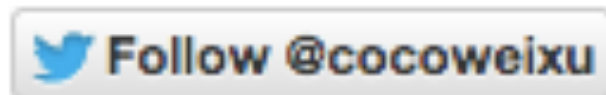


# Social Media & Text Analysis

lecture 8 - Automatic Summarization for Twitter



**CSE 5539-0010 Ohio State University**

**Instructor: Wei Xu**

**Website: [socialmedia-class.org](http://socialmedia-class.org)**

# Homework #3

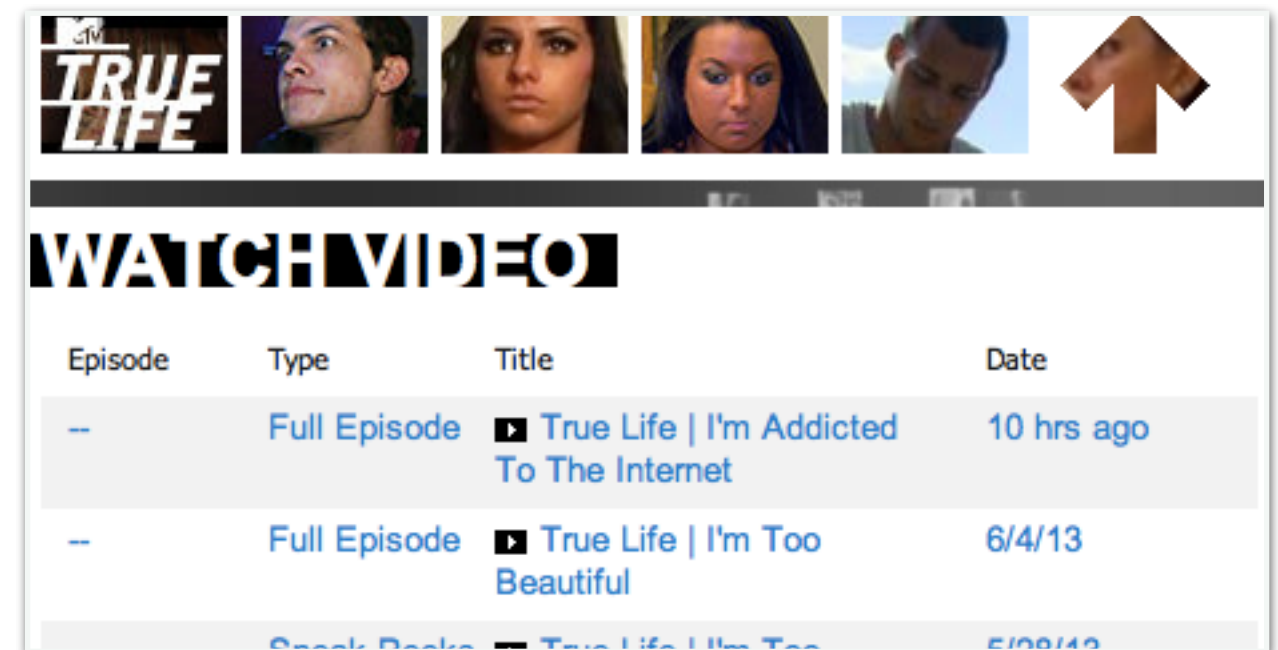
- Implementing a simplified **word2vec** algorithm, including:
  - softmax function
  - neural network basics
  - and word2vec
- Group of 2
- Estimated to release in the 1st week of Nov
- Due in 3~4 weeks

# Summarization



A screenshot of a Twitter feed with eight tweets. The tweets are from various users, some with profile pictures and some with blurred names. The tweets discuss the TV show 'True Life' and its themes of internet addiction and social media. The tweets are as follows:

- Living in the gangsters paradise #truelife (49s)
- True life: we're addicted to #froyo. @lindsayvfenton #bff (50s)
- #nw True Life I'm addicted to the internet gawd this is the story of my life. (53s)
- Watching True Life I am addicted to the Internet...social networks have ruined our society (us) (54s)
- I'm watching True Life: I'm addicted to the Internet and it reminded me of @babyydani (54s)
- Wow RT @jewdith123: OMFG this foo on true life said he has to post 10 shirtless pictures of himself (cont) (58s)
- @DC\_Blackburn haha true! Life in shorts is miles better (1m)
- @chelceebastien holla at us boys! Why aren't we famous? True life: 2 fab 4 Yuma (1m)



