



Event Detection on Social Media Streams

Antonia Saravanou, Ph.D. Student

antoniasar@di.uoa.gr, www.di.uoa.gr/~antoniasar

KDDLab, dept. of Informatics and Telecommunications, University of Athens, Greece

Motivation

We study the problem of event detection in a social stream



Events could be news, disasters, concerts, sports happenings, etc.

Our goals:

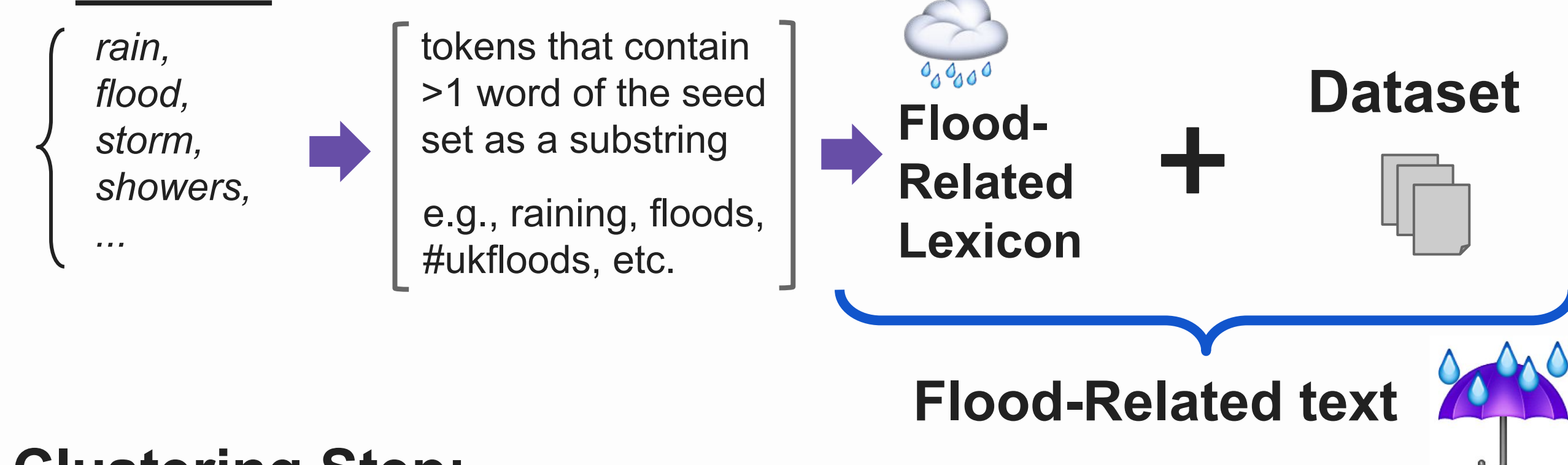
- Identify the event
- Monitor the *evolution*, the *duration* and the *location* of the event
- Inform users

Content-based event detection

Identifying topics in a stream of text

A) Filtering Step:

Create Flood-Related Lexicon & Extract Flood-Related text **seed set**



B) Clustering Step:

- Find areas

K-Means using the GPS coords → Clusters as Voronoi polygons

- Cluster areas

K-Means using:

1. **count**(*d*)

