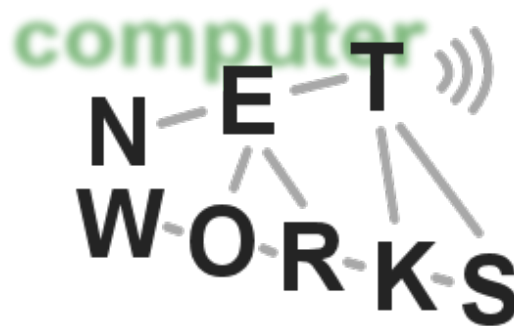


# Social Networks: Applications

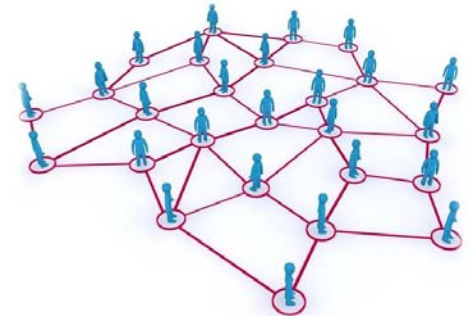
Advanced Computer Networks  
Summer Semester 2012



# Decentralized Data Routing in Mobile Social Networks

# Mobile Social Network (MSN)

- MSN:
  - A virtual network consisting a set of mobile individuals
  - Decentralized:
    - No central component, no radio tower
    - Only local information is known for each node
  - Communication:
    - Two nodes can communicate only when they move into each other's communication range
    - Multi-hop: message is forwarded to the destination via multiple relays
- It can be modeled as a weighted social graph.
  - Mobile devices -> nodes
  - Encounters of nodes -> edges
  - Weight -> number of encounters



# Data Routing in MSNs

- Network topology is dynamic due to mobility
  - The connections between nodes are intermittent
  - No fixed end-to-end path (routing table doesn't work!)
  - Data routing is in a store-carry-and-forward manner
- Information can be used for routing:
  - Social, Location, Mobility
- **Question: how to devise decentralized routing scheme?**

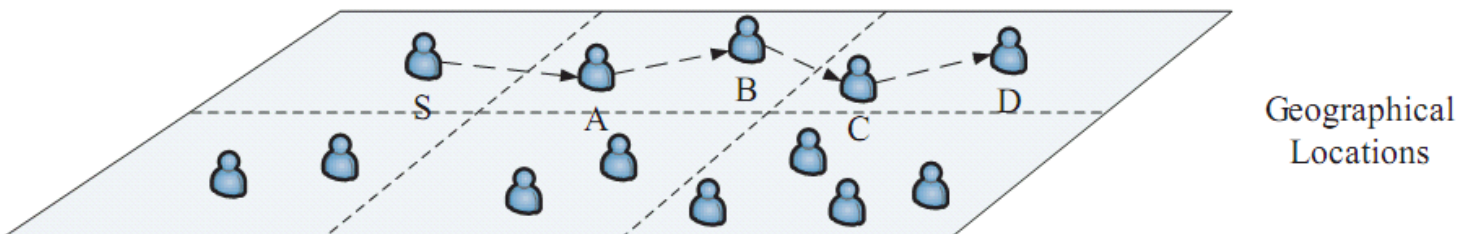


# Location Information

- The **location information** is described as the geographical positions of nodes.
  - A sequence of time-varying events like:  
node ID, location ID, start time and time duration
- Given the basic location information, we can obtain
  - Geographical distance
  - Statistical location information
    - Time proportion that a node stays at a location

# Geographical Distance Based Strategy

- **Geographical Distance**: the Euclidean distance between two nodes
  - Given the position of  $r_i$  and  $r_j$ ,  $d_{ij} = ||r_i - r_j||$
- **Strategy**: choose the node having closer distance with destination as the relay
  - Example: GPRS routing [1]
- Advantages and disadvantages
  - **Pros**: it can geographically approach the destination
  - **Cons**:
    - Geographical information need dedicate device like GPS
    - The position of destination is hard to know

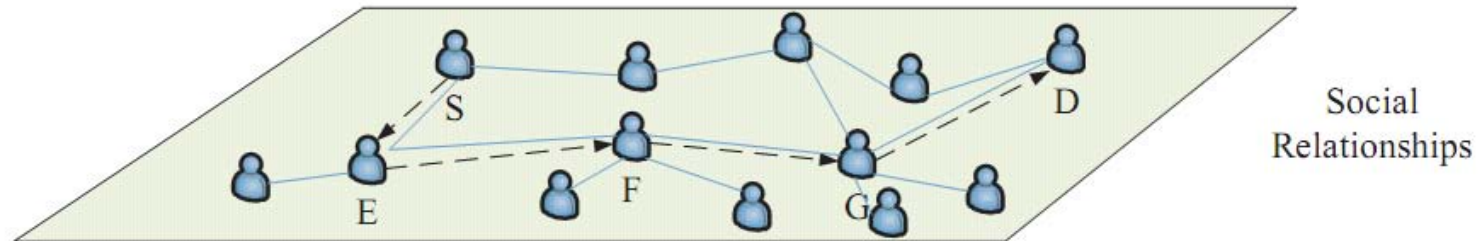


# Social Information

- **Social information** in a MSN indicates the structural status of a node in social graph
  - **Social degree**: the number of friends
  - **Weight of an edge**: the strength of the social tie between a pair of nodes
  - **Common friends**: the friends two nodes share, which describes the similarity of them in social structure.

# Social Centrality Based Strategy

- Social centrality
  - The quantification of the relative importance of nodes in the social network
  - Many definitions, can be simply defined as the node degree (the number of friends)
- Strategy: forward data to nodes with higher social centrality
  - Example: BubbleRap[2]
- Advantages and disadvantages
  - Pros: simple, decentralized, easy to implement
  - Cons: could meet a dead end (i.e. the node with highest centrality value does not connected to the destination)



[2] P. Hui, J. Crowcroft, and E. Yoneki. Bubble rap: social-based forwarding in delay tolerant networks. MobiHoc '08, pages 241–250, New York, NY, USA, 2008. ACM.



# Social Similarity Based Strategy

- Social Similarity

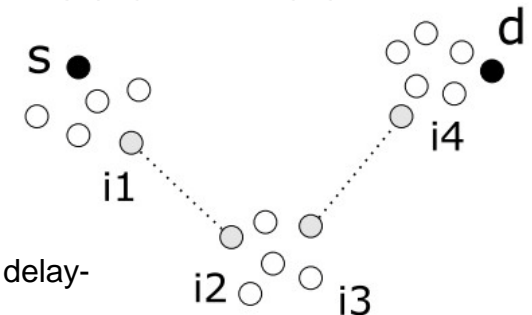
- Indicates the trustiness and cohesive of social links
- Evaluates by the number of common friends of two nodes  
$$S_{ij} = |F_i \cap F_j|$$

- Strategy: forward data to the nodes with higher social similarity to the destination

- Example: SimBet[3]