A screenshot of a YouTube video player interface. The video player shows a list of episodes from the TV show 'True Life'. The episodes are as follows:

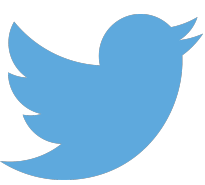
Episode	Type	Title	Date
--	Full Episode	True Life   I'm Addicted To The Internet	10 hrs ago
--	Full Episode	True Life   I'm Too Beautiful	6/4/13
--	Full Episode	True Life   I'm Too Beautiful	6/4/13

## SUMMARY:

I'm watching **true life** "I'm addicted to Internet" ... while I'm on mine lol

Okay these girls on **True Life** I'm Too Beautiful are not that pretty

# Summarization

- Given a (or a set) of documents, generate a short summary
-  Given a (large) set of topically and temporally clustered tweets, select a few representative tweets as the summary.

# Previous Work

Selected Work	Size of Input	Length of Summary
Wei et al. (2012)	average 10k tweets	10 tweets
Inouye & Kalita (2011)	approximately 1500 tweets	4 tweets ❖
Rosa et al. (2011)	average 410 tweets	1, 5, 10 tweets
Liu et al. (2011)	average 1.7k tweets	about 2 or 3 tweets ★
Takamura et al. (2011)	2.8k - 5.2k tweets	26 - 41 tweets ★

❖ Human annotators strongly prefer different numbers of tweets in a summary for different topics.

★ Used the length of human reference summaries to decide the length of system outputs, which information is not available in practice.

# Research Questions

- What is the perfect length of multi-tweet summary?
- Will IE help summarization on Twitter?
  - noisy text: performance of IE?
  - short context: still need in-depth event analysis?
  - redundant: is word enough?

# SumBasic

- Intuition:

words occurring frequently in the documents occur with higher probability in the human summaries than words occurring less frequently

# SumBasic

- a very simple but strong summarization algorithm [Nenkova and Vanderwende, 2005]
- Intuition:  
  
words occurring frequently in the documents occur with higher probability in the human summaries than words occurring less frequently



# SumBasic

- Step 1: computes the probability of each word **w** :

$$P(w) = \frac{n(w)}{\sum_i w_i}$$

- Step 2: computes the salience score of each sentence **S** :

$$Score(S) = \sum_{w \in S} \frac{P(w)}{|\{w \mid w \in S\}|}$$

- Step 3: pick the highest scored sentence into summary
- Step 4: for each word in sentences chosen at step 3, update their probability:

$$P_{new}(w) = P_{old}(w) \cdot P_{old}(w)$$

- Step 5: repeat Step 2~4 until reach desired length of summary

# Varied-length Summary

- For a set of topically clustered tweets, amount of information varies greatly:
  - from very repetitive to very discrete
  - e.g.

album release of a less notable singer

VS.

album release of a famous/controversy singer

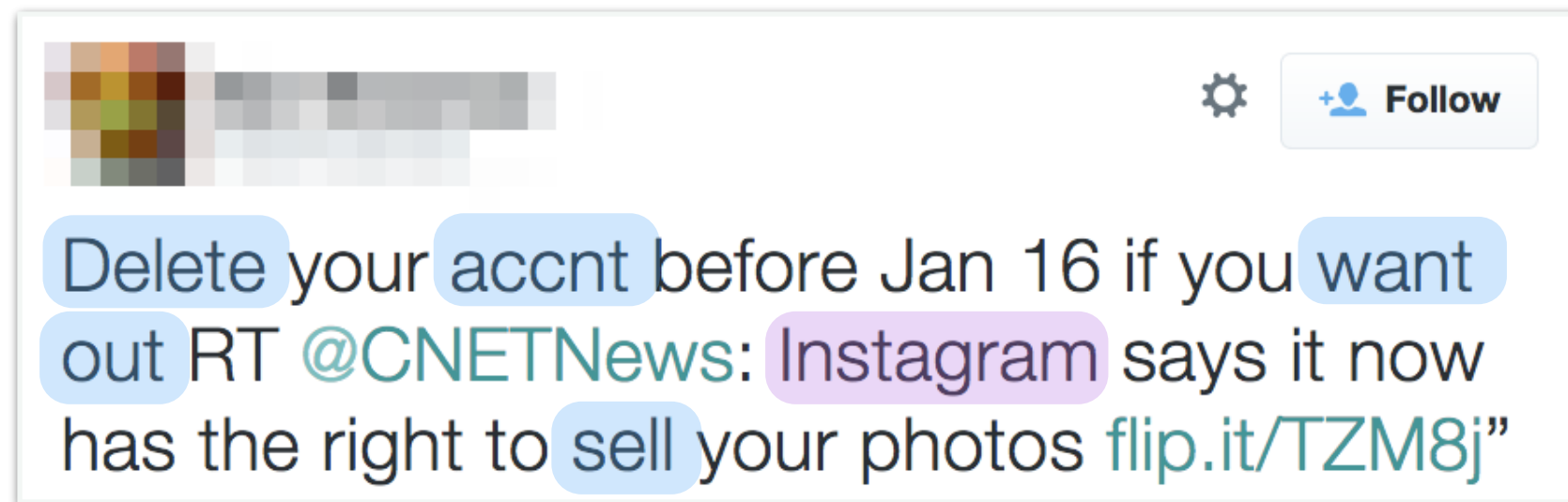
# Information Extraction (IE)