2. **ratio**(*d*) = $\text{count}(d) / \sum \text{count}(d')$, for all *d'*

3. **speed**(*d*) = $\text{ratio}(d) - \text{ratio}(d-1)$



K = 500

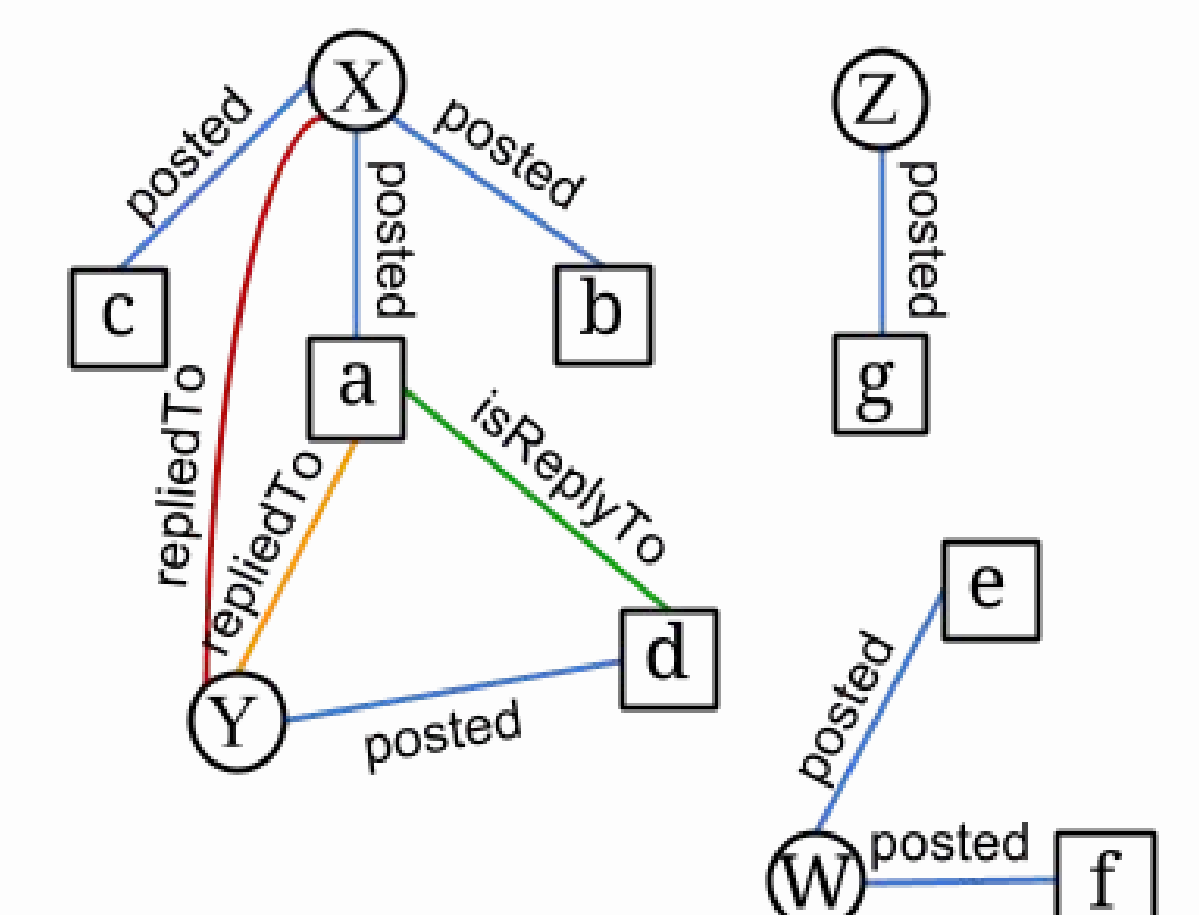
Graph-based event detection

Model network as a graph

A) Building Step:

Construct the network graph

- userX posted tweetA
- userY posted tweetD
- tweetD isReplyTo tweetA
- userY repliedTo userX
- userY repliedTo tweetA



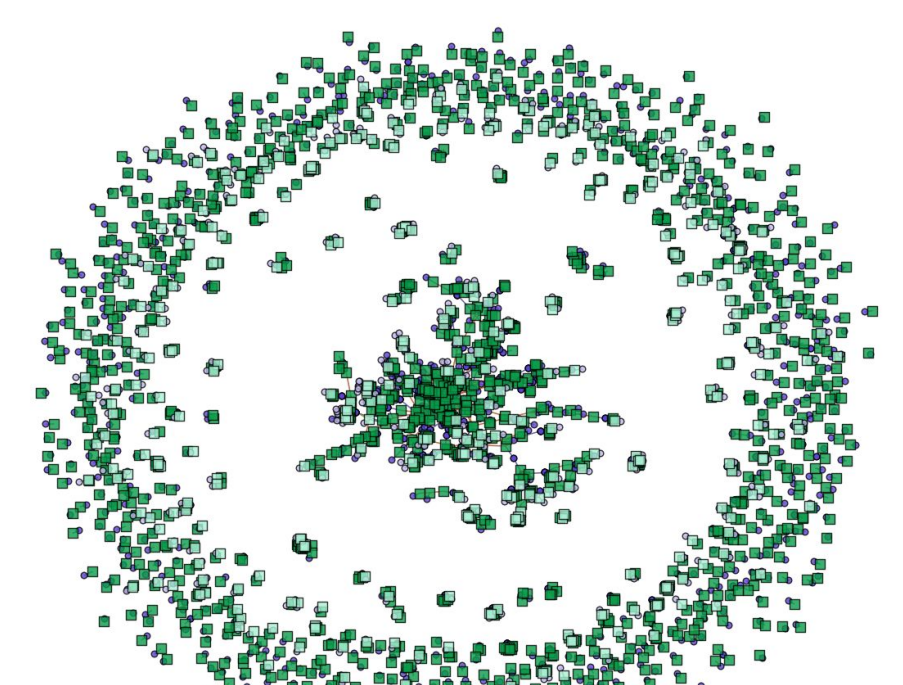
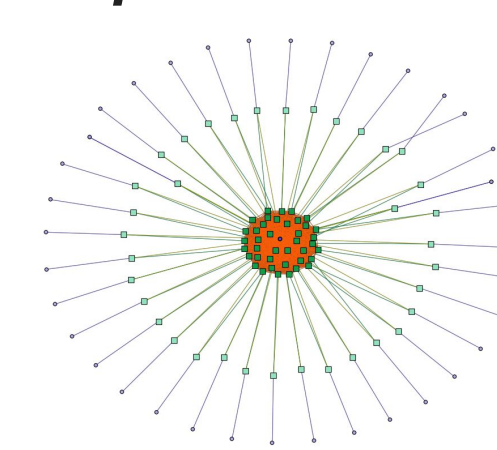
B) Filtering Step:

- **Large Connected Components** indicate large-scale conversations:

LCCs → event candidates

prune *spam* *LCCs* → star graph structure

e.g.,

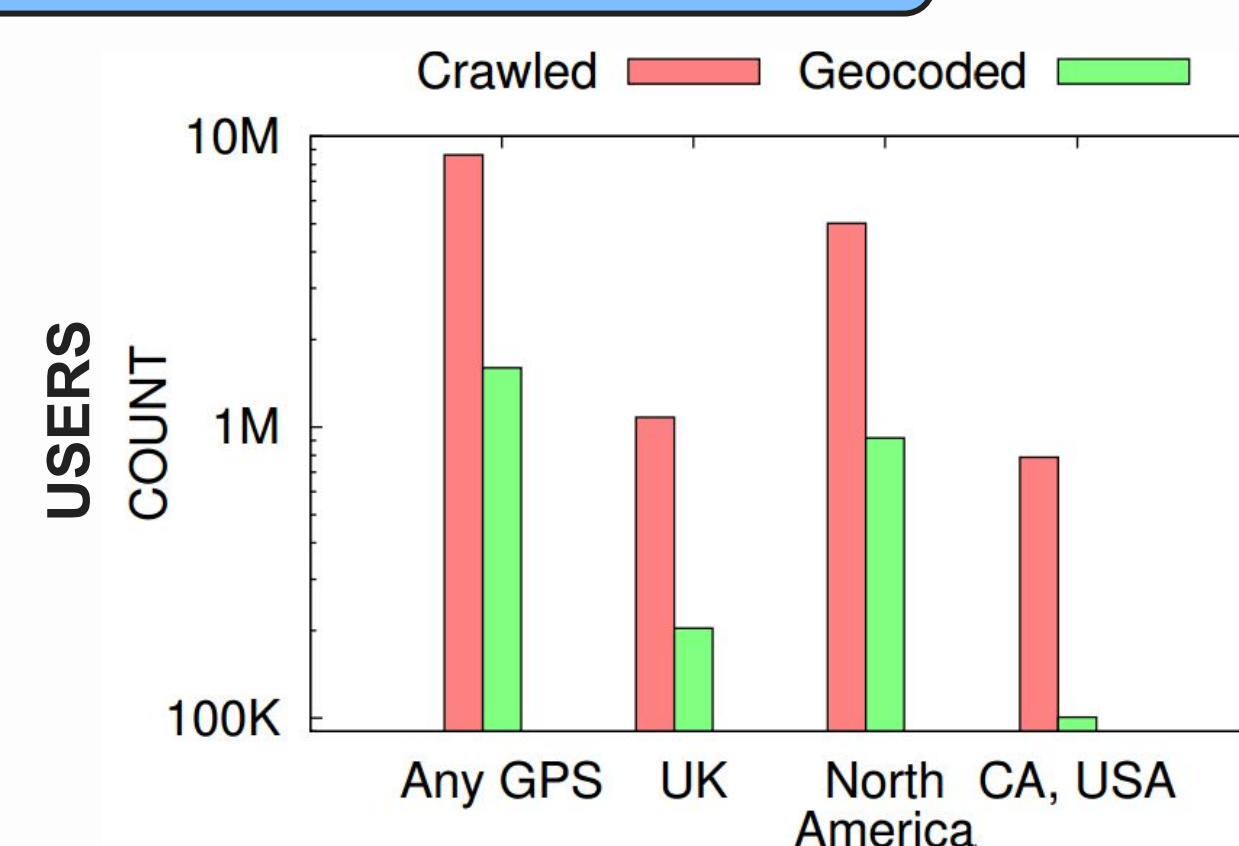


Facet Crawler

To collect data from Twitter, we use our custom made crawler

different queries → different constraints from Twitter API
e.g., limited number of queries per 15 min timecap

Our crawler can handle different query types with different parameters



example:

Crawl tweets with *location*:

- 2D bounding box with GPS (**green**), GPS enabled
- custom geocoding (**red**)
e.g., location in user profile

