxy

| Input Daemon | Summary Model Cache |

| Input Daemon | Summary Model Cache |

| Input Daemon | Summary Model Cache |

Summary Model for Interval t

Distributed Summary Model
(Interval t)

Online Summaries
(Interval t)

Summary Model for Interval t

Offline Job
Adjust Summaries
(Interval t-1)

| Interval  t |
| Interval t-1 |
| Interval  2 |
| Interval  1 |

| Interval  t |
| Interval t-1 |
| Interval  2 |
| Interval  1 |

| Interval  t |
| Interval t-1 |
| Interval  2 |
| Interval  1 |

## Distributed Storage

9

# Summary Model

- Represented as a binary tree
  - Constructed using a divide and conquer paradigm.
- Current set of vectors partitioned
  - Into two disjoint sets using *2-means clustering.*
  - Centroids of the two sets become *two children.*
    - Each set contains points nearest to the centroid of that set, compared to the centroid of the other set.
  - Both sets scheduled for further division in *parallel.*
  - Distance function must satisfy *Triangle Inequality.*
- Generated once per batching interval
  - Updates summary model based on the incoming objects in the last interval.
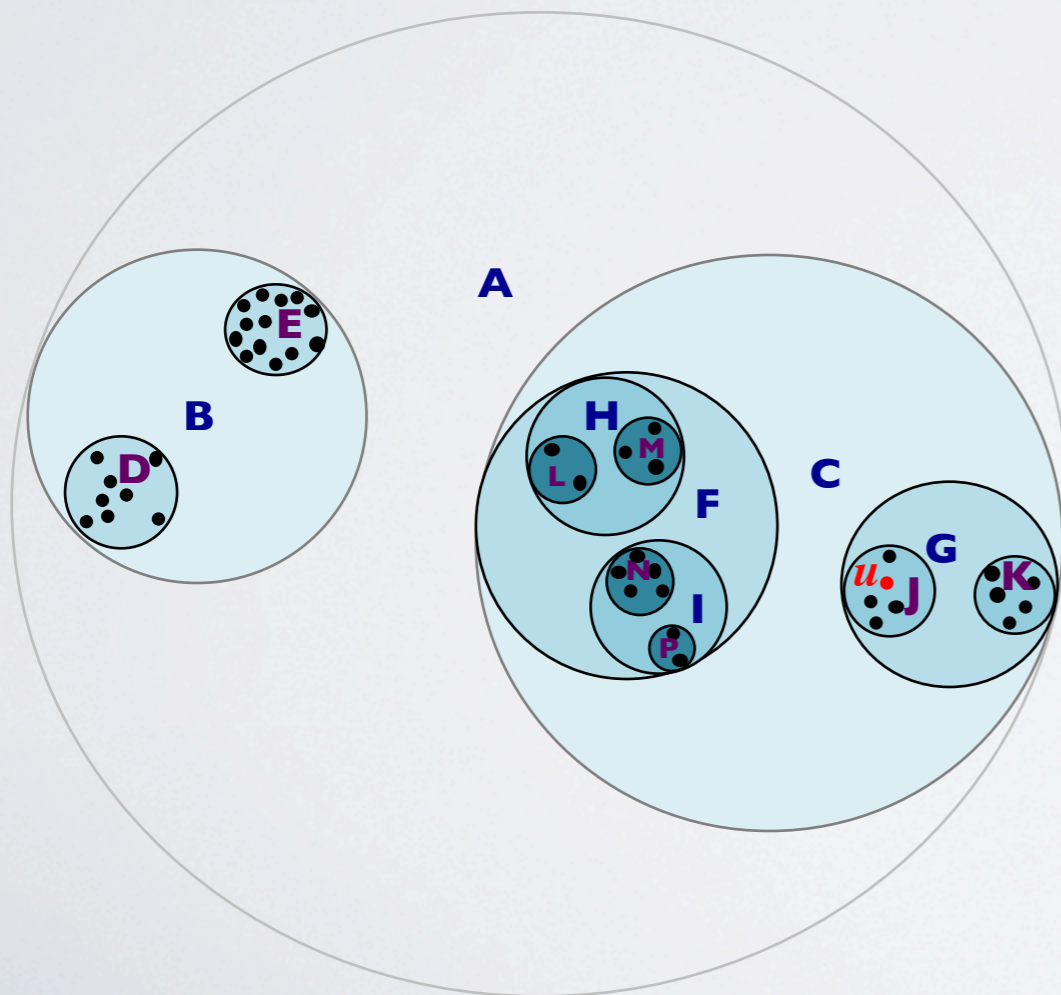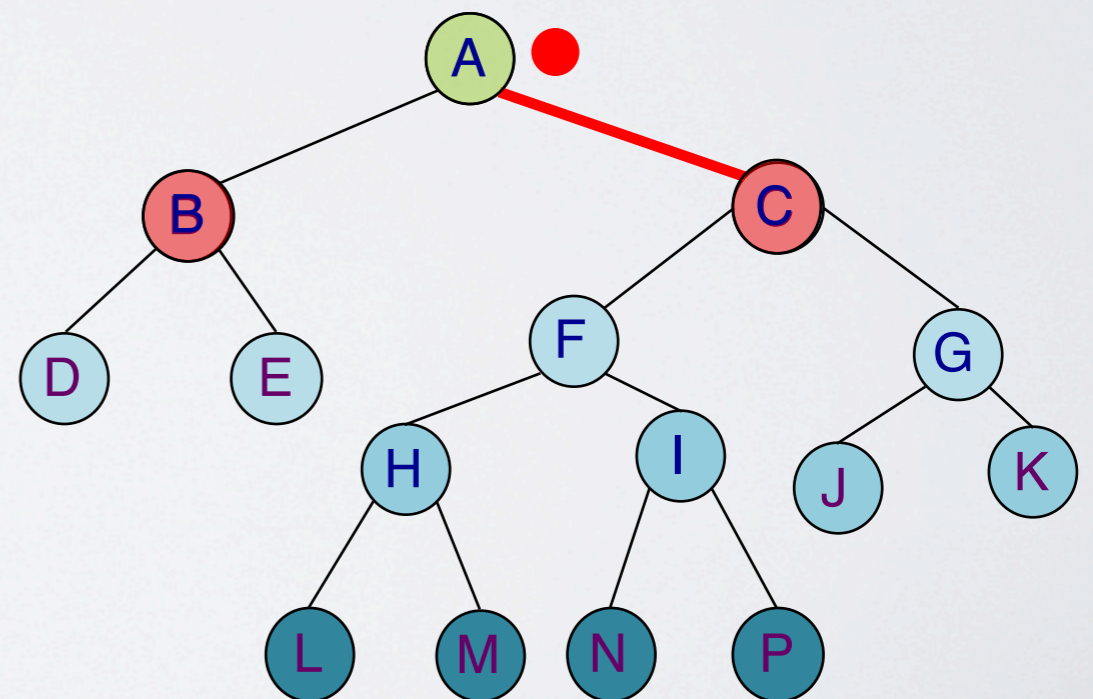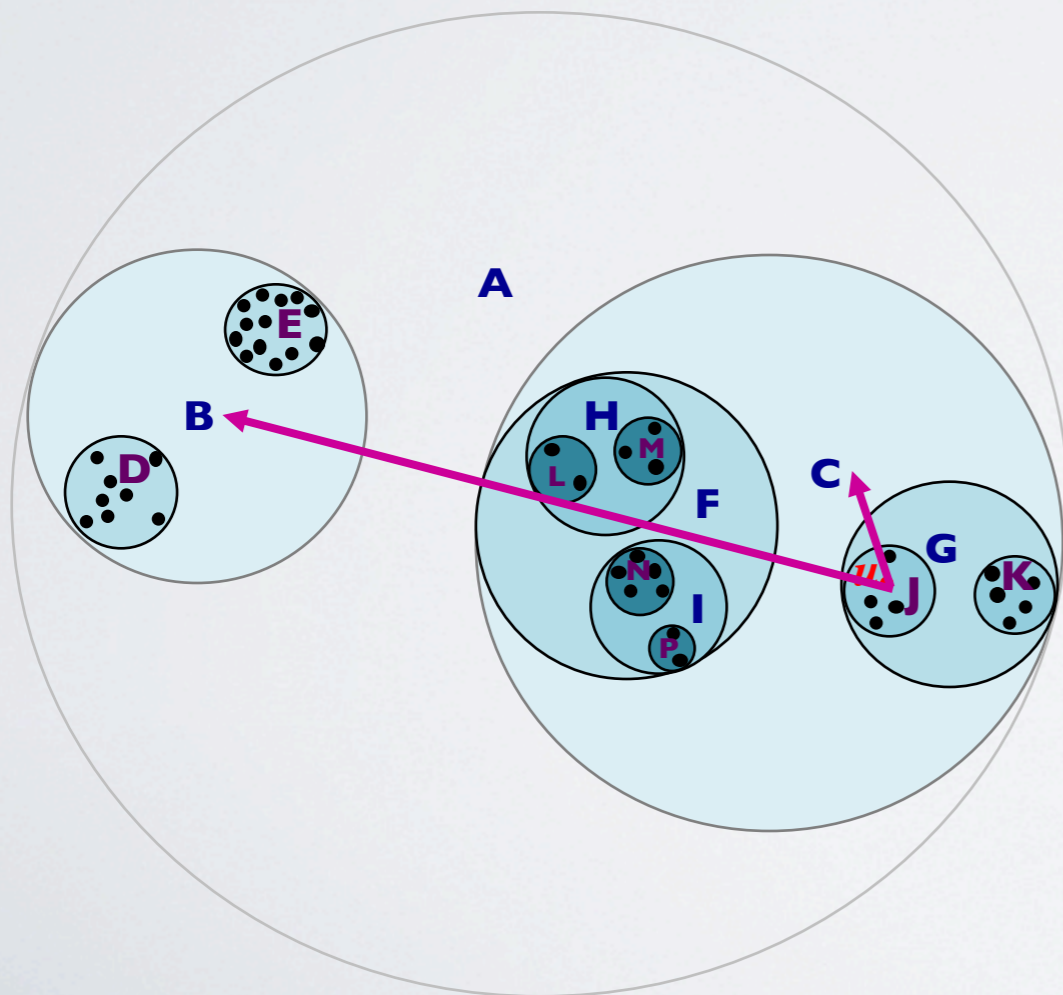
# Insertion

