

Philipp Zormeier

Event Detection in Social Streams

14.11.2012
„Mining Volatile Data“





I. Introduction

II. Social stream model

III. Unsupervised approach

IV. Supervised approach

V. Performance evaluation

VI. Summary



What is Event Detection?

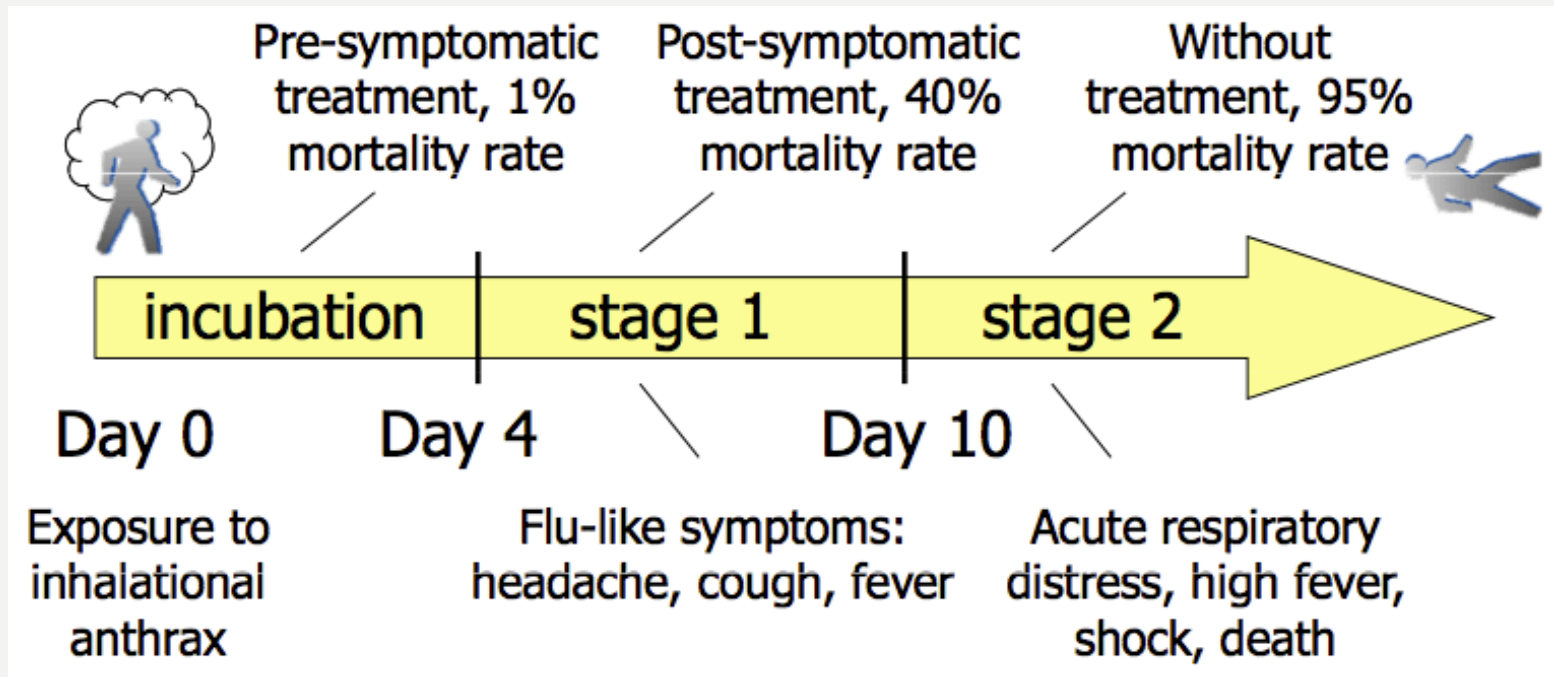
- Analysis of monitoring data
- Detection of interesting changes
- Characterization of the event





Applications of Event Detection?

Famous example:





Social Streams

- People post about their lives
- Important events are captured in bursts of posts
- Continuous interaction
- Posts contain **temporal**, **structural** and **content** information



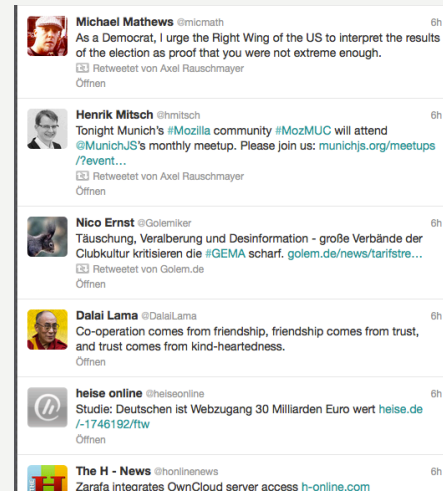
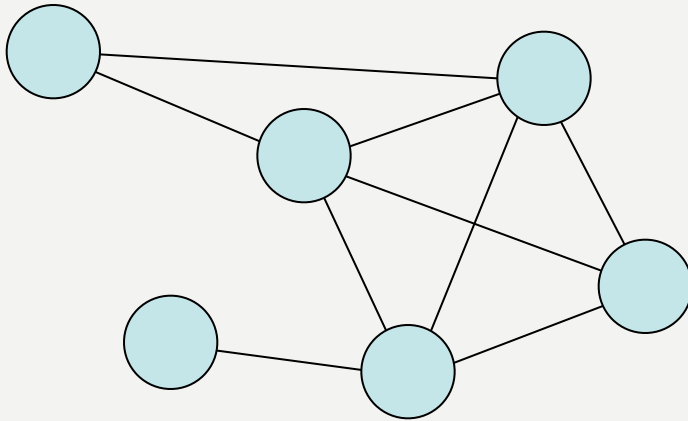


Key Challenges

- (i) Ability to use both the content and the structure of the interactions for detection
- (ii) Ability to use temporal information
- (iii) Ability to handle very large and massive volumes of text documents under the one-pass constraint

What do we call a Social Stream?