- Named Entity [Ritter et al. 2011]
- Event Phrases [Ritter et al. 2012]

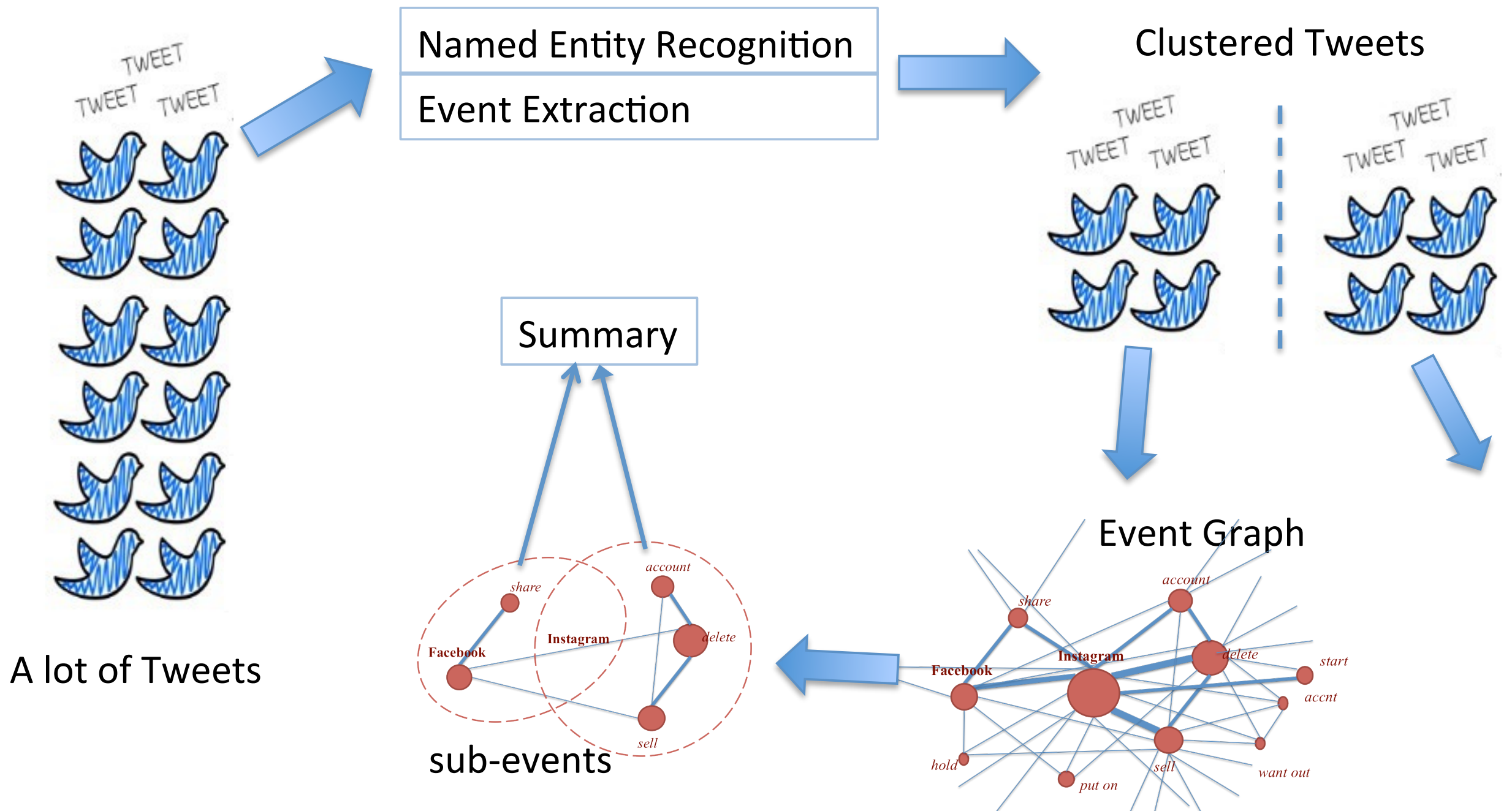


# Information Extraction (IE)

- Named Entity [Ritter et al. 2011]
- Event Phrases [Ritter et al. 2012]