- For example, consider two dimensional space and Euclidean distance
- Insert object u to appropriate cluster

# Insertion
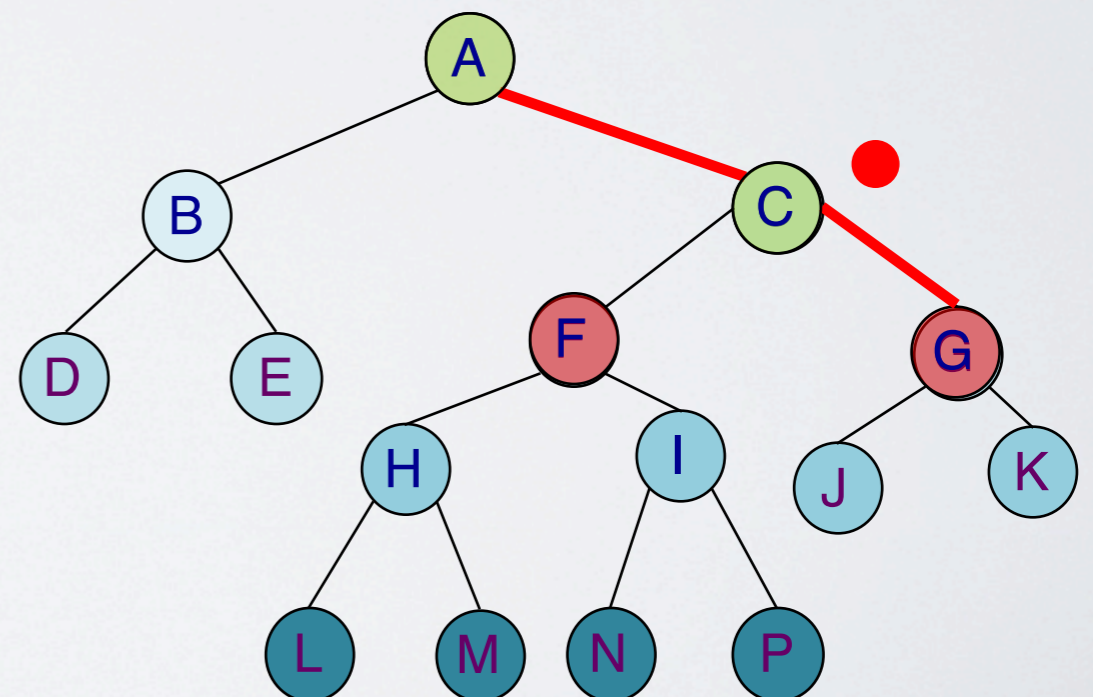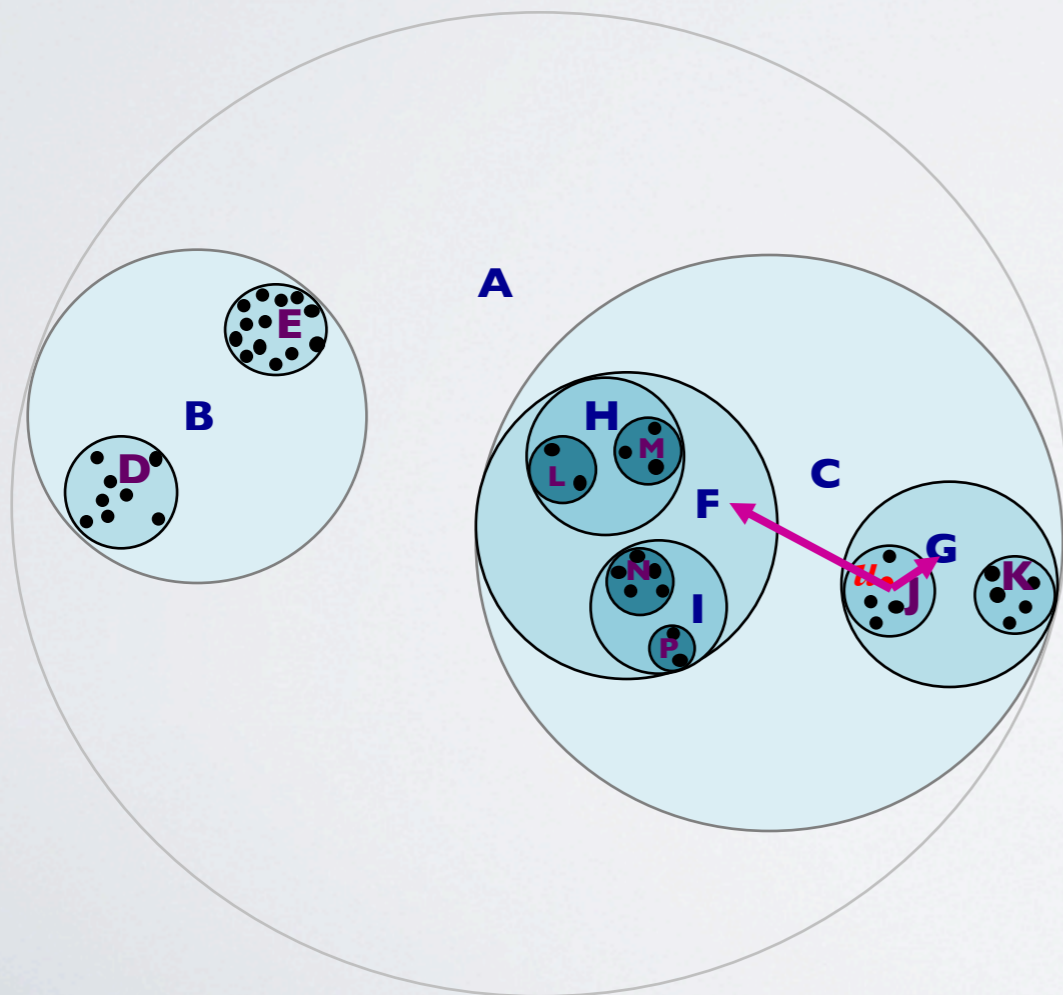
- u arrives at root, A