- Advantages and disadvantages

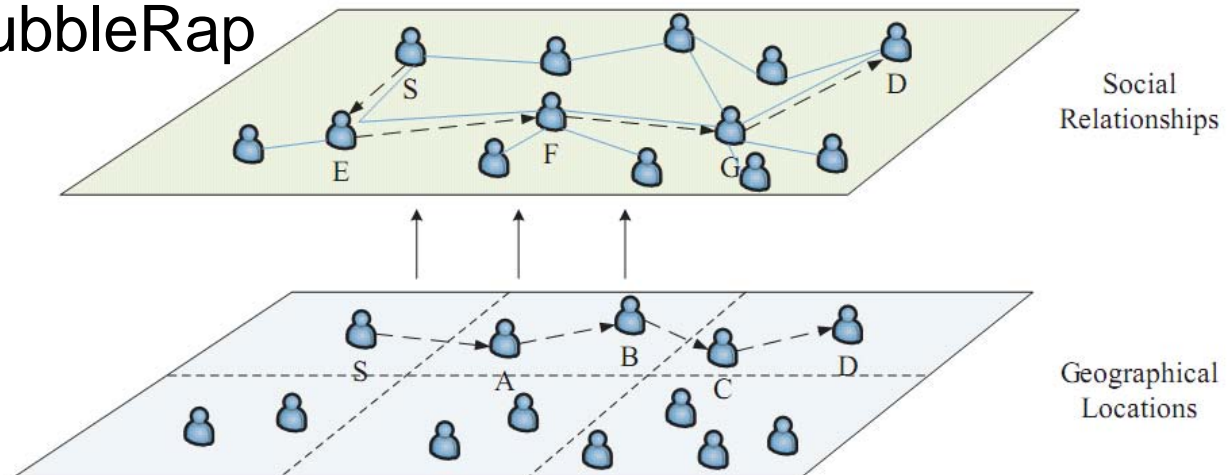
- Pros: reduce dead end
- Cons: need to exchange friends list of each node pairs



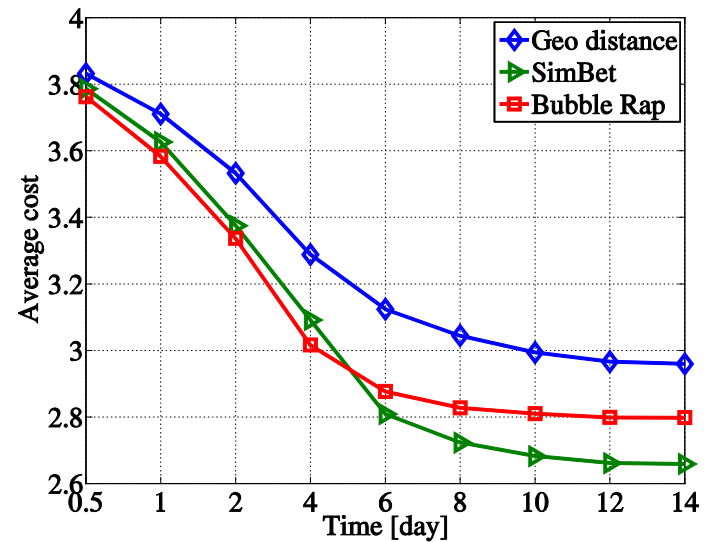
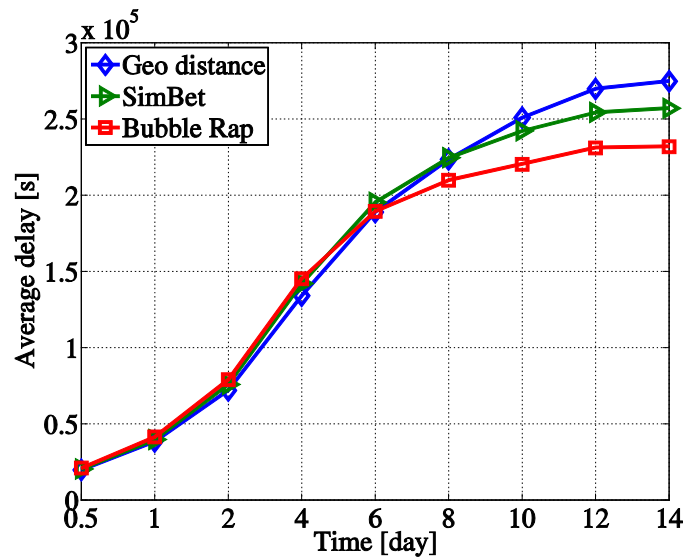
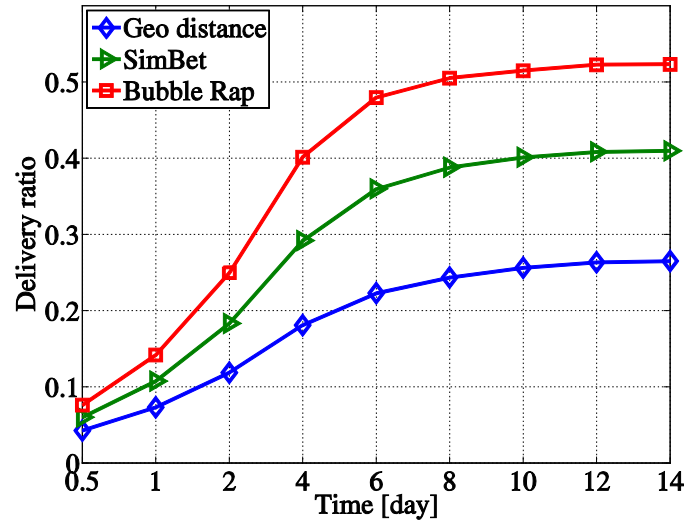
[3] E. M. Daly and M. Haahr. Social network analysis for routing in disconnected delay-tolerant manets. MobiHoc '07, pages 32–40, New York, NY, USA, 2007.

# Comparison: Strategies

- Two level of information: location info and social info
- Location-based
  - Distance: GPSR
- Social-based
  - Similarity: SimBet
  - Centrality: BubbleRap



# Comparison: Performance



# Decentralized Routing in MSNs: Summary

- Both location information and social information can be used to design decentralized routing strategies in MSNs
- Simple location based strategies **do not perform better** than social based strategies in delivery ratio, delay and path length
- The power of social links:
  - Location information is expensive and highly private
  - Social information is public available and easy to use
- Choose the **appropriate** strategy based on real applications

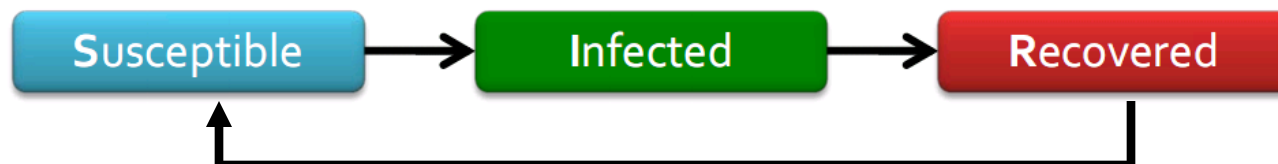
# Epidemics

# Spread of Contagious Diseases

- Spread of contagious diseases
  - Can pass explosively through a population
  - Determined by the properties of the virus: including its contagiousness, the length of its infectious period, and its severity
  - Also affected by network structures within the population it is affecting
- The spread of computer viruses
- Diffusion of ideas through social networks

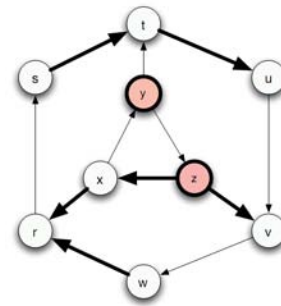
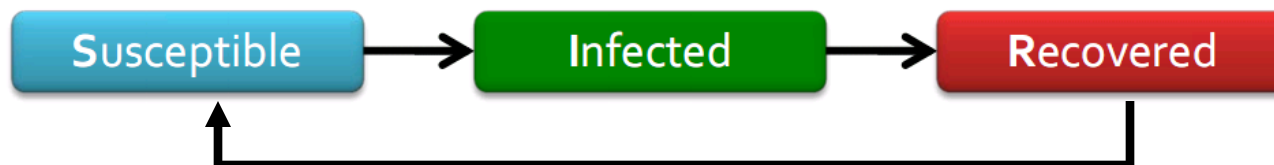
# The SIRS Epidemic Model

- Each node goes through the following potential stages: S-I-R-S
  - **Susceptible**: Before the node has caught the disease, it is susceptible to infection from its neighbors.
  - **Infectious**: Once the node has caught the disease, it is infectious and has some probability of infecting each of its susceptible neighbors.
  - **Removed**: After a particular node has experienced the full infectious period, this node is removed from consideration, since it no longer poses a threat of future infection.
  - **Susceptible**: after the removed stage, it returns to the Susceptible stage



## ○ Process

- Initially, some nodes are in the **I** state and all others are in the **S** state.
- Each node  $v$  that enters the **I** state remains infectious for a fixed number of steps  $t_I$ .
- During each of these  $t_I$  steps,  $v$  has a probability  $p$  of passing the disease to each of its susceptible neighbors.
- After  $t_I$  steps, node  $v$  is no longer infectious. It then enters the **R** state for a fixed number of steps  $t_R$ . During this time, it cannot be infected with the disease, nor does it transmit the disease to other nodes.
- After  $t_R$  steps in the **R** state, node  $v$  returns to the **S** state.