Geocoding can extract an additional **10%**

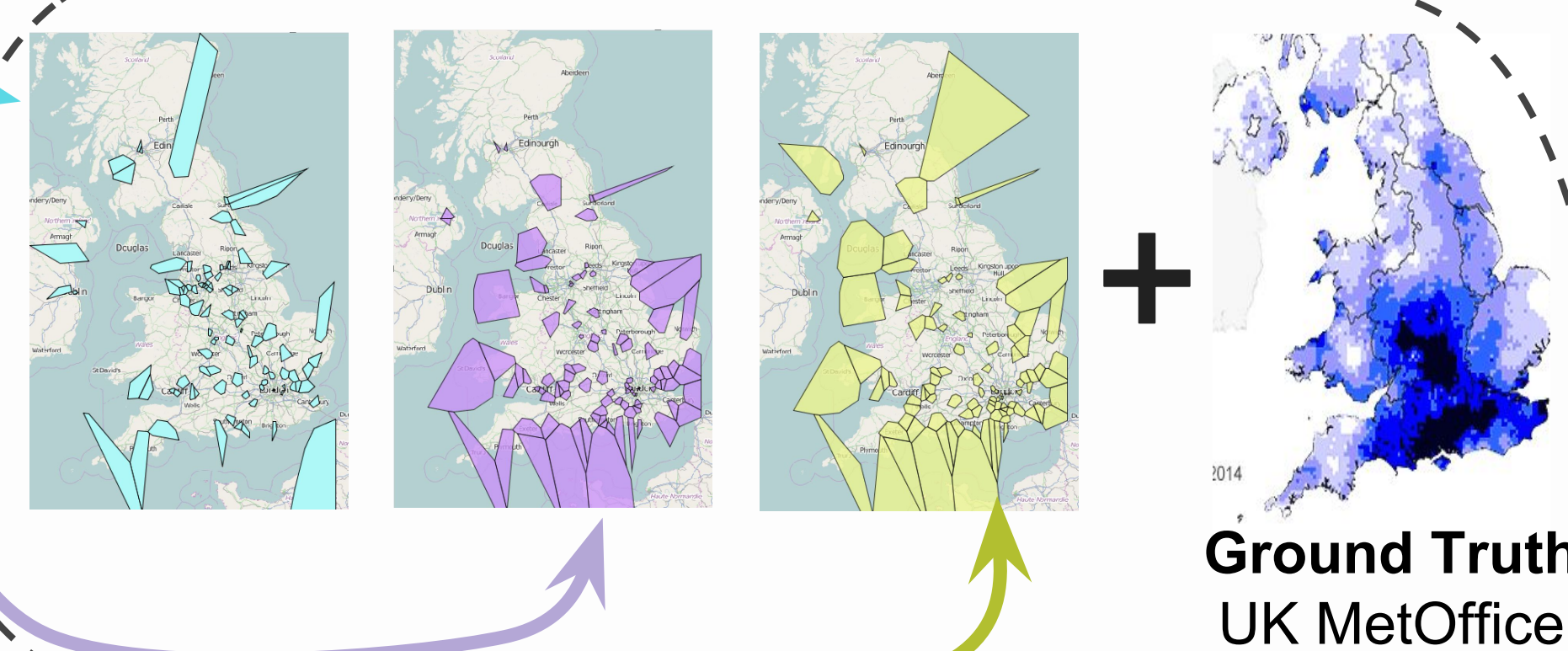
Experiments on content-based

- Collection of **public** tweets from the UK
- > 2.3M geotagged tweets