- u compared to B and C
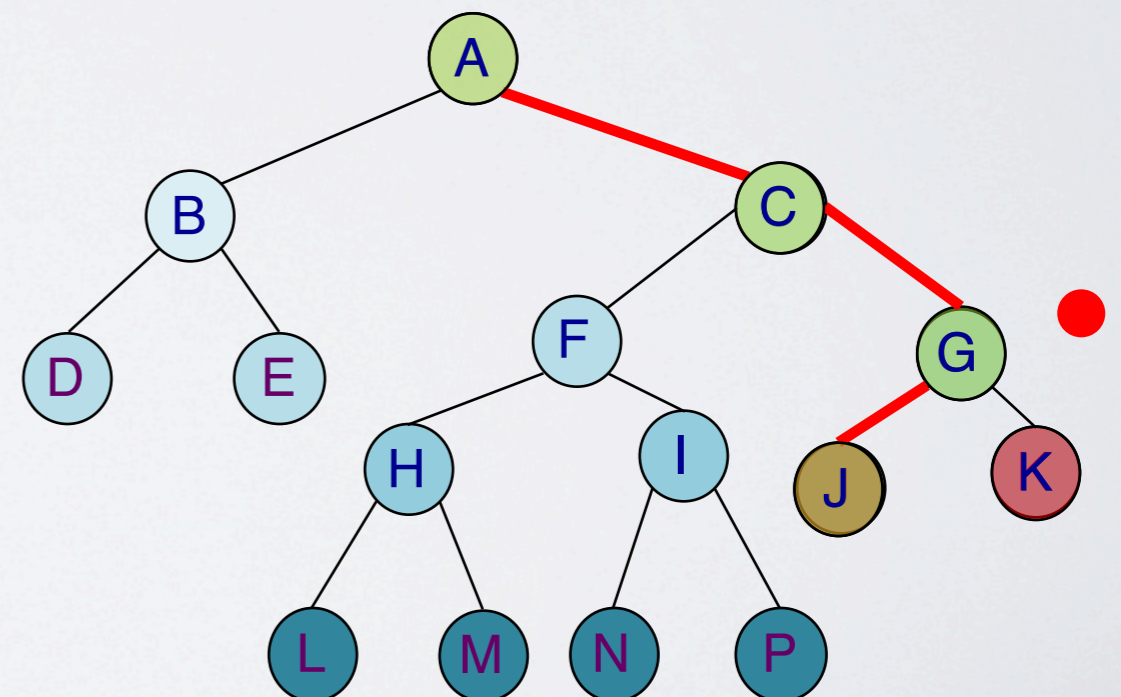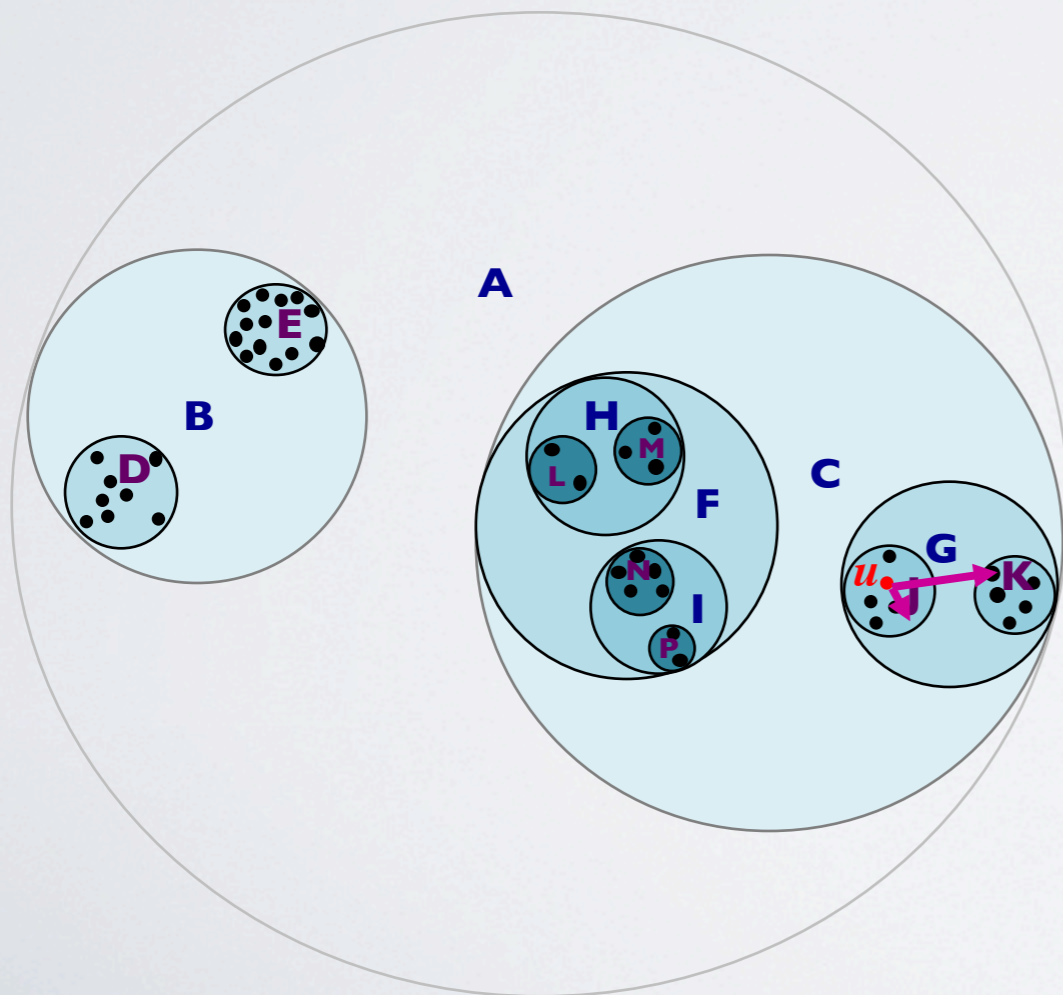- u pushed to C