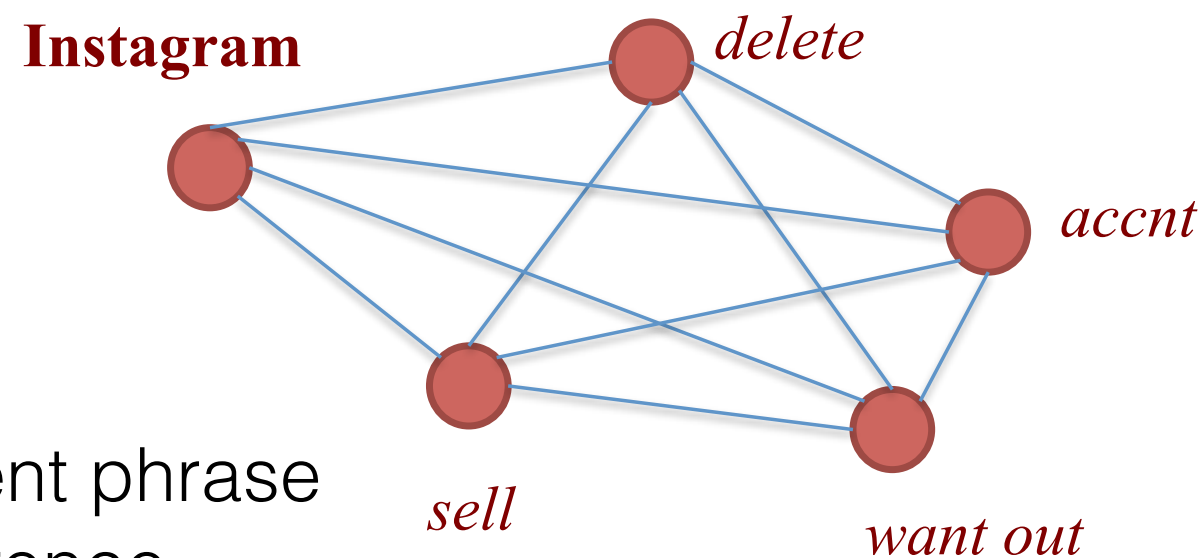
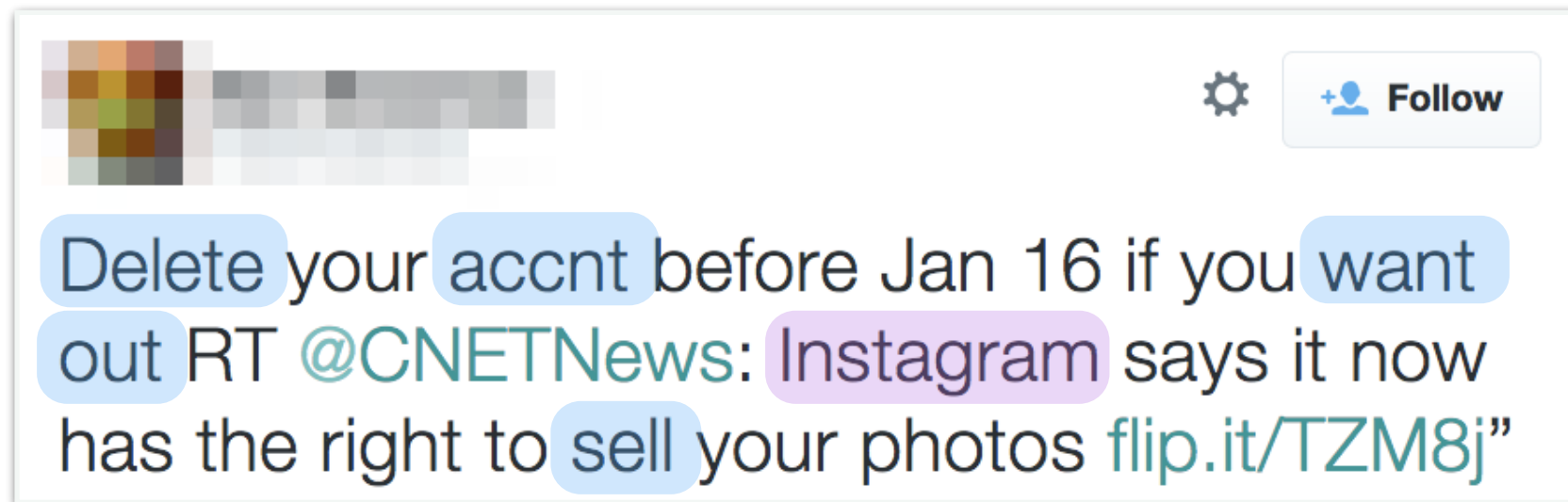
# System Overflow