Social Stream = Structure + Content



Michael Mathews @micmath 6h
As a Democrat, I urge the Right Wing of the US to interpret the results of the election as proof that you were not extreme enough.
Retweeted von Axel Rauschmayer
Offnen

Henrik Mitsch @hmitsch 6h
Tonight Munich's #Mozilla community #MozMUC will attend @MunichJS's monthly meetup. Please join us: munichjs.org/meetups/?event...
Retweeted von Axel Rauschmayer
Offnen

Nico Ernst @Golemiker 6h
Täuschung, Verabberung und Desinformation - große Verbände der Clubkultur kritisieren die #GEMA scharf. goiem.de/news/tarifstre...
Retweeted von Golem.de
Offnen

Dalai Lama @DalaiLama 6h
Co-operation comes from friendship, friendship comes from trust, and trust comes from kind-heartedness.
Offnen

heise online @heiseonline 6h
Studie: Deutschen ist Webzugang 30 Milliarden Euro wert heise.de/-1746192/ftw
Offnen

The H - News @honlineews 6h
Zarafa integrates OwnCloud server access h-online.com

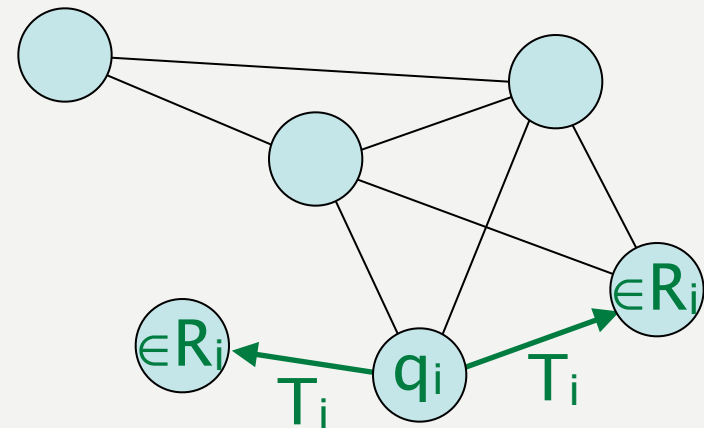


Definition of a Social Stream

Continuous and temporal Sequence of objects $S_1 \dots S_r \dots$
such that each $S_i = (q_i, R_i, T_i)$ contains..

- a text document T_i
- an origination node $q_i \in N$
- a set of receiver nodes

$$R_i \subseteq N \quad (\forall r \in R_i \quad (q_i, r) \in A)$$



Graph $G = (N, A)$



Examples: Facebook



Lars Butnotleast

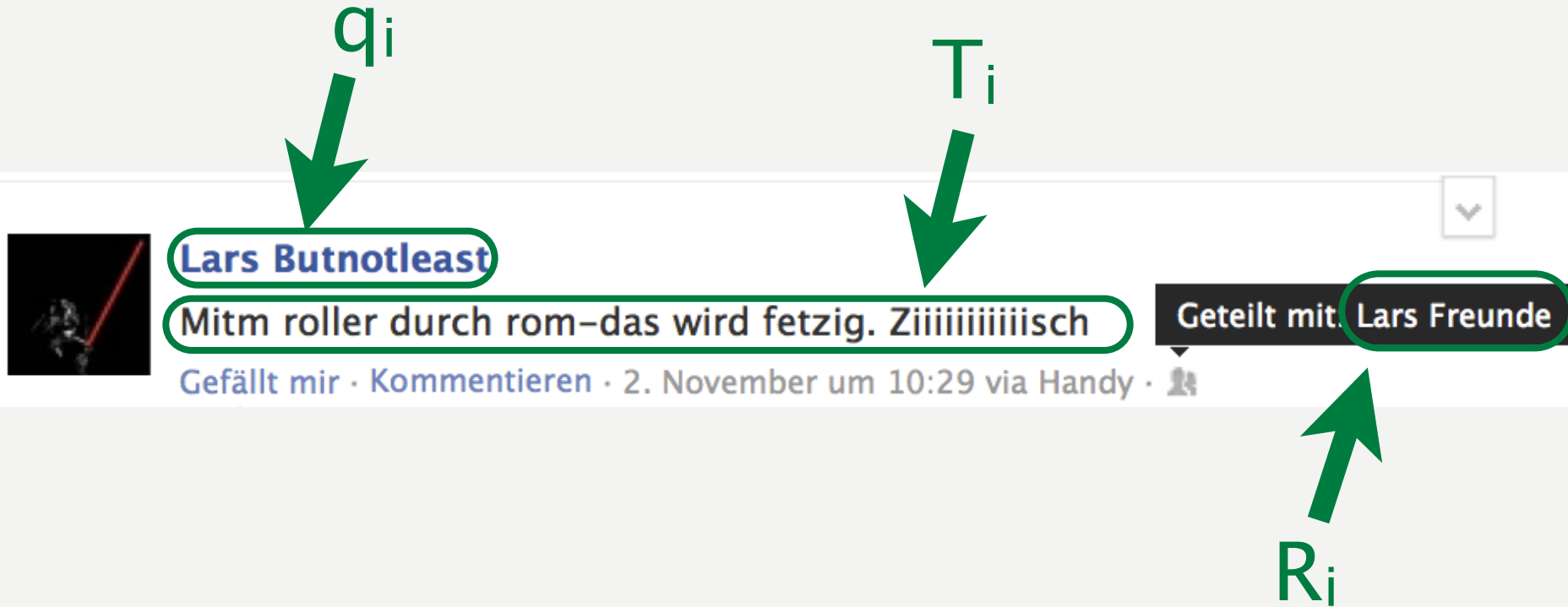
Mitm roller durch rom–das wird fetzig. Ziiiiiiiiisch

Gefällt mir · [Kommentieren](#) · 2. November um 10:29 via Handy ·

Geteilt mit: Lars Freunde



Examples: Facebook





Examples: Twitter



rammdoesig @rammdoesig

7 Nov

Mich interessiert eigentlich nur noch wie Springfield und Southpark abgestimmt haben, dann schlafen. [#uswahl](#) [#election2012](#)

Öffnen



Examples: Twitter

 q_i  T_i **rammdoesig** @rammdoesig

7 Nov

Mich interessiert eigentlich nur noch wie Springfield und Southpark abgestimmt haben, dann schlafen. #uswahl #election2012

Öffnen

R_i not visible here (all followers)



Examples: Email

Example



Von:

Philipp Zormeier +

An:

spam@siebenfarben.de +

Hi,

This email is just an example.