# Insertion

- u compared to F and G
- u pushed to G

# Insertion

- u compared to J and K
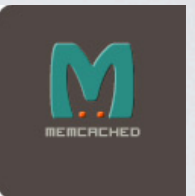- u pushed to J
- u assigned to cluster J

# Twitter Application

- Twitter based Fact-finder application
  - Tweets crawled using search api during events.
    - Crimea, Sandy, Fukushima, Egypt Revolution etc.
  - More than 4 million tweets.
- Customization Interface
  - Tweets converted to high dimensional vectors of tokens and frequencies.

$$d(\mathbf{u}, \mathbf{v}) = 1 - \frac{\mathbf{u}.\mathbf{v}}{\mathbf{u}^2 + \mathbf{v}^2 - \mathbf{u}.\mathbf{v}}$$

  - Used Tanimoto distance function.
- Summary Query Interface
  - Queries SocialTrove using keywords.
  - Ranks matching summaries according to credibility.

# Implementation

- Summary Model
  - Written in Java
  - MapReduce job to generate model runs on Spark.
  - Broadcasts serialized model every interval.
- Client Proxy
  - Java / Python – based on different client applications.
  - Serves matching objects from distributed memcache
- Input Proxy
  - Python – periodically polls twitter via search api
- Distributed Storage
  - HDFS (Hadoop Distributed File System)

# SocialTrove Testbed

- Runs on UIUC Green Data Center
    - http://greendatacenters.web.engr.illinois.edu
    - Capacity: 40 large nodes, 40 small nodes
- Testbed Setup
    - 8 nodes
    - Intel Xeon E5-2620 – 12 logical cores
    - 16 GB Ram per node
    - 1 Gbps Ethernet – Star topology
    - 1 TB disk
- One driver machine
    - apollo3.cs.illinois.edu
    - Intel Xeon E5-2690 – 40 logical cores
    - 128 GB Ram
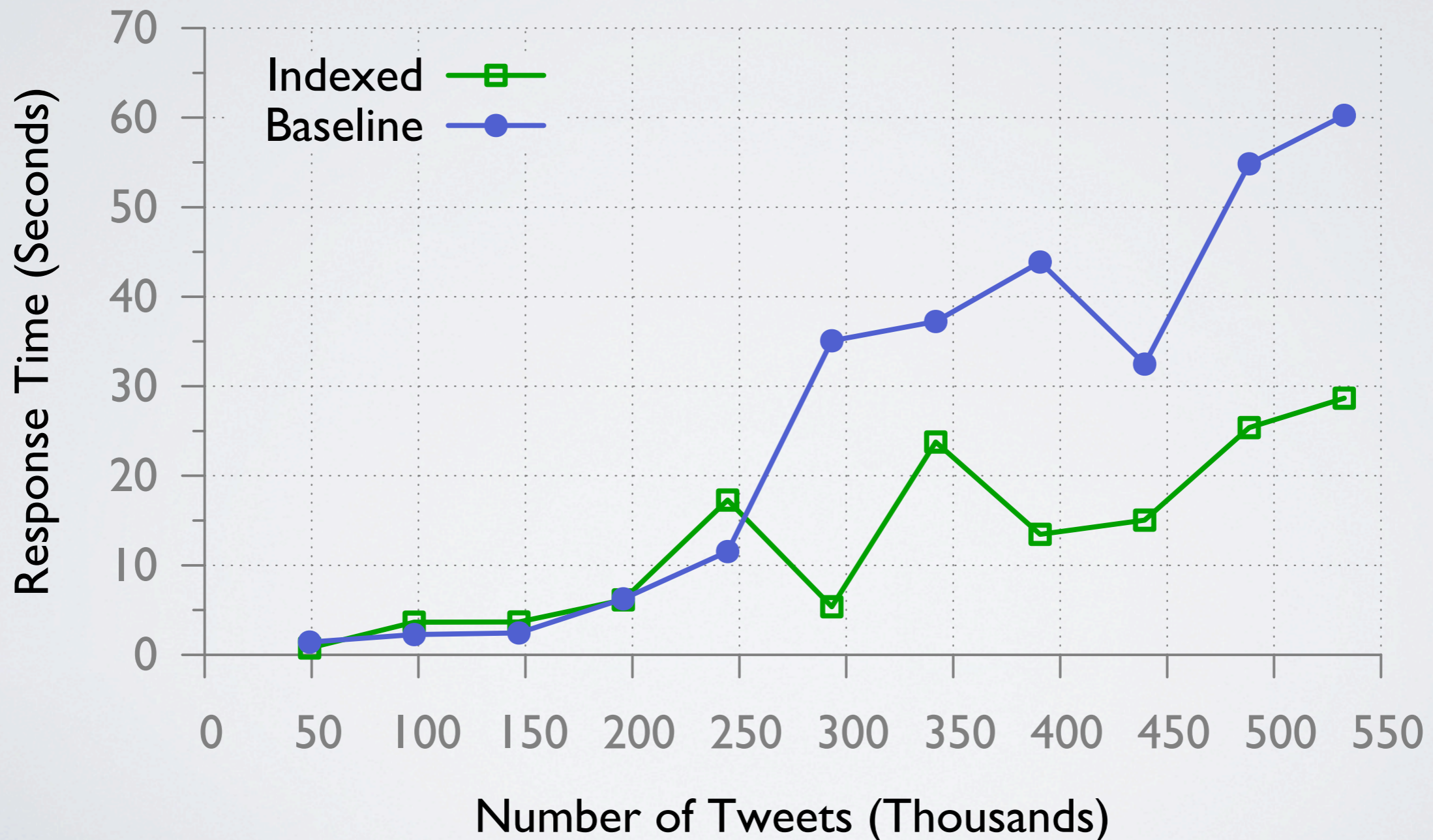    - 16 TB Raid-5 disk

# Evaluation

- SocialTrove compared to 4 schemes
  - **Baseline**
    - No summary model computed.
    - Load balanced assignment to storage.
  - **Indexing**
    - No summary model computed.
    - Tweets indexed by keywords.
  - **Summary Baseline**
    - Computes a flat summary model in advance.
  - **Summary Indexing**
    - Computes a keyword indexed summary model.