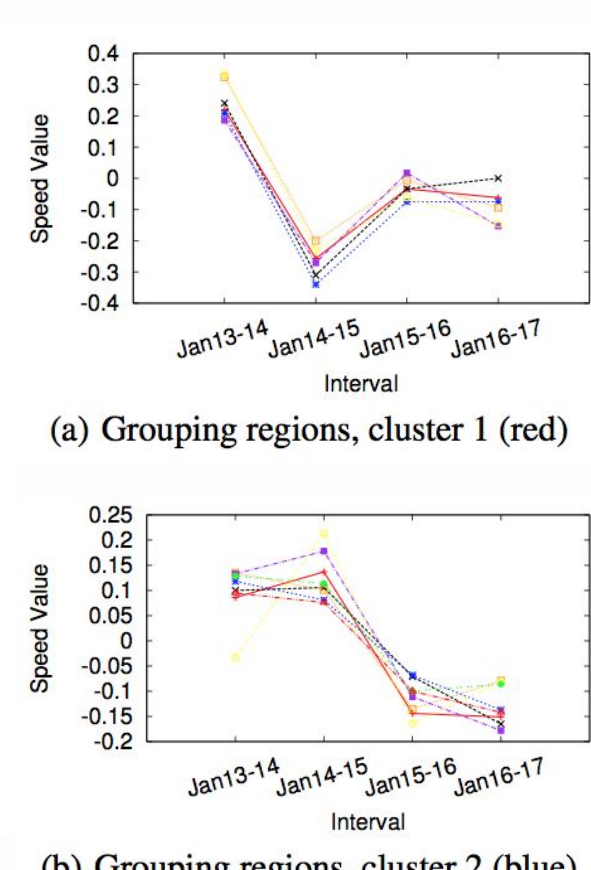
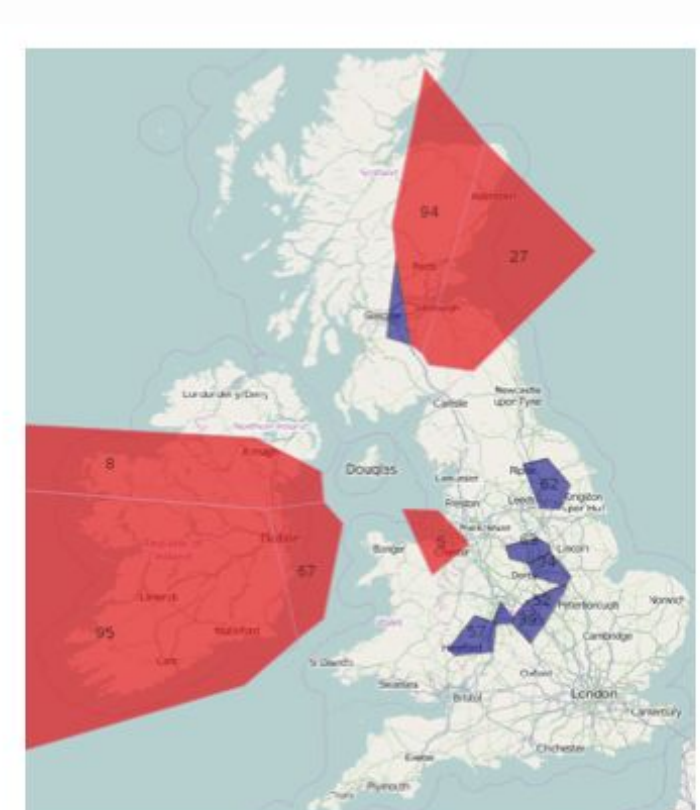
Top 100 areas based on:

- number of **all** tweets
- number of **flood-related** tweets
- signal-to-noise ratio**:

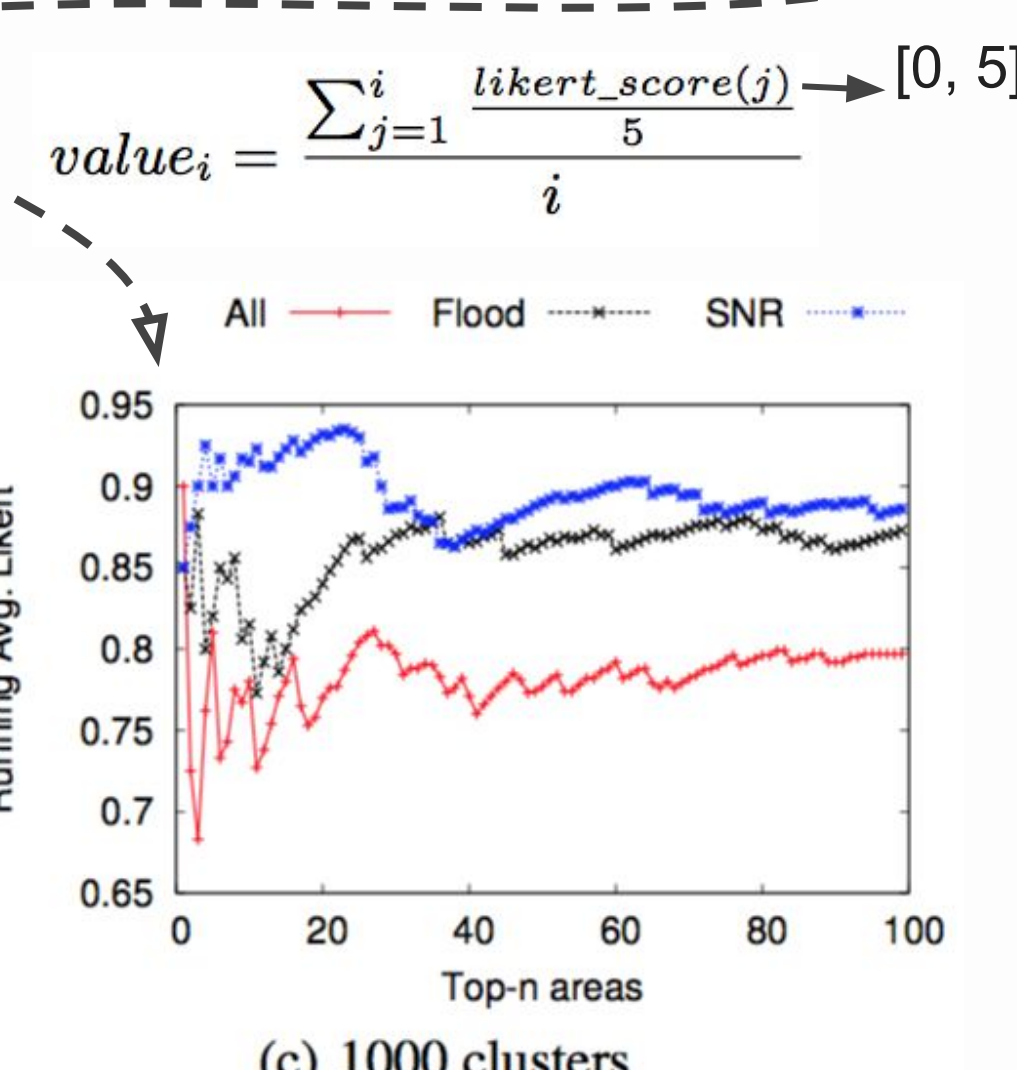
$$\text{score}(r) = \frac{\# \text{flood-related tweets in } r}{\# \text{tweets in } r}$$



Visualization of 2 clusters



- **Speed** feature has better performance
- **Red cluster**: mostly unaffected, *speed* decreases
- **Blue cluster**: affected, *speed* increases