Bye,

Philipp



Examples: Email





Data organization

Data points in Clusters

Cluster ~ Topics

New Point in Cluster ~ The actor talks about the topic

Many new points in Cluster ~ Something happened?

New Cluster ~ A new topic comes up



Social stream clustering

Partitioning the stream objects $S_1 \dots S_r \dots$

into k clusters $C_1 \dots C_k$, such that

- for all i : S_i belongs to at most one cluster
- the similarity function uses content and network structure



Cluster summaries

$$\Psi(C_i) = (V_i, \eta_i, W_i, \Phi_i)$$

$V_i = \{j_{i1}, \dots, j_{is}\}$: Set of nodes

$\eta_i = v_{i1} \dots v_{is}$: Node frequencies

$W_i = \{l_{i1}, \dots, l_{iu}\}$: Set of words

$\Phi_i = \varphi_{i1} \dots \varphi_{iu}$: Word frequencies



Similarity of stream objects to clusters

Overall similarity:

$$\text{Sim}(S_i, C_r) = \lambda \cdot \text{SimS}(S_i, C_r) + (1 - \lambda) \text{SimC}(S_i, C_r)$$

$$\lambda \in [0, 1]$$



Similarity of stream objects to clusters

Structure:

$$\text{SimS}(S_i, C_r) = \frac{\sum_{t=1}^{s_r} b_t \cdot v_{rt}}{\sqrt{\|R_i \cup \{q_i\}\|} \cdot \left(\sum_{t=1}^{s_r} v_{rt} \right)}$$

b binary vector: $b_t = 1$, if node $j_{rt} \in R_i \cup \{q_i\}$.



Similarity of stream objects to clusters

Content: $\text{SimC}(S_i, C_r) = \text{TF-IDF}$



Similarity of stream objects to clusters

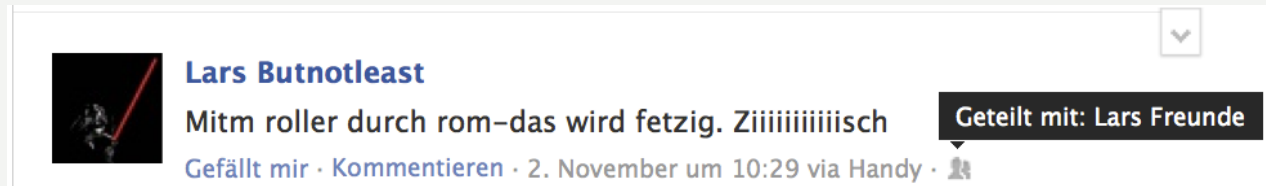
Overall similarity:

$$\text{Sim}(S_i, C_r) = \lambda \cdot \text{SimS}(S_i, C_r) + (1 - \lambda) \text{SimC}(S_i, C_r)$$

$$\lambda \in [0, 1]$$



Assignment to clusters



New object

- » Find closest cluster
- » Similar enough? ($\text{Sim}(S_i, C_r) > \mu - 3 \cdot \sigma$)
 - » no \rightarrow create new cluster
 - » yes \rightarrow assign to cluster



Case 1: Creation of clusters

A new data point is placed in its own new cluster.
The **most stale** cluster is replaced.

most stale = last recent updated

This is a **novel event**.



Case 2: Assignment to clusters

The data point is assigned to an existing cluster.

- Update cluster summary
- Check, if **evolution event** occurred:
 - ⇒ Calculate fractional cluster presence



Case 2: Assignment to clusters

Fractional cluster presence $F(t - H, t, C_i)$

- Data arrival ratio in time period $(t - H, t)$
- Time horizon H has to be chosen

Evolution event:
$$\frac{F(t_c - H, t_c, C_i)}{F(t(C_i), t_c - H, C_i)} \geq \alpha$$

threshold α

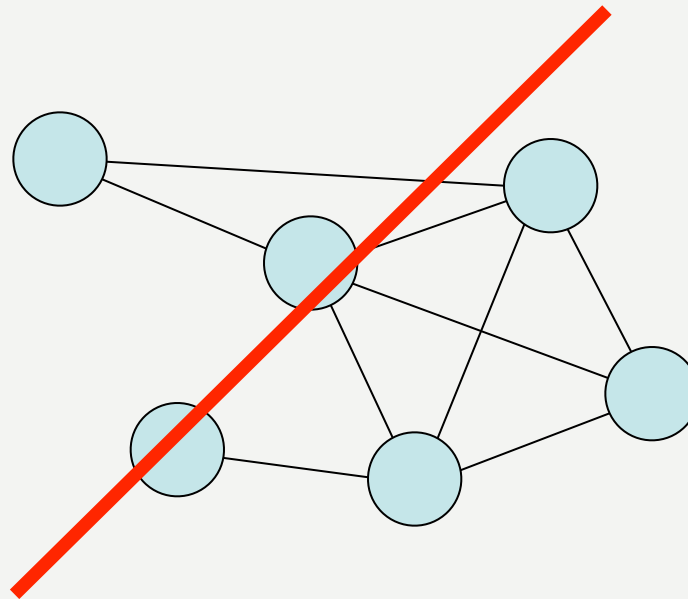


```
while(end of stream not reached)
  i = i + 1;
  Receive next object Si;
  for each cluster C1 compute Sim(Si, C1);
  Let r be index of most similar cluster Cr;
  if(Sim(Si, Cr) < μ - 3 · σ)
    then replace most stale cluster;
    else add Si to Cr and update Ψ(Cr);
  Update μ and σ;
```



```
while(end of stream not reached)
  i = i + 1;
  Receive next object  $S_i$ ;
  for each cluster  $C_1$  compute  $\text{Sim}(S_i, C_1)$ ;
  Let  $r$  be index of most similar cluster  $C_r$ ;
  if( $\text{Sim}(S_i, C_r) < \mu - 3 \cdot \sigma$ )
    then replace most stale cluster;
    ⇒ novel event
  else add  $S_i$  to  $C_r$  and update  $\Psi(C_r)$ ;
    ⇒ check for evolution event
  Update  $\mu$  and  $\sigma$ ;
```

Performance Issues



In reality it is a bit more complex..