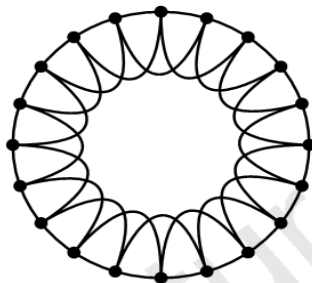




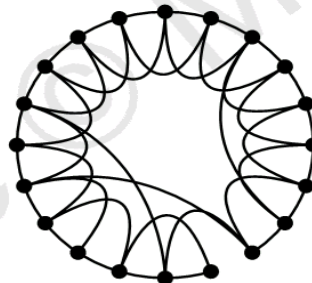
# Small-World Contact Networks

- **Recap: The 1-dimensional Watts-Strogatz Model**
- Starting from a ring lattice with  $n$  vertices and  $k$  edges per vertex.
  - Regular network with high clustering coefficient
- We rewire each edge at random with probability  $p$  ( $0 \leq p \leq 1$ ).
  - $p=0$ : regular network
  - $p=1$ : random network
  - $0 < p < 1$ : small world network

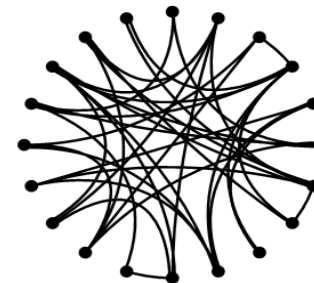
Regular



Small-world

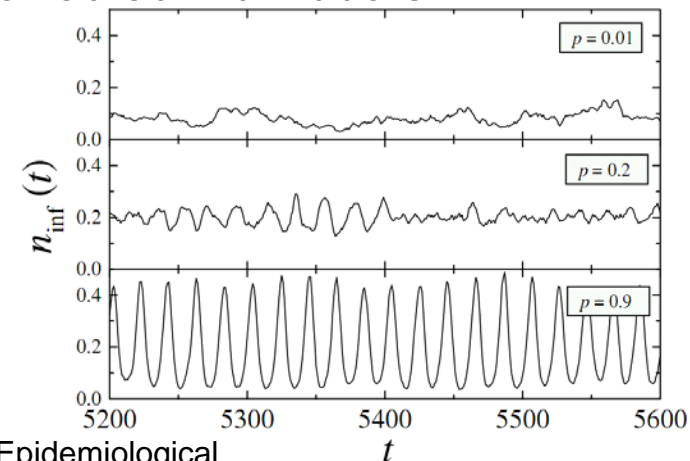
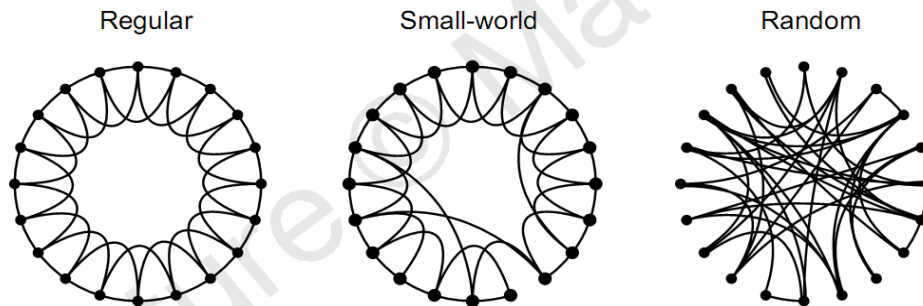


Random



# Small World Effect in the SIRS Model [4]

- Different behavior is observed depending on the value of  $p$ 
  - When  $p$  is small ( $p=0.01$ )
    - Disease transmission through the network occurs mainly via the **short-range local edges**
    - Flare-ups of the disease in one part of the network **never become coordinated** with other parts
  - When  $p$  increases ( $p=0.2$ )
    - These flare-ups start to **synchronize**
    - Oscillations intermittently appear and then disappear
  - For very large values of  $p$  ( $p=0.9$ )
    - There are **clear waves** in the number of affected individuals



[4] Marcelo Kuperman and Guillermo Abramson, Small World Effect in an Epidemiological

Model, PHYSICAL REVIEW LETTERS, Vol. 86, No. 13, 2001, 2909-2012.

FIG. 1. Fraction of infected elements as a function of time.

# Case Study: Tracking Flu Using Twitter [5]

- Collecting information about epidemics:
  - The **location, timing and intensity** of an epidemic
  - Information is collected from school and workforce absenteeism figures, phone calls and visits to doctors and hospitals
- Gathering this information is a **difficult, resource-demanding, time-consuming** procedure
- Use of search engine data to detect Influenza-like Illness (ILI)
  - Geographic clusters with a heightened proportion of health-related queries
- **Using Twitter to detect ILI?**

# Data

- Twitter, UK
  - Daily average of 160,000 tweets
  - 24 weeks from 06/22/2009 to 12/06/2009
  - Twitter geolocation (geographical coordinates).
- Official health reports
  - Health Protection Agency (HPA), UK.
  - Region A = Central England & Wales
  - Region B = South England
  - Region C = North England
  - Region D = England & Wales
  - Region E = Wales & Northern Ireland