# Response Time

## without summary model

Response time is very low with on demand summarization

# Response Time

## with summary model



4 million tweets

Legend:
- SocialTrove
- Summary Indexed
- Summary Indexed Diverse
- Summary Baseline
- Summary Baseline Diverse

Y-axis: Response Time (Milliseconds) — 0.1, 1, 10, 100, 1000, 10000

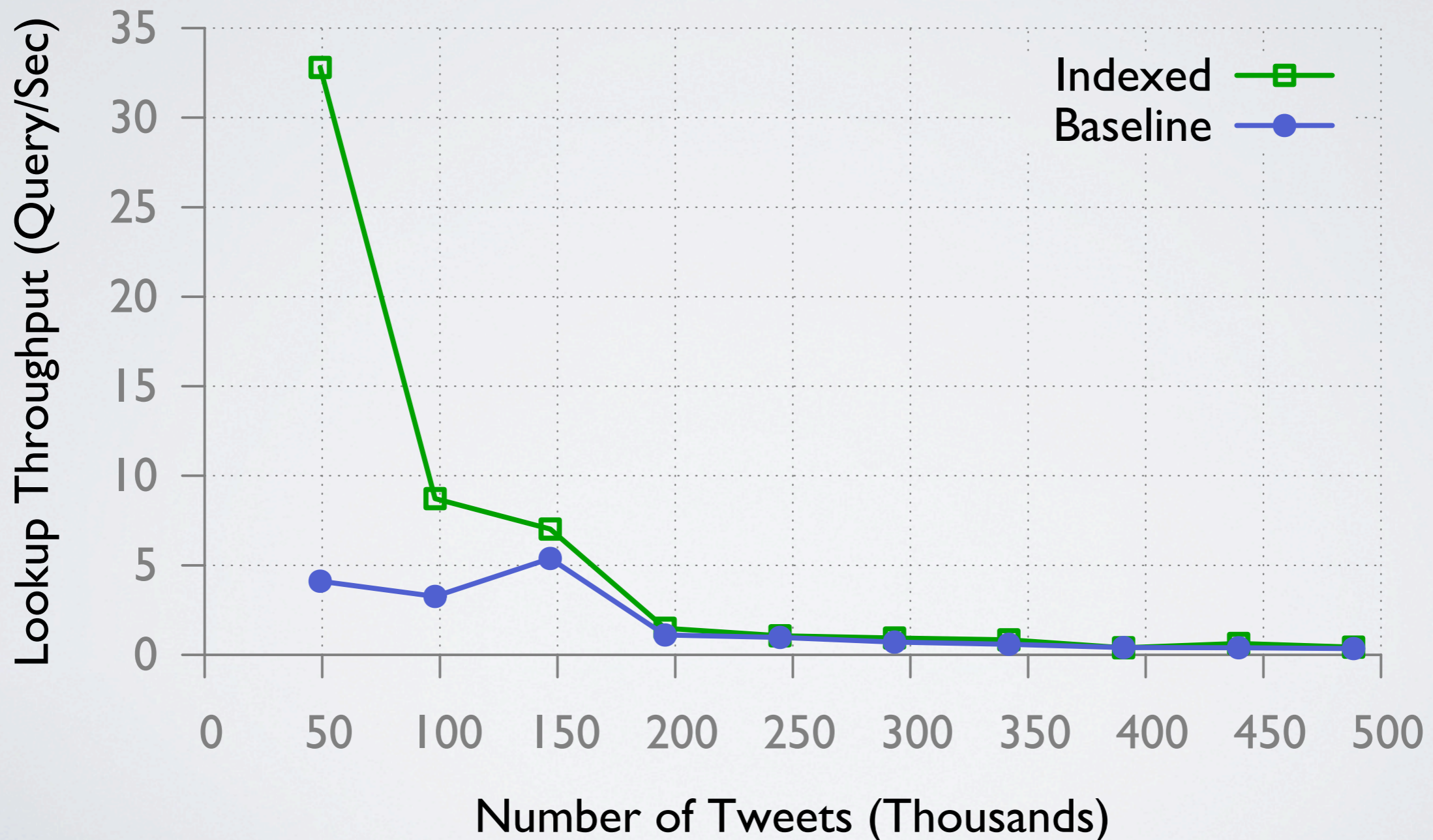X-axis: Load (Query/Second) — 0K, 5K, 10K, 15K, 20K, 25K, 30K

# Lookup Throughput
## **without** summary model

On demand summarization is not scalable

# Lookup Throughput
## with summary model

### 4 million tweets



Legend:
- SocialTrove
- Summary Indexed
- Summary Indexed Diverse
- Summary Baseline
- Summary Baseline Diverse