Event Detection in Social Streams

II. Social stream model





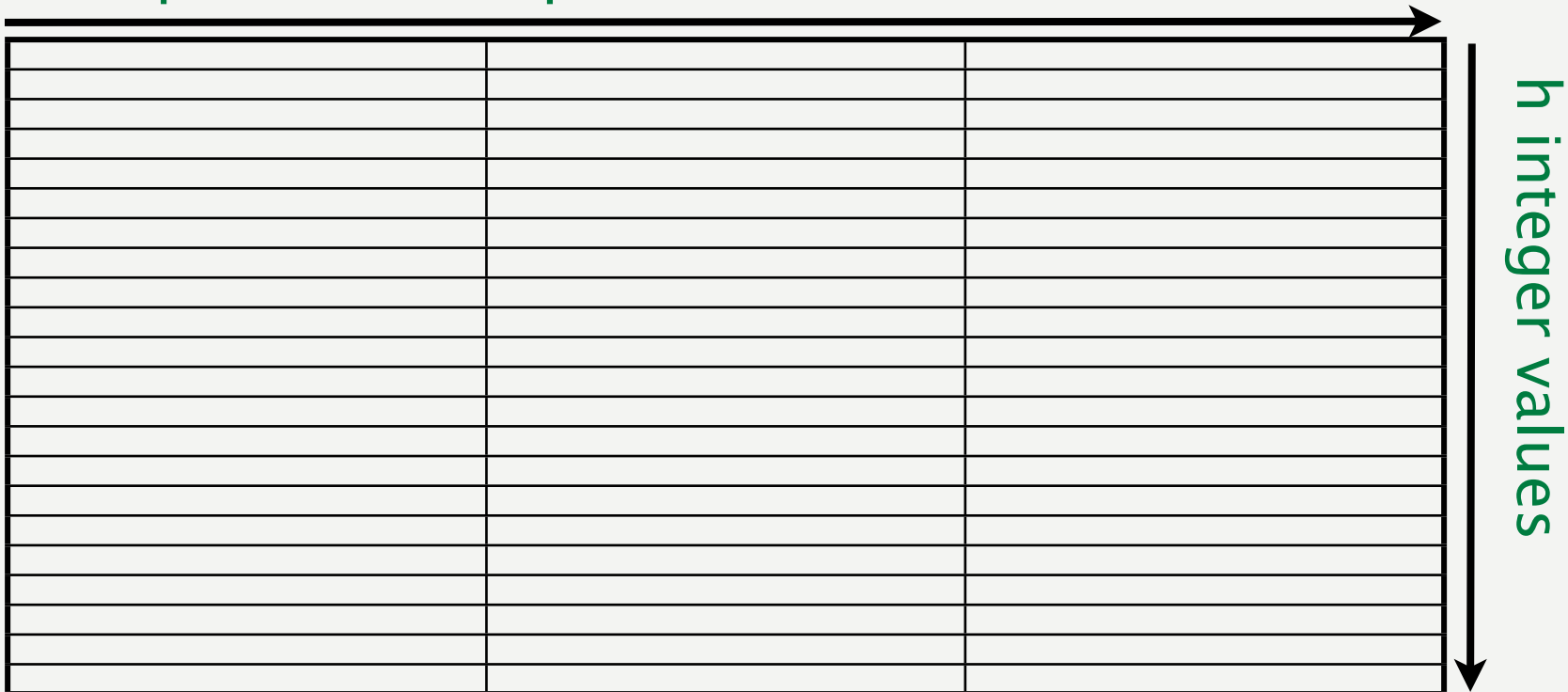
Sketch-based Speedup

- Speed up node counting
- Using Count-min sketch (Hash based)
- Main idea: Estimate counts



Sketch-based Speedup

w pairwise independent hash functions





Supervised Event Detection

Assumptions:

- ✓ Known event E
- ✓ Access to history of the stream
- ✓ Information about set of relevant posts



Variation of Social Stream Clustering

Cluster replacement is not allowed

Event Signature $V(E)$

Vector of relative distribution of event-specific stream objects to clusters

Horizon Signature over $(t_c - H, t_c)$

Vector of relative distribution of arriving points



Supervised Event Detection

- ✓ Calculate dot product of horizon signature and event signature
- ✓ Events are signaled at certain alarm level



Data sets

1) **Twitter Social Stream**

1,628,779 tweets

59,192,401 nodes (avg ~84 nodes/tweet)

2) **Enron Email Stream (filtered)**

349,911 emails

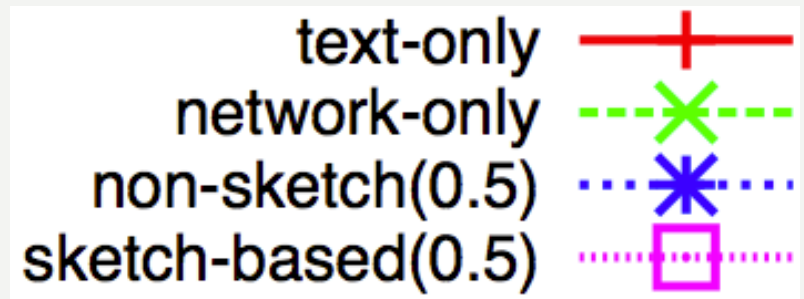
29,083 nodes (avg ~3.62 receivers/email)



Clustering performance

$$\text{Sim}(S_i, C_r) = \lambda \cdot \text{SimS}(S_i, C_r) + (1 - \lambda) \text{SimC}(S_i, C_r)$$