# Official Health Reports

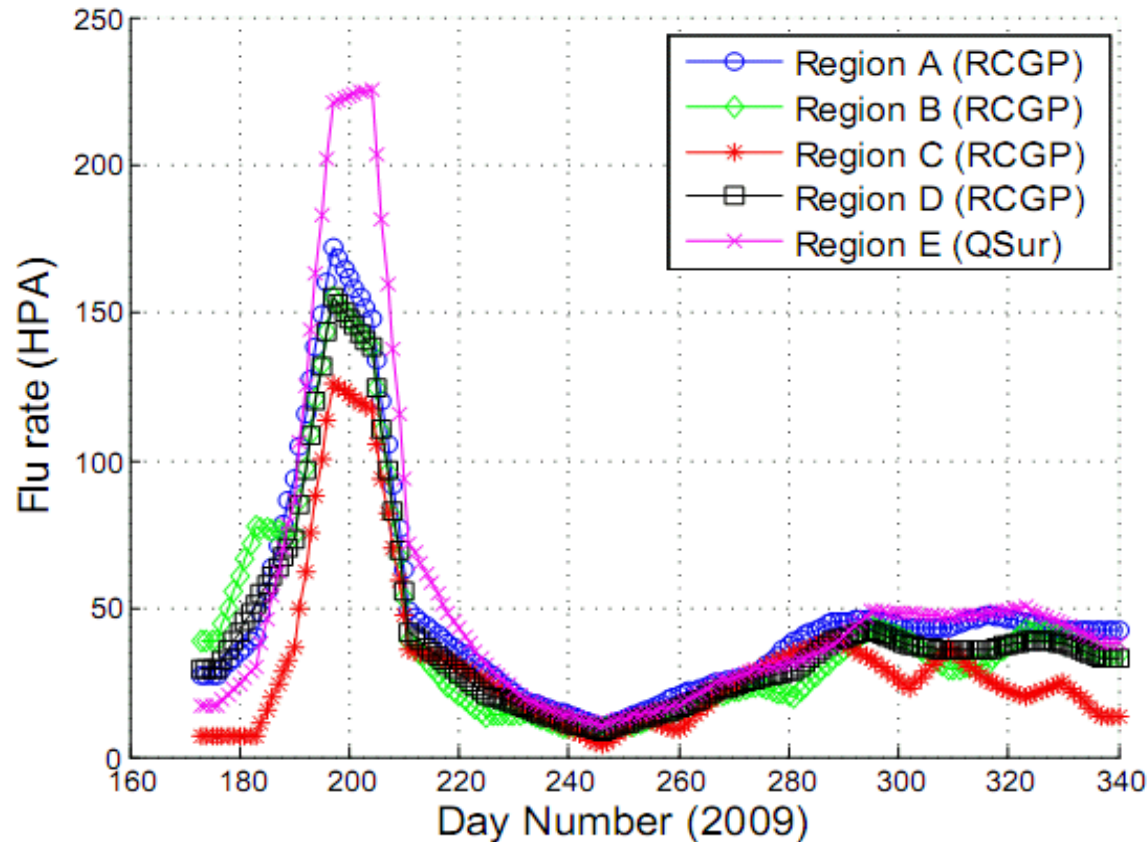
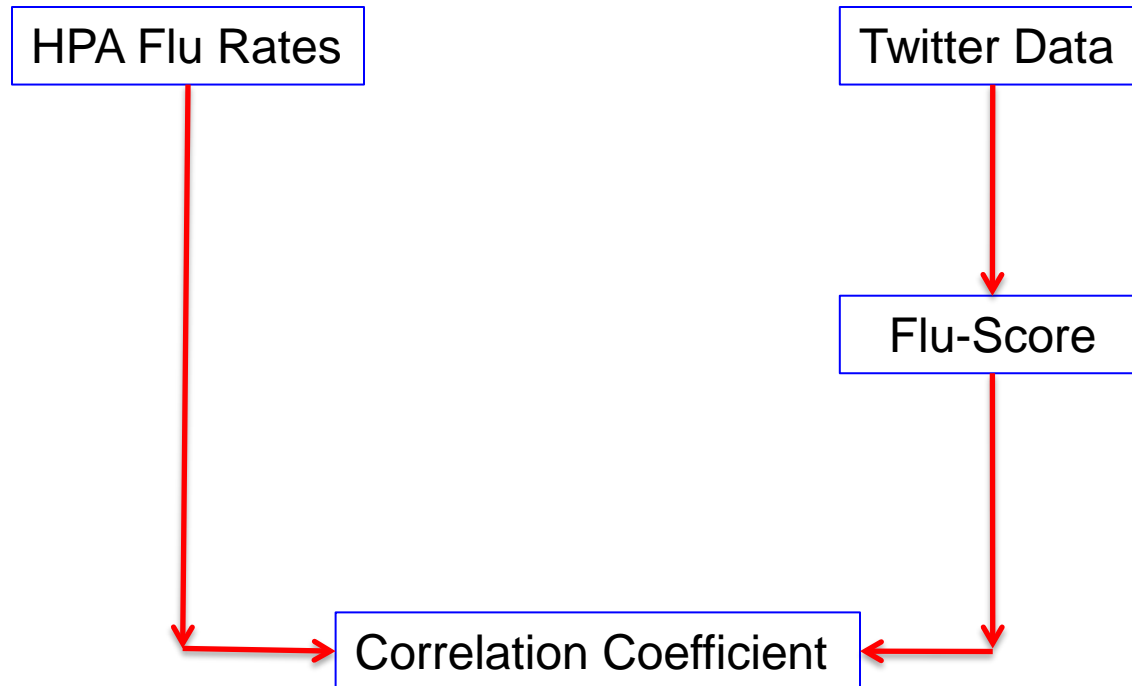


Fig. 1: Flu rates from the Health Protection Agency (HPA) for regions A-E (weeks 26-49, 2009). The original weekly HPA's flu rates have been expanded and smoothed in order to match with the daily data stream of Twitter (see section III-B).

# Methodology



# Computing Flu-scores

- The **daily set** of Tweets:

$$\mathcal{T} = \{t_j\}, \text{ where } j \in [1, n]$$

- **Textual markers**: expressing illness symptoms, e.g. fever, temperature, sore throat, infection, headache

- A **set of textual markers**:  $\mathcal{M} = \{m_i\} \quad i \in [1, k]$

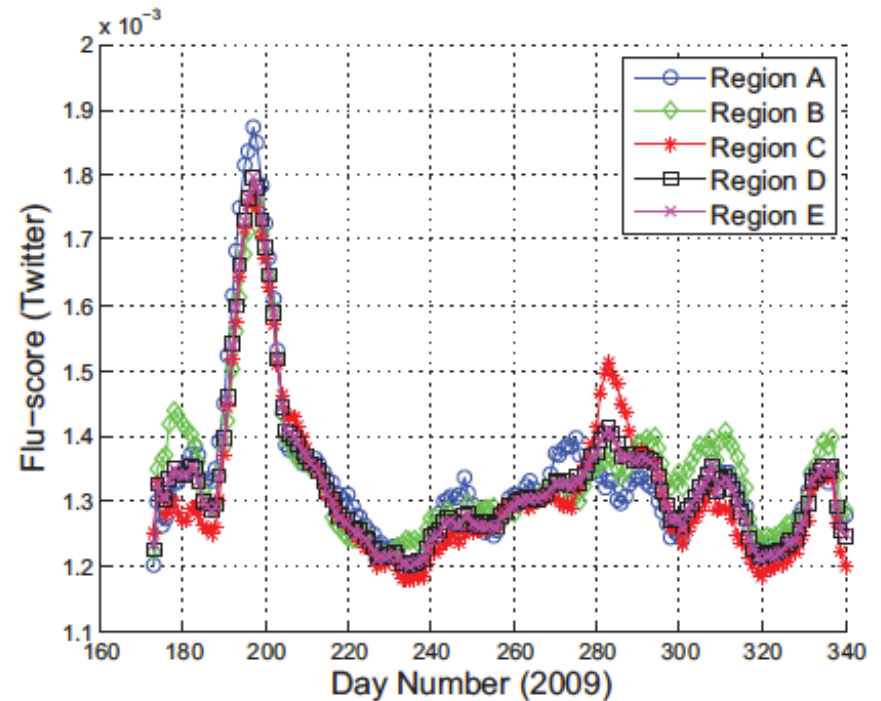
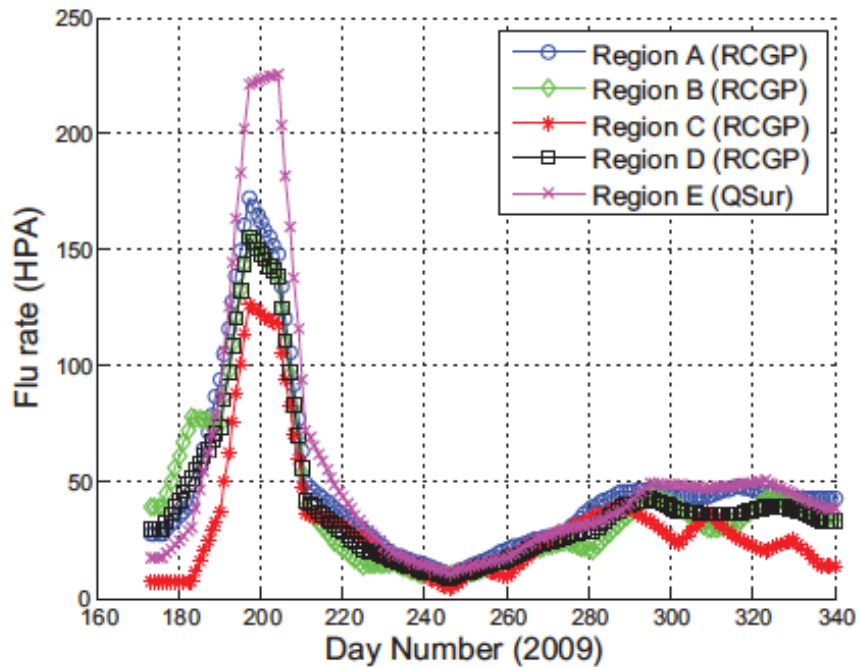
- Let  $m_i(t_j)=1$  if  $m_i$  appears in tweet  $t_j$ , otherwise=0