Experiments on graph-based

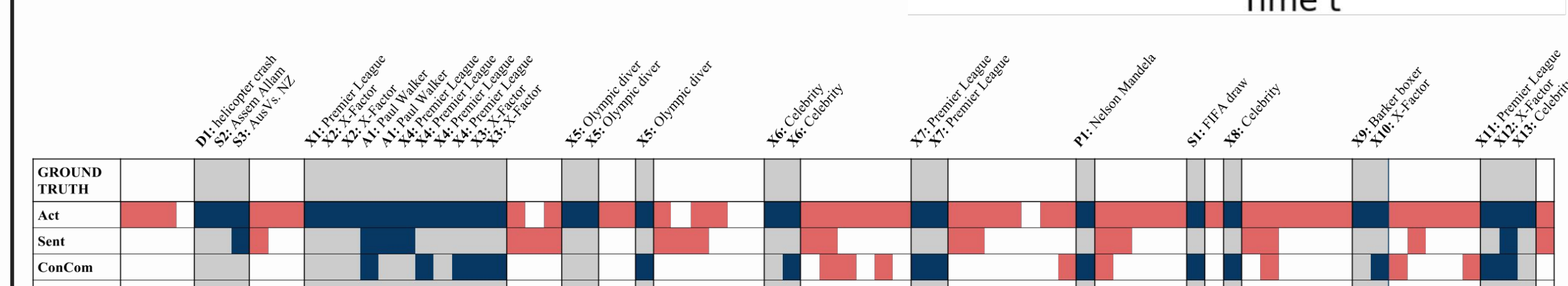
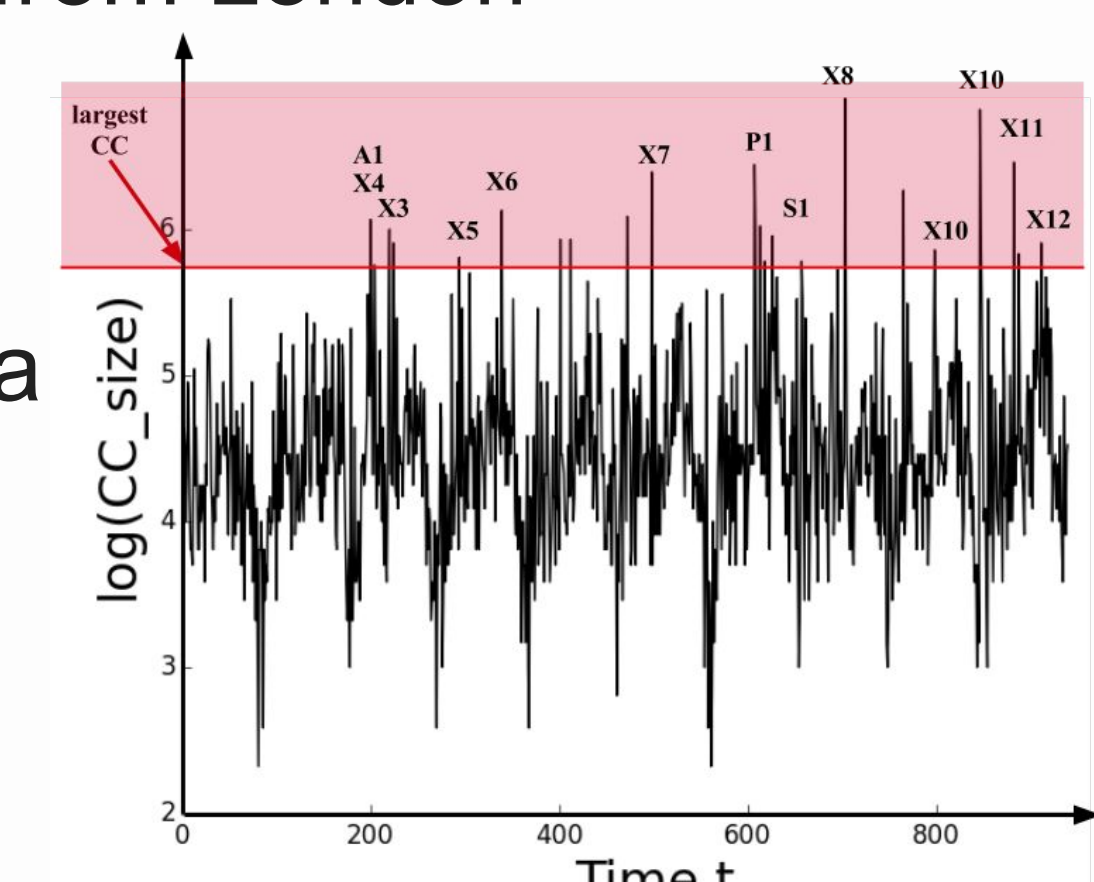
- Collection of **public** tweets from London
- ~ 700K geotagged tweets

ground truth:

- *automatically* from Wikipedia
- manual investigation

data into 15-min segments

12 out of 19 events ✓



Comparison:

- **Act**: number of tweets
- **Sent**: negative / positive
- **ConCom**: our method

	jaccard similarity	precision	recall	f-score
Act	0.5	0.42	0.86	0.56
Sent	0.47	0.23	0.17	0.2
ConCom	0.72	0.65	0.52	0.58