- ✓ $\lambda = 0$ text-only
- ✓ $\lambda = 1$ network-only
- ✓ $\lambda = 0.5$ non-sketch
- ✓ $\lambda = 0.5$ sketch-based



sketch table: $h = 262,213$, $w = 2$



Effectiveness of the clustering

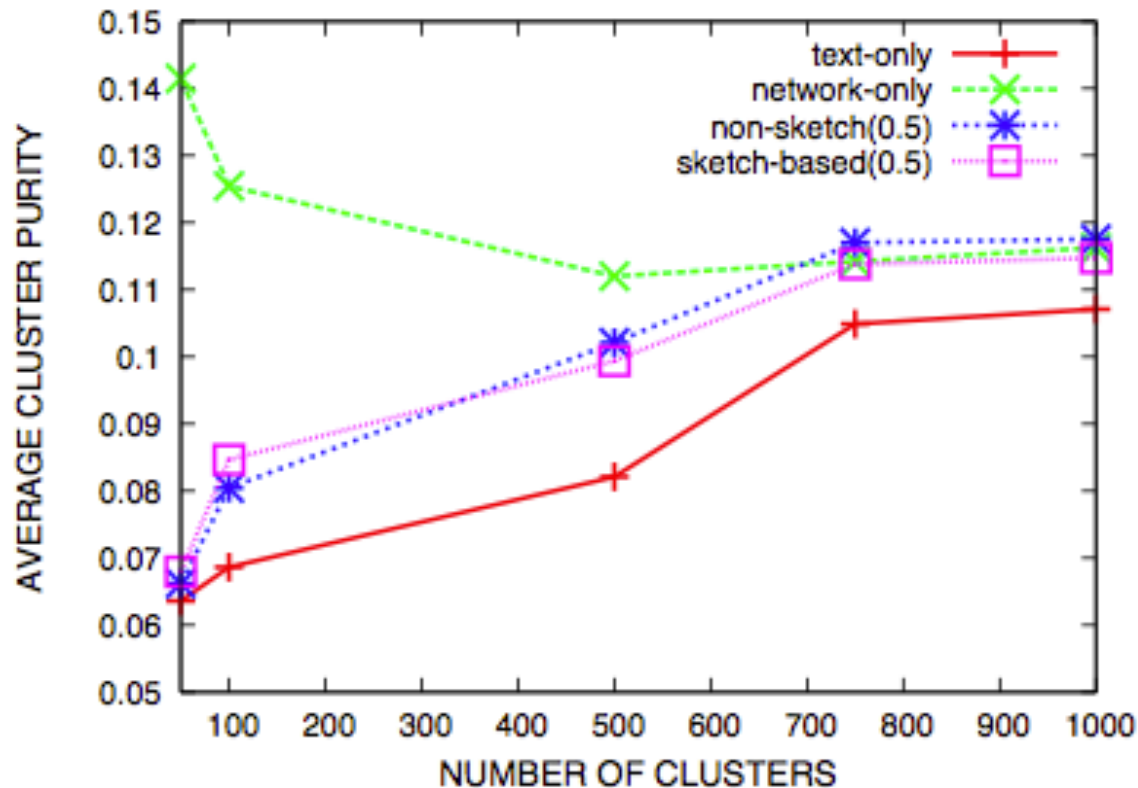
- Assumption: Frequent hash tag → meaningful event
- Test for purity of clusters on dominant hash tags
- Purity = Fraction of objects with dominant tag