- The flu-score of a tweet  $t_j$ :  $s(t_j) = \frac{\sum_i m_i(t_j)}{k}$

- The flu-score of one day Twitter corpus

$$f(T, M) = \frac{\sum_j s(t_j)}{n} = \frac{\sum_j \sum_i m_i(t_j)}{k \bullet n}$$

# Flu-rate vs Flu-score





# Correlation Coefficient

- A measure of the correlation (linear dependence) between two variables X and Y
- Giving a value between +1 and -1 inclusive

- Definition  $\rho_{X,Y} = \text{corr}(X, Y) = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{E[(X - \mu_X)(Y - \mu_Y)]}{\sigma_X \sigma_Y}$ ,

- For a sample:

$$r_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{(n - 1)s_x s_y} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}}$$

Correlation	Negative	Positive
None	-0.09 to 0.0	0.0 to 0.09
Small	-0.3 to -0.1	0.1 to 0.3
Medium	-0.5 to -0.3	0.3 to 0.5
Strong	-1.0 to -0.5	0.5 to 1.0

# Correlations between Twitter Flu-scores and HPA Flu rates

- Strong correlation is observed!
- It indicates linear correlation between Twitter flu-scores and HPA flu rates, thus the flu-scores can be used to predict flu rates!

Region	HPA Scheme	Corr. Coef.
A	RCGP	0.8471
B	RCGP	0.8293
C	RCGP	0.8438
D	RCGP	0.8556
E	QSur	0.8178

# Extensions

- Learning HPA's flu rates from Twitter flu-scores
  - Linear regression is used to build a weighted Twitter flu-scores to model flu rates
- Automatic extraction of ILI textual markers
  - Creating candidate markers from:
    - Encyclopedic reference
    - Informal references
  - Forming the flu-subscores with time series.

# Tracking Flu: Summary

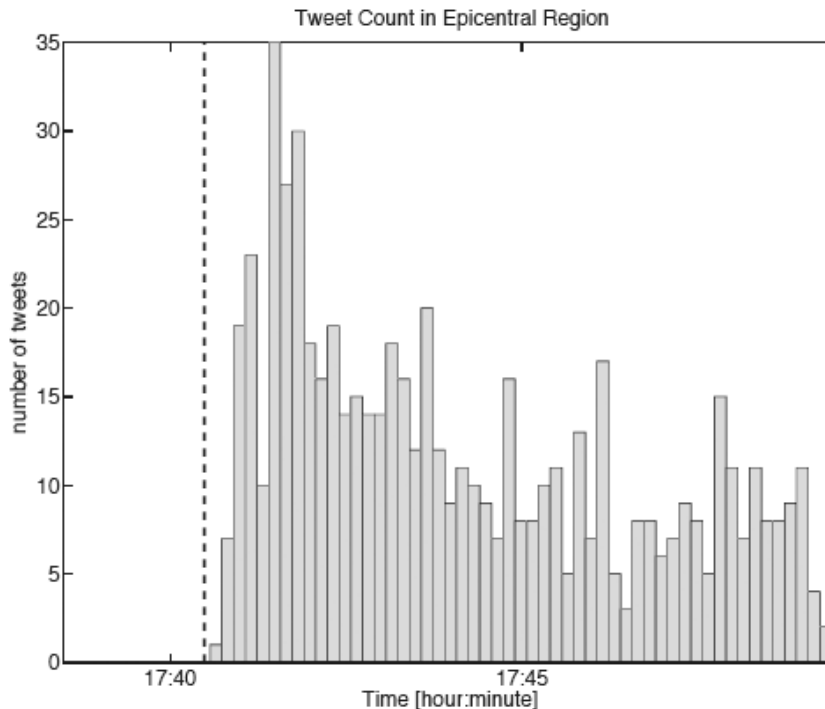
- Tracking the flu outbreak in the UK using Twitter messages.
- High correlation between the flu-score and the HPA flu rates, greater than 95%.
- Advantages:
  - Less resource-demanding: only monitoring Twitter website, can be done automatically
  - More faster: can be done efficiently, while official reports need to delay 1 or two weeks
- Disadvantages
  - Still need sample data from official statistics for learning
  - Not everyone post their disease: could be not accurate
  - Not suitable for all sort of contagious diseases: some of them are not discussed in Twitter for privacy reason

# Tracking Earthquake Using Twitter [6,7]

[6] Paul Earle, Michelle Guy, Richard Buckmaster, Chris Ostrum, Scott Horvath, and Amy Vaughan, OMG Earthquake! Can Twitter Improve Earthquake Response? U.S. Geological Survey, 2010

[7] Takeshi Sakaki, Makoto Okazaki, and Yutaka Matsuo. 2010. Earthquake shakes Twitter users: real-time event detection by social sensors. In Proceedings of the 19th international conference on World wide web (WWW '10). ACM, New York, NY, USA, 851-860.

- Subsequent earthquakes generated volumes of earthquake-related tweets
- Access to firsthand accounts of earthquake shaking within seconds of an earthquake is intriguing



Tweet count following the 2009 MW 4.3 Morgan Hill, CA, earthquake. Tweets are binned in 10-second intervals, and the dashed line marks the origin time of the earthquake. After the earthquake, the tweet frequency quickly rose above the background level of less than one per hour to about 150 per minute.

# Using Twitter for Earthquake Detection and Assessment

- **Real-time** nature of Microblogging
  - We can know what happens around other users in realtime
- **Public tweets** are stored in an openly searchable database
- **Earthquake Reports**
  - U.S. response time: **1.5-20 min**
- Earthquake detection using Twitter
  - The typical delay for tweet transmission is **5 seconds**
  - Earthquake could be detected in under **a minute**
- Twitter could be faster than Earthquake Wave!
  - An earthquake propagates at about **3–7 km/s** (for 100km, about **20s**)
  - **Early alarm is possible?**