Y-axis: Lookup Throughput (Query/Sec) — 0K, 5K, 10K, 15K, 20K, 25K

X-axis: Number of Summaries requested from 4 Million tweets — 10, 20, 30, 40, 50, 60, 70, 80, 90, 100

# Model Generation Time



8 worker machines

| Number of Workers | Model Generation Time (4 million tweets) |
|---|---|
| 4 | 337 minutes |
| 8 | 239 minutes |
| 11 | 195 Minutes |

# Ranked Summary

**Protests – Real-time Summary - Thu 25 Jun 2015 07:19:02 GMT**

Thu 25 Jun 2015
07:19:02 GMT
Score : 25
Sources

OG: "Chinese war veterans protest Beijing over missing benefits: http://t.co/T4x3mgtXhB"

Thu 25 Jun 2015
07:19:02 GMT
Score : 23
Sources

Massive protest by AAP Youth Wing against smriti Irani. @AamAadmiParty @DrKumarVishwas @dilipkpandey @ArvindKejriwal http://t.co/EkMbK5MU2w

Thu 25 Jun 2015
07:19:02 GMT
Score : 16
Sources

Video: Activists protest Israeli killer drones at Paris Air Show http://t.co/q0fKmmYBzk

Thu 25 Jun 2015
07:19:02 GMT
Score : 15
Sources

On the day of the #SaveILF protest against cuts for the disabled, I cannot begin to describe how angry this makes me: http://t.co/RsRDhjjuvF

Thu 25 Jun 2015
07:19:02 GMT
Score : 13
Sources

Thousands of Buddhists protest to get their Religious Freedom #DalaiLama #Glastonbury http://t.co/HJrXjGoxAs #UK

# Ranked Summary

| Thu 25 Jun 2015 07:29:15 GMT<br>Score : 21<br>Sources | #Cholera outbreak in #SouthSudan is wake up call to gov't & humanitarian agencies http://t.co/LL0xRLqMlW http://t.co/2O4u40lMkZ |
|---|---|
| Thu 25 Jun 2015 07:29:15 GMT<br>Score : 21<br>Sources | MSG' movie showcases the need for citizens to adopt various humanitarian causes and support welfare activities #MSGFilmDVDLaunched |
| Thu 25 Jun 2015 07:29:15 GMT<br>Score : 19<br>Sources | There is no humanitarian crisis in Calais just a bunch of thugs trying to enter the UK illegally. The law should deal with ALL criminals.... |
| Thu 25 Jun 2015 07:29:15 GMT<br>Score : 18<br>Sources | #MSGFilmDVDLaunched the purpose and motivation behind the film @MSGTheFilm is to promote Humanitarian Activities |
| Thu 25 Jun 2015 07:29:15 GMT<br>Score : 18<br>Sources | Vayapam scam in MP, Lalit Gate in Rajasthan,Humanitarian scam in Delhi, fakes in debate......good speed |
| Thu 25 Jun 2015 07:29:15 GMT<br>Score : 11<br>Sources | Dear media, as soon as you are done with worshipping BBC, pls make sure to report the humanitarian crises in Karachi due to the Heatwave. |
| Thu 25 Jun 2015 07:29:15 GMT<br>Score : 10<br>Sources | 9 humanitarian INGOs refused registration. Ironically, the UN banned JuD was allowed to work in #Awaran #Balochistan http://t.co/GE3mtPfCj8 |

26

# Ranked Summary

Thousands at Moscow rally against Russian intervention in #Ukraine

Man in Ukraine plays the piano to help calm down a riot. http://t.co/fdNAc0cfJ2

For Crimea, Google Shows Different Borders Based on Your Location: Russia's Minister of Communications and Mass Media http://t.co/vlHGYIibOC

Militants in eastern #Ukraine were equipped with Russian weapons and the same uniforms as those worn by Russian forces that invaded Crimea.

50,000 #Ukraine supporters march in Moscow protest Russia's intervention in #Crimea.

Some russian tanks on ukrainian border already painted with 'peacekeeping' slogans. How much longer until the 'humanitarian intervention'?

I've been speaking to @BarackObama about the situation in Ukraine. We are united in condemnation of Russia's actions. http://t.co/7Rk2k8iOIK

Ukrainian Defense Ministry says its lone submarine has been taken by Russians.

Ukraine prepares armed response as city seized by pro-Russia forces

Ukraine crisis: Nato warns Russia against further intervention - BBC News

# Conclusions

- Summarizes social content stream with arbitrary granularity, in real time.
- Caches a summary model in memory.
  - Scale-out input and client query proxies.
- High throughput and low response time.
  - Thanks to asynchronous updates to summary model, and temporal locality in social stream.
- Outperforms traditional indexing methods.
- Limitation: Evaluation performed only in the context of tweets.