Efficiency of the clustering

- Number of stream objects processed by time



Effectiveness of the clustering

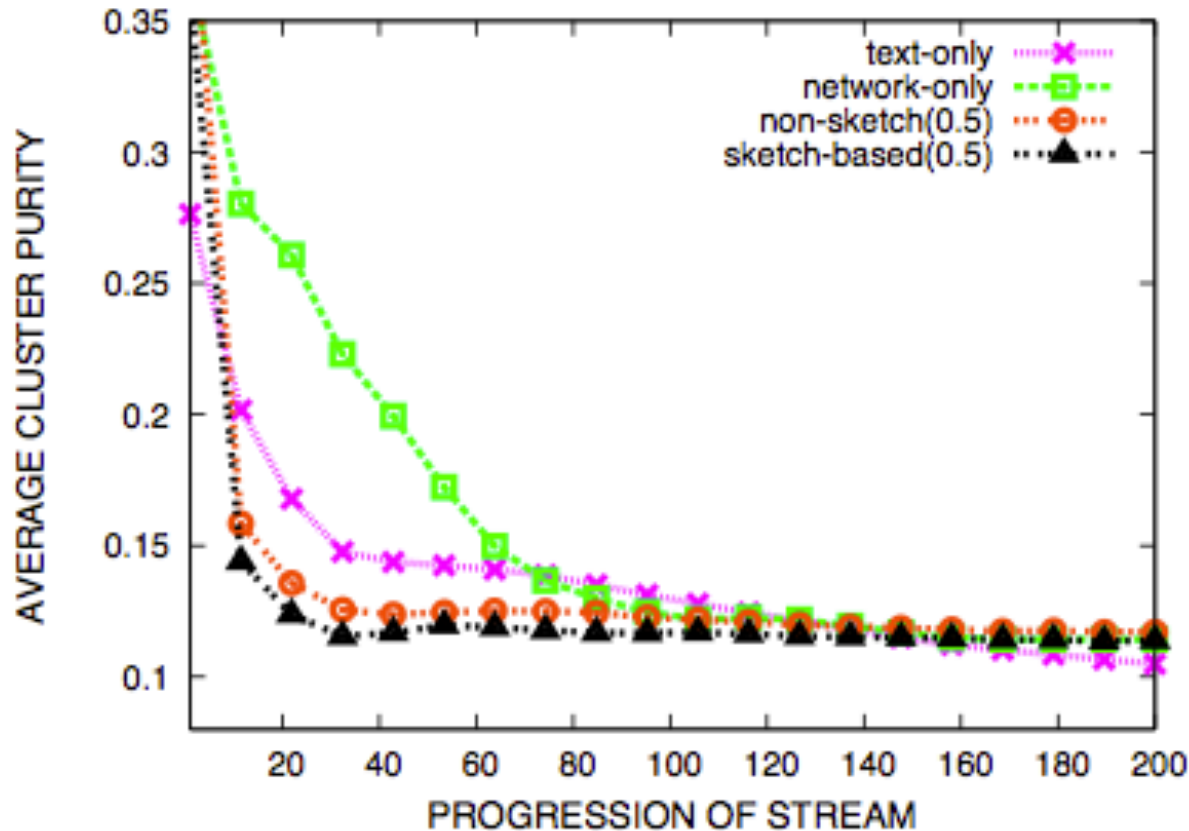


sketch
 $w = 262,213$
 $h = 2$

Effectiveness of the clustering

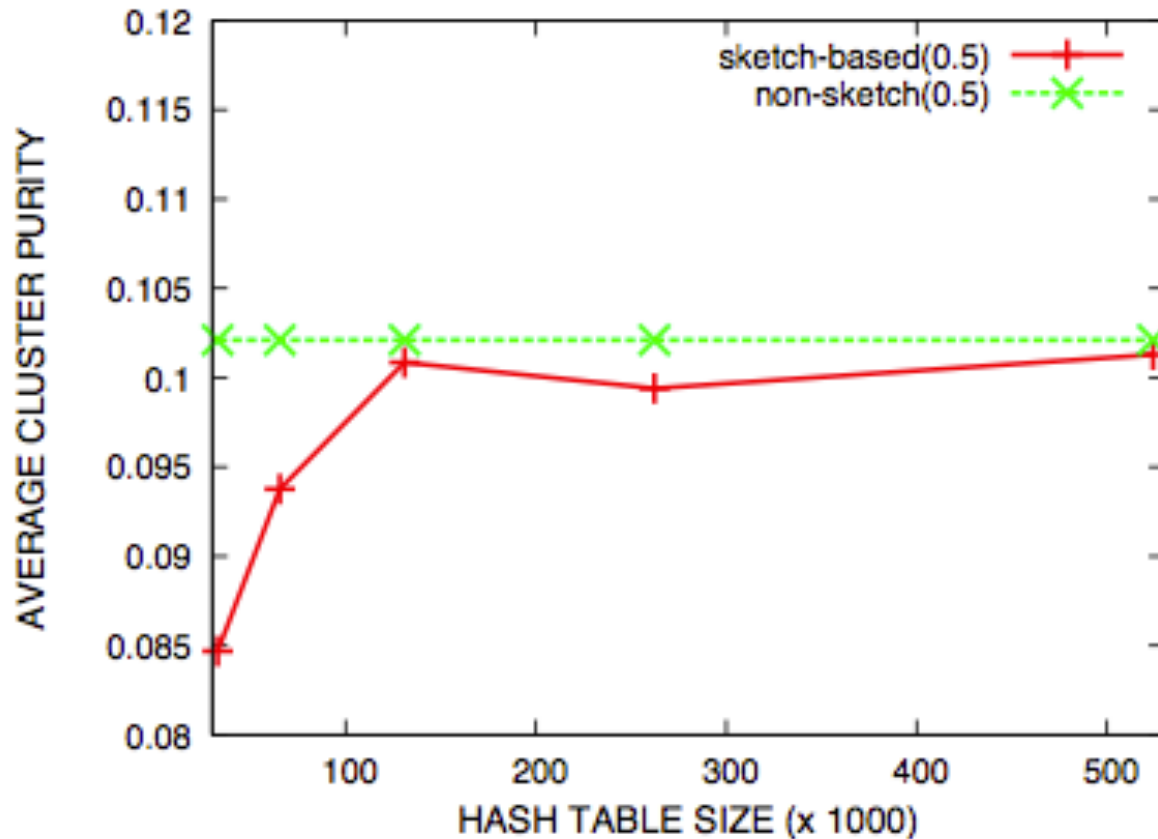
clusters
 $k = 750$

sketch
 $w = 262,213$
 $h = 2$





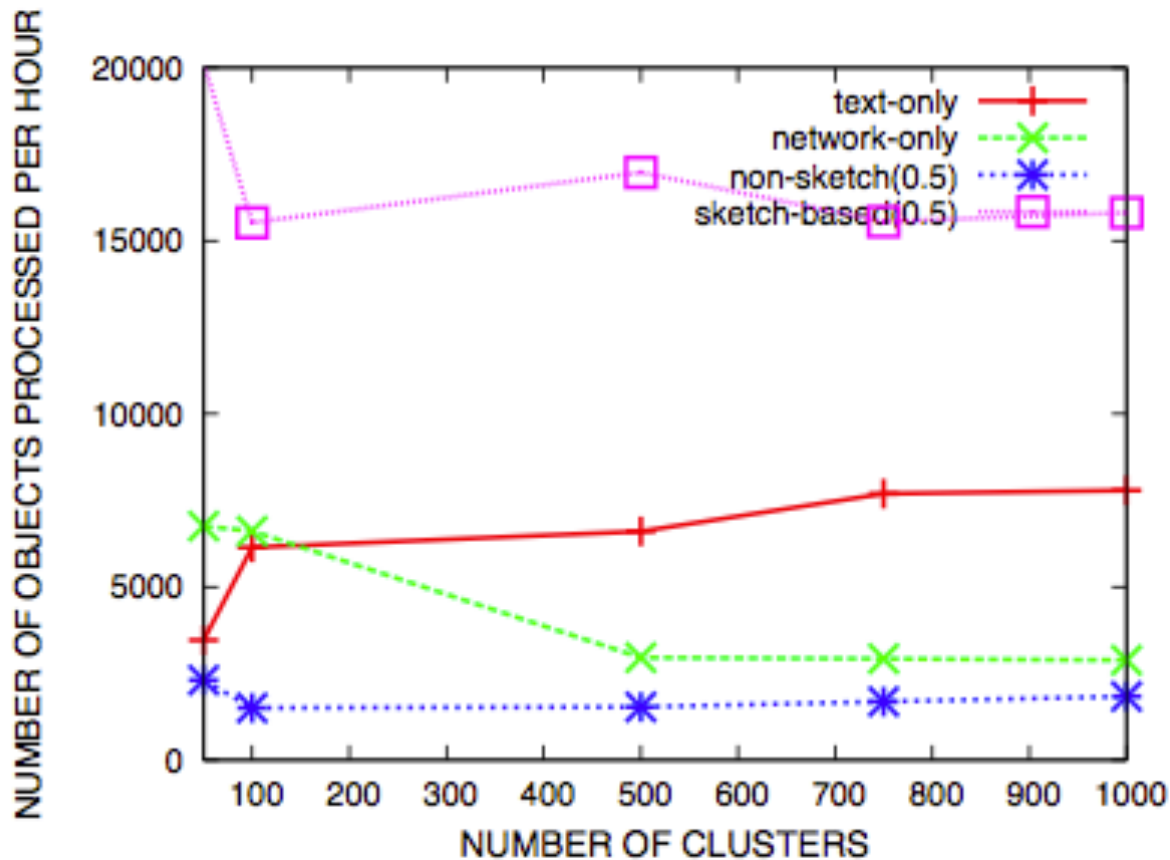
Effectiveness of the clustering