# Mapping the Felt Region

- Morgan Hill, 30 March 2009, MW 4.3

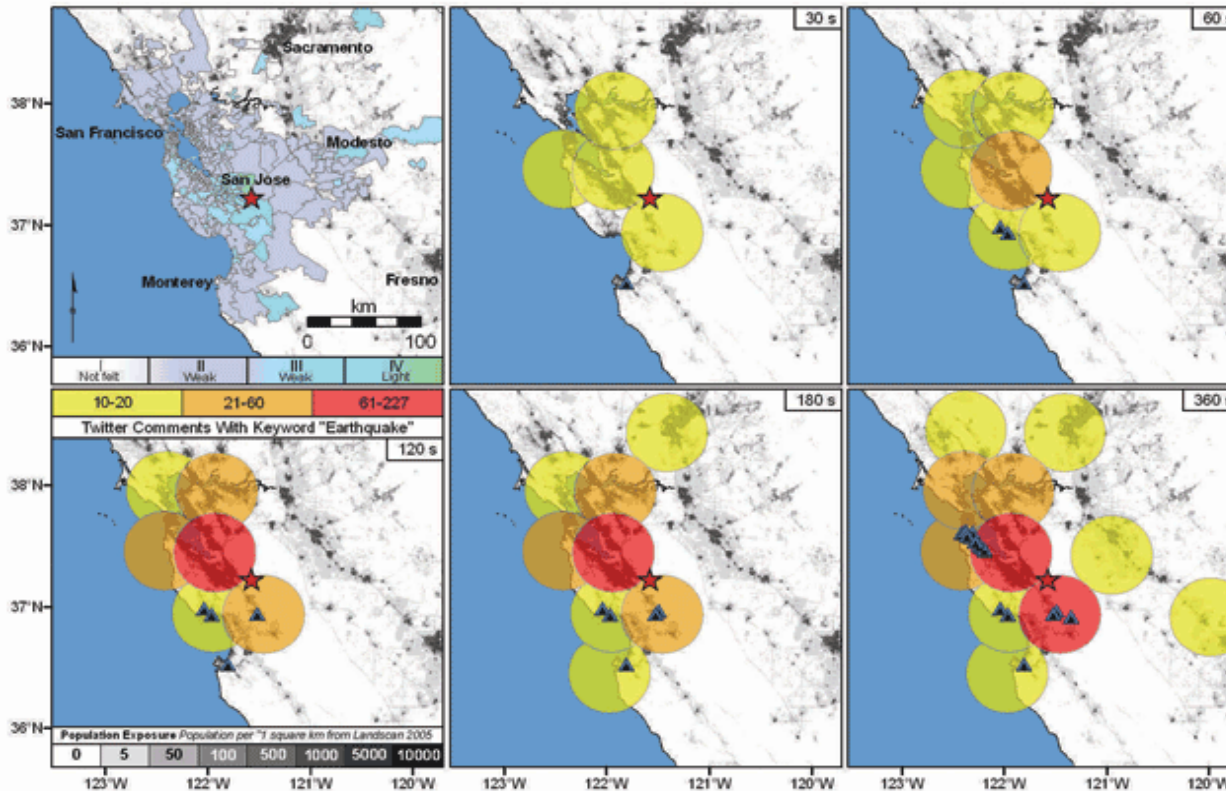
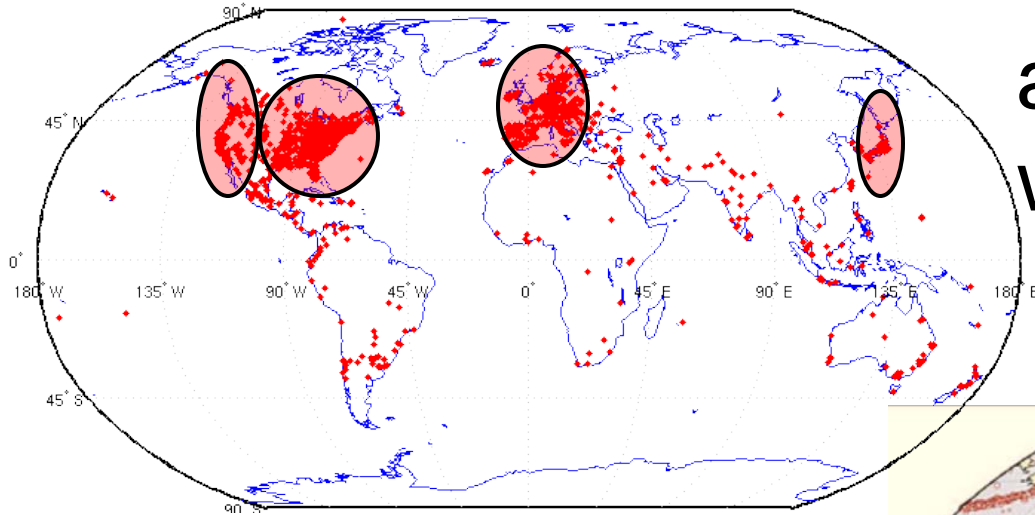


Figure 2. Comparison of the intensity map produced by the USGS DYFI system (upper left) and the geocoded tweet counts for the 30 March 2009 MW 4.3 Morgan Hill earthquake. The extent of the geographic tweet search is indicated by the size of the circles, which are color coded by the number of tweets. The population is shown in the background as gray scale and the tweets with exact latitude and longitude geo-references are shown as black triangles with blue outlines. The different panels show the integrated tweet count at discrete times after the earthquake as indicated in the upper right corner of the maps.

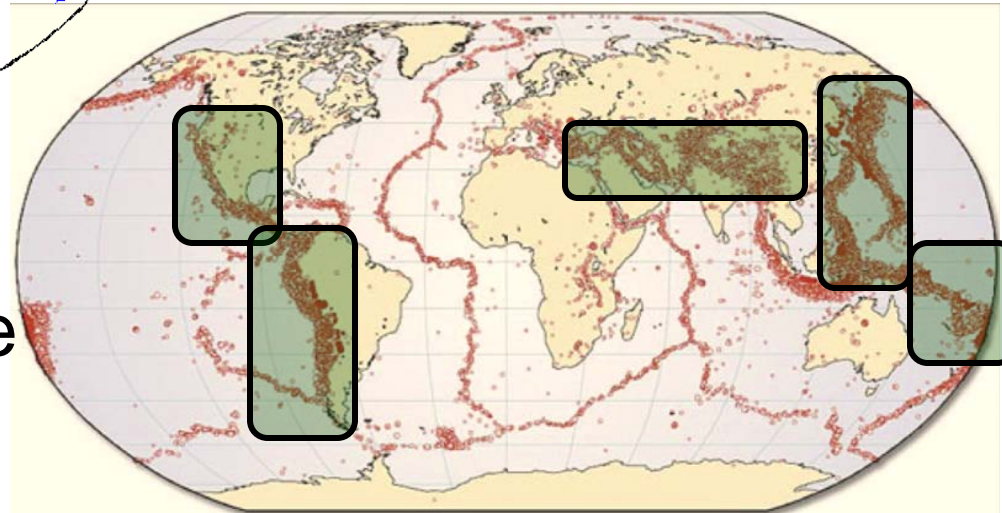


# Twitter and Earthquakes in Japan



a map of Twitter user world wide

a map of earthquake occurrences world wide



The intersection is regions with many earthquakes and large twitter users in Japan.

Other regions:

Indonesia, Turkey, Iran, Italy, and Pacific coastal US cities

# Event Detection Using Twitter

- Do semantic analysis on Tweet
  - To obtain tweets on the target event precisely
- Regard Twitter user as a sensor
  - To detect the target event
  - To estimate location of the target

# Semantic Analysis on Tweet

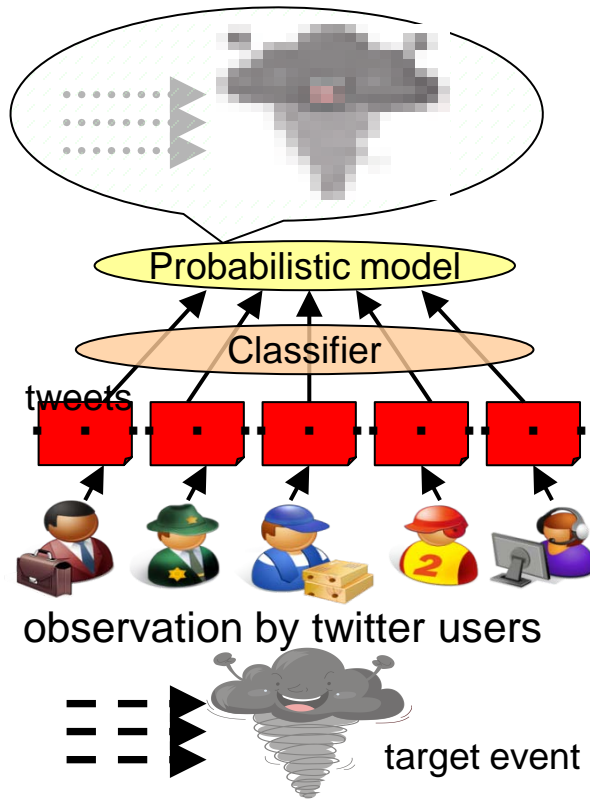
- ▶ Search tweets including keywords related to a target event
  - ▶ Example: In the case of earthquakes
    - ▶ “shaking”, “earthquake”
- ▶ Classify tweets into a positive class or a negative class
  - ▶ Example:
    - ▶ “Earthquake right now!!” ---positive
    - ▶ “Someone is shaking hands with my boss” --- negative
  - ▶ Create a classifier

# Semantic Analysis on Tweet

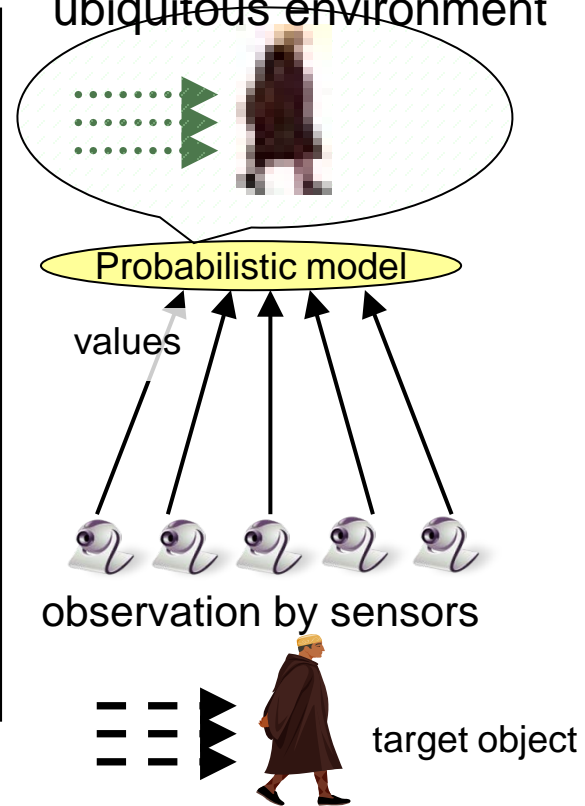
- ▶ Create classifier for tweets
  - ▶ use **Support Vector Machine(SVM)** - a machine learning algorithm
- ▶ Features (Example: I am in Japan, earthquake right now!)
  - ▶ **Statistical features** (7 words, the 5<sup>th</sup> word)  
the number of words in a tweet message and the position of the query within a tweet
  - ▶ **Keyword features** ( I, am, in, Japan, earthquake, right, now)  
the words in a tweet
  - ▶ **Word context features** (Japan, right)  
the words before and after the query word

# Tweet as a Sensory Value

Event detection from twitter

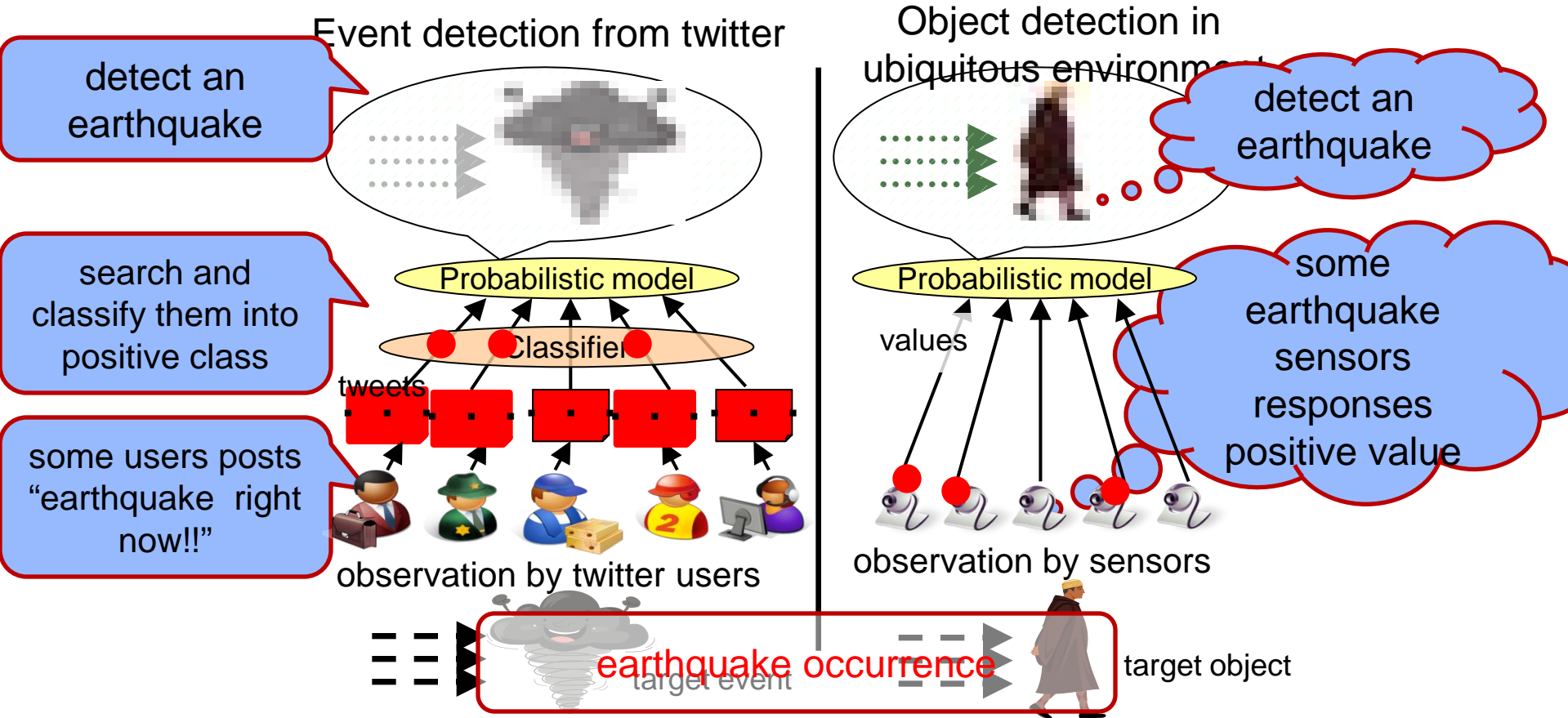


Object detection in ubiquitous environment



the correspondence between **tweets processing** and **sensory data detection**

# Tweet as a Sensory Value



We can apply methods for sensory data detection to tweets processing

# Tweet as a Sensory Value

- ▶ We make two assumptions to apply methods for observation by sensors
- ▶ Assumption 1: Each Twitter user is regarded as a sensor
  - ▶ a tweet → a sensor reading
  - ▶ a sensor detects a target event and makes a report probabilistically
  - ▶ Example:
    - ▶ make a tweet about an earthquake occurrence
    - ▶ “earthquake sensor” return a positive value
- ▶ Assumption 2: Each tweet is associated with a time and location
  - ▶ a time : post time
  - ▶ location : GPS data or location information in user’s profile

Processing time information and location information, we can detect target events and estimate location of target events