clusters
 $k = 500$



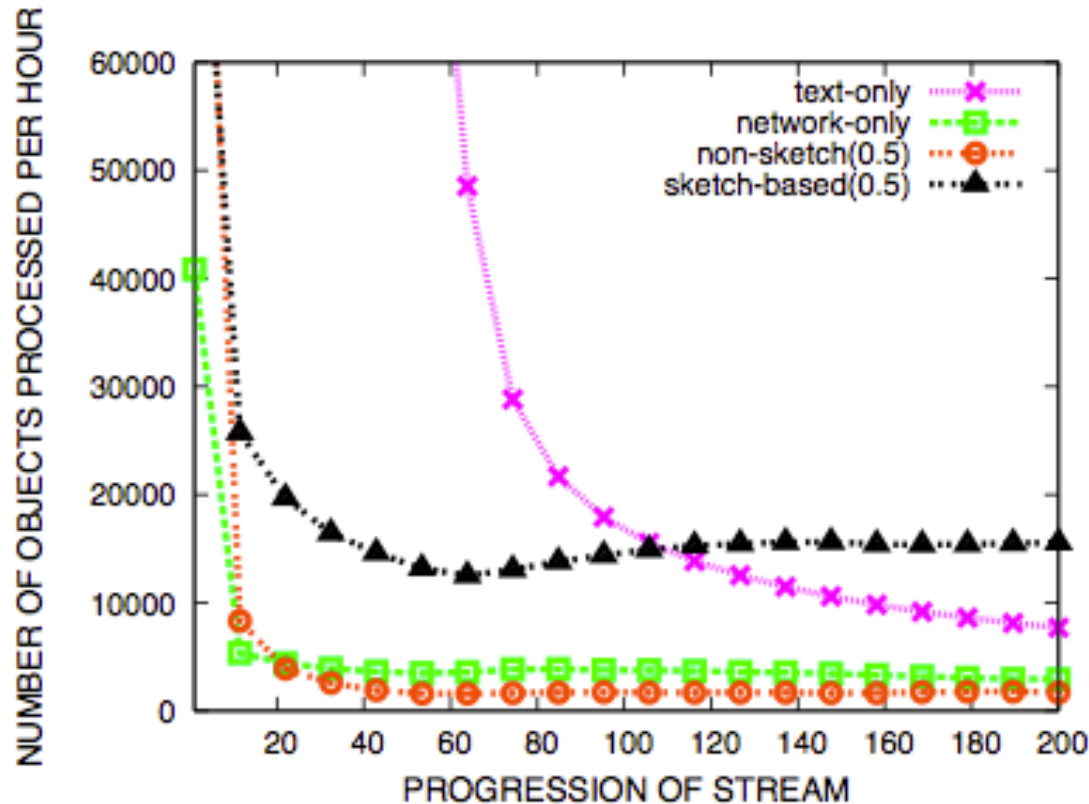
Efficiency of the clustering



sketch
 $w = 262,213$
 $h = 2$



Efficiency of the clustering

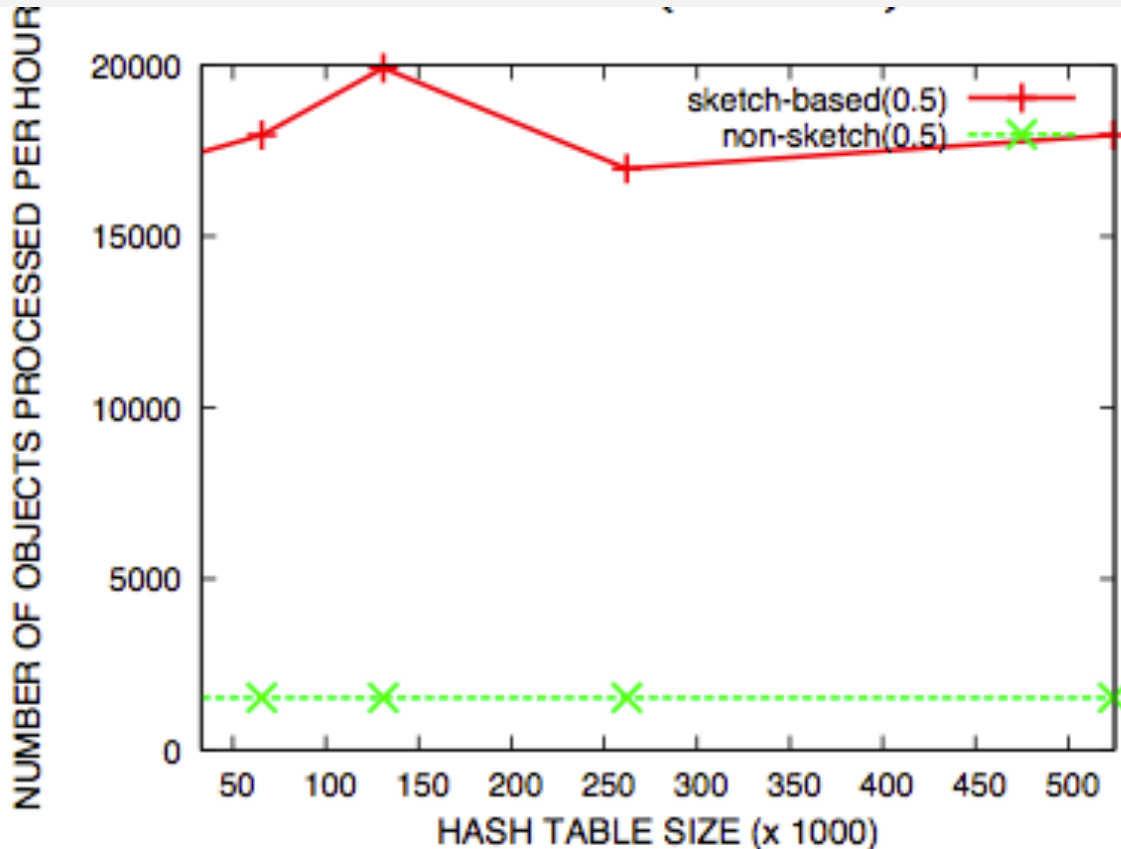


clusters
 $k = 750$

sketch
 $w = 262,213$
 $h = 2$



Efficiency of the clustering



clusters
 $k = 500$



Effectiveness of the event detection (unsupervised)