# Earthquake Location Estimation

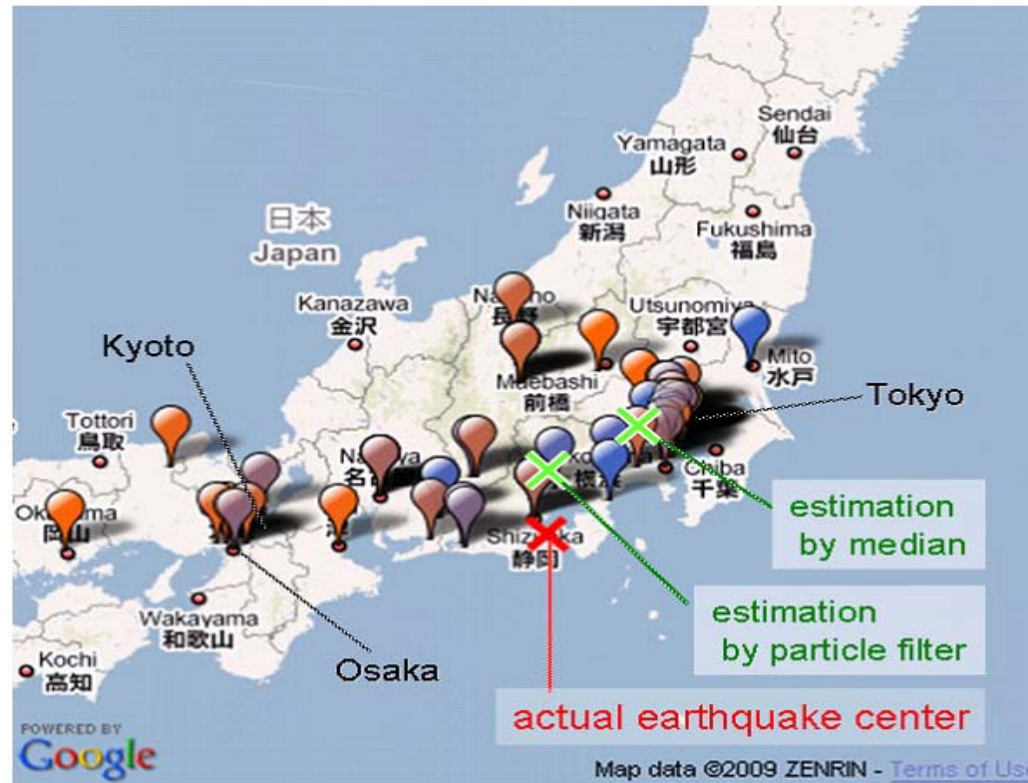


Figure 9: Earthquake location estimation based on tweets. Balloons show the tweets on the earthquake. The cross shows the earthquake center. Red represents early tweets; blue represents later tweets.



# Earthquake Reporting System

- Toretter ( <http://toretter.com>)
  - Earthquake reporting system using the event detection algorithm
  - All users can see the detection of past earthquakes
  - Registered users can receive e-mails of earthquake detection reports

Published	Location	Title	Screen_name	URL
2009-08-11 05:09:57	Saitama, Japan	地震おおいわー	tondol	<a href="http://twitter.com/tondol">http://twitter.com/tondol</a>
2009-08-11 05:08:56	unknown	地震。	tr0ly	<a href="http://twitter.com/tr0ly">http://twitter.com/tr0ly</a>
2009-08-11 05:08:53	iPhone: 35.509506,139.615601	揺れたね	Hakkan	<a href="http://twitter.com/Hakkan">http://twitter.com/Hakkan</a>
2009-08-11 05:08:53	Mie Prefecture	すごい地震だ [mb]	narude531 masu	<a href="http://twitter.com/narude531 masu">http://twitter.com/narude531 masu</a>
2009-08-11 05:08:52	Kawasaki city	地震だ！！	yaketazamma	<a href="http://twitter.com/yaketazamma">http://twitter.com/yaketazamma</a>
2009-08-11 05:08:52	unknown	地震こわいですかんへん	wzcc	<a href="http://twitter.com/wzcc">http://twitter.com/wzcc</a>
2009-08-11 05:08:52	Kansai	あら、地震？	HARU_IRO	<a href="http://twitter.com/HARU_IRO">http://twitter.com/HARU_IRO</a>
2009-08-11 05:08:52	Sakado, Saitama, Japan	地震だ	d_wackys	<a href="http://twitter.com/d_wackys">http://twitter.com/d_wackys</a>
2009-08-11 05:08:51	unknown	愛知も揺れたw	edoman	<a href="http://twitter.com/edoman">http://twitter.com/edoman</a>
2009-08-11 05:08:51	unknown	また地震 長いな	lauk.az	<a href="http://twitter.com/lauk.az">http://twitter.com/lauk.az</a>
2009-08-11 05:08:51	JP	地震なる	echomitt	<a href="http://twitter.com/echomitt">http://twitter.com/echomitt</a>

# Earthquake Reporting System

- Effectiveness of alerts of this system
  - Alert E-mails urges users to prepare for the earthquake if they are received by a user shortly before the earthquake actually arrives.
- Is it possible to receive the e-mail before the earthquake actually arrives?
  - An earthquake is transmitted through the earth's crust at about 3~7 km/s.
  - a person has about **20~30 sec** before its arrival at a point that is 100 km distant from an actual center

# Results of Earthquake Detection

- In all cases, we sent E-mails before announces of JMA
- In the earliest cases, we can sent E-mails in 19 sec.

Date	Magnitude	Location	Time	E-mail sent time	time gap [sec]	# tweets within 10 minutes	Announce of JMA
Aug. 18	4.5	Tochigi	6:58:55	7:00:30	95	35	7:08
Aug. 18	3.1	Suruga-wan	19:22:48	19:23:14	26	17	19:28
Aug. 21	4.1	Chiba	8:51:16	8:51:35	19	52	8:56
Aug. 25	4.3	Uraga-oki	2:22:49	2:23:21	31	23	2:27
Aug.25	3.5	Fukushima	2:21:15	22:22:29	73	13	22:26
Aug. 27	3.9	Wakayama	17:47:30	17:48:11	41	16	1:7:53
Aug. 27	2.8	Suruga-wan	20:26:23	20:26:45	22	14	20:31
Ag. 31	4.5	Fukushima	00:45:54	00:46:24	30	32	00:51
Sep. 2	3.3	Suruga-wan	13:04:45	13:05:04	19	18	13:10
Sep. 2	3.6	Bungo-suido	17:37:53	17:38:27	34	3	17:43

# Discussion

- Advantages
  - No need of dedicate devices
  - Provide a fast detection and assessment of earthquake
  - Possibility of early alarm
- Limitations of earthquake detection with Twitter
  - Need enough Twitter samples
    - Depend on the population distribution of Twitter
    - If the center of a target event is in an oceanic area, it's more difficult to locate it
    - The number of tweets maybe not as large as we had anticipated.
  - Could be not accurate
    - Interfered by retweets (not in the earthquake area)
    - Incorrect tweet geolocations
  - Could be unstable
    - The network service is not reliable when earthquake happens
    - Vulnerable to hacker attacks
  - Still not fast enough

# Other Possible Applications

- Social life
  - Detect the hot news in the world
- Economy
  - Detect the trend of stocks
- Politics
  - Predict and evaluate president selection and other political events
- Science
  - Discover patterns of social interactions and influences

○ ...

# Overview: Social Network Analysis

- Concepts
  - Social network
  - Characteristics of social networks
  - Connectivity
  - Giant component
  - Community
  - Betweenness
  - The Bow-Tie Structure of the Web
  - Information Cascade
  - Power-law distribution
  - Diameter
  - Clustering coefficient
  - Six degrees of separation and Milgram experiment
  - Regular Network, Small-World Network, Random Network
  - Geographical distance
  - Social centrality
  - Social similarity
  - Mobile social network

# Overview: Social Network Analysis

- Models and Algorithms
  - Community detection algorithm
  - Analytical model for information cascade
  - Rich get richer model
  - The Watts-Strogatz model
  - Analysis of decentralized search
    - Inverse-Square Principle
  - The SIRS Epidemic Model

# Overview: Social Network Analysis

- Applications
  - Decentralized routing
  - Epidemics
    - Flu detection
  - Tracking Earthquake Using Twitter



- About the Exam
- Q & A
  - Office time: Jul 12, 12:00-12:40, Room 3.110
- Contact:
  - Wenzhong Li
    - [wenzhong.li@informatik.uni-goettingen.de](mailto:wenzhong.li@informatik.uni-goettingen.de)
    - <http://cs.nju.edu.cn/lwz/>

**Thank You!**