- Examined within a case study
- Detection of novel and evolution events succeeded
- Detection even with foreign language content
- Connection of related events succeeded

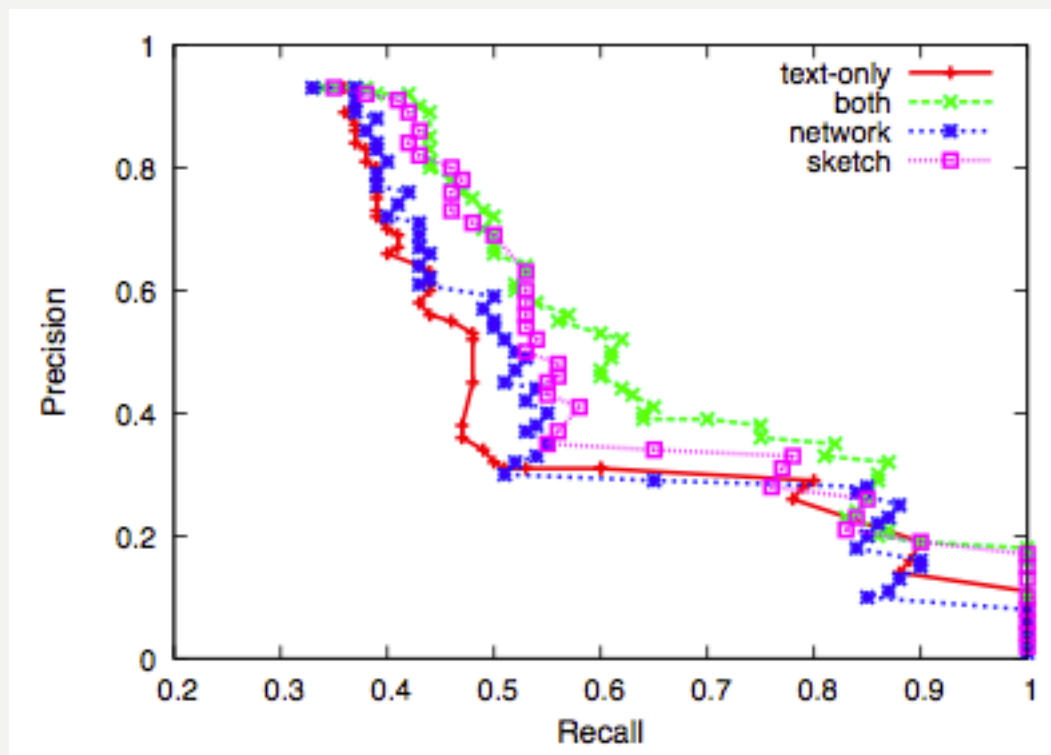


Effectiveness of the event detection (supervised)

- For each period of 5 minutes an event bit was set
- Continuous alarm signal fires at given threshold t
- Variation in threshold t
- Results in tradeoff between precision and recall



Effectiveness of the event detection (supervised)



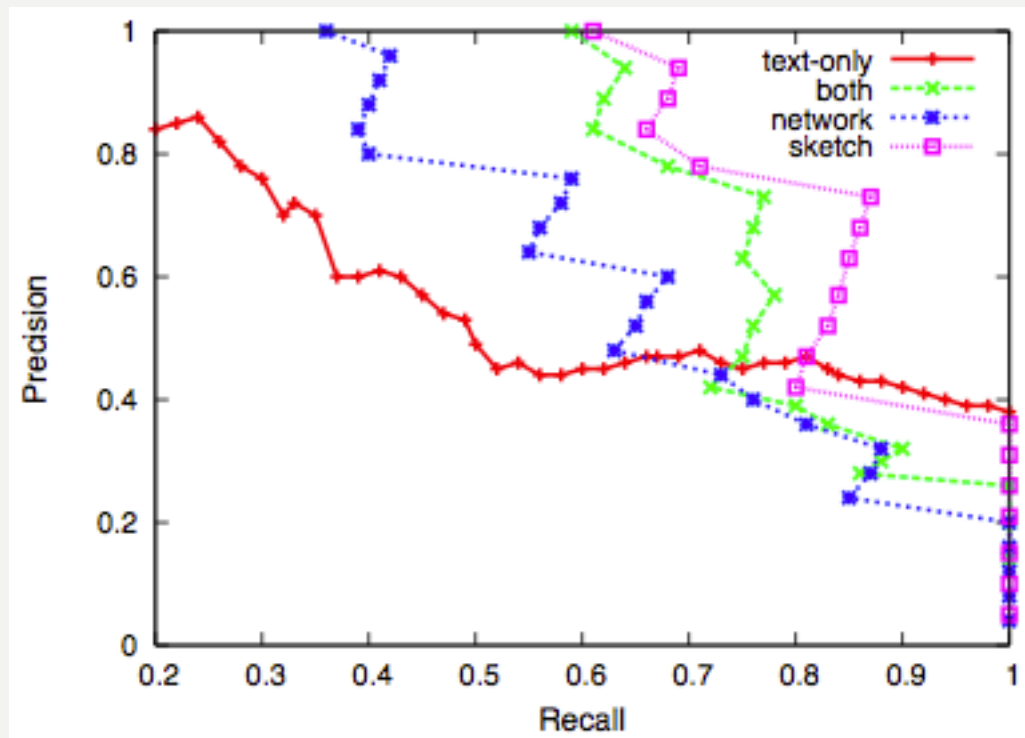
clusters
 $k = 750$

sketch
 $w = 262,213$
 $h = 2$

Japan Nuclear Crisis



Effectiveness of the event detection (supervised)



clusters
 $k = 750$

sketch
 $w = 262,213$
 $h = 2$

Uganda protests



- ++ Usage of content and structure
- + Efficient methods (effective speed up)
- + Effective methods
- Novel event criteria questionable
- Focus of paper on clustering rather than event detection