Wei Xu, Alan Ritter, Ralph Grishman.

"A Preliminary Study of Tweet Summarization using Information Extraction" in LASM (2014)

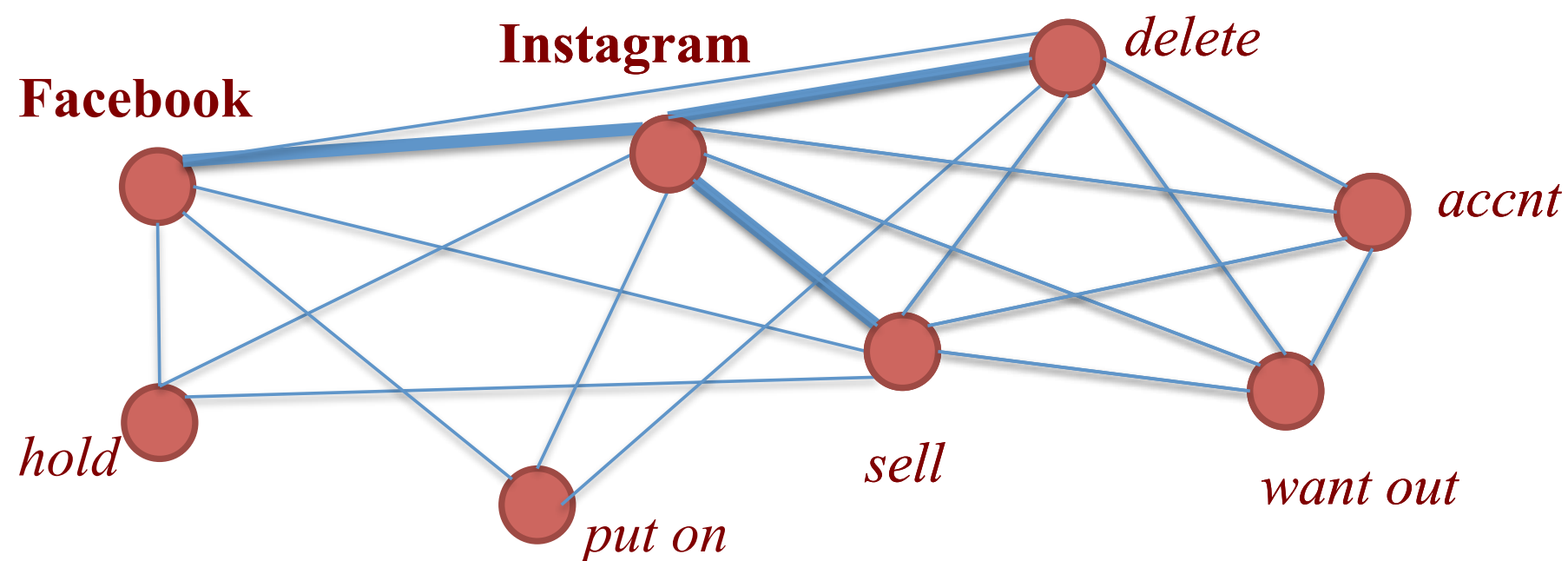
# Event Graph



**Node** - named entities + event phrase

**Edge** (weighted) - co-occurrence

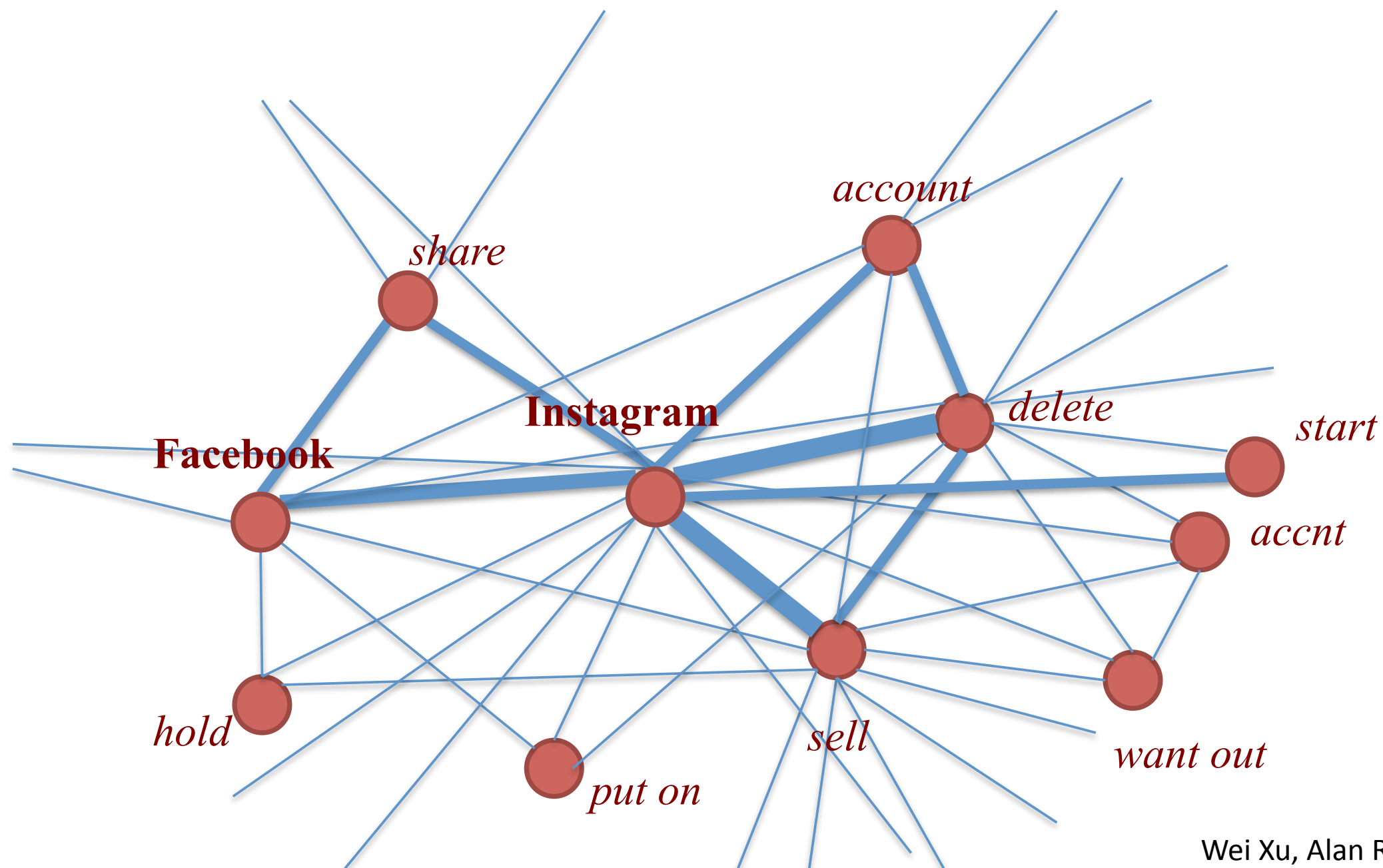
# Event Graph



Wei Xu, Alan Ritter, Ralph Grishman.

“A Preliminary Study of Tweet Summarization using Information Extraction” in LASM (2014)

# Event Graph



Wei Xu, Alan Ritter, Ralph Grishman.

“A Preliminary Study of Tweet Summarization using Information Extraction” in LASM (2014)



# PageRank

- a graph-based ranking algorithm
- a trademark of Google
- Idea: web surfing / random walk


The importance of a webpage is defined recursively and depends on the number and importance of all webpages that link to it.

- also used for local graph partitioning

# PageRank

- Saliency score of nodes:

$$Score(u) = (1 - d) + d \times \sum_{v \in Adj(u)} \frac{Score(v)}{|Adj(v)|}$$

adjacent nodes

- directed graph
- iterate towards converge
- initial rank of node does not matter
- only edges matter
- total weight of the graph stays the same

# PageRank → Event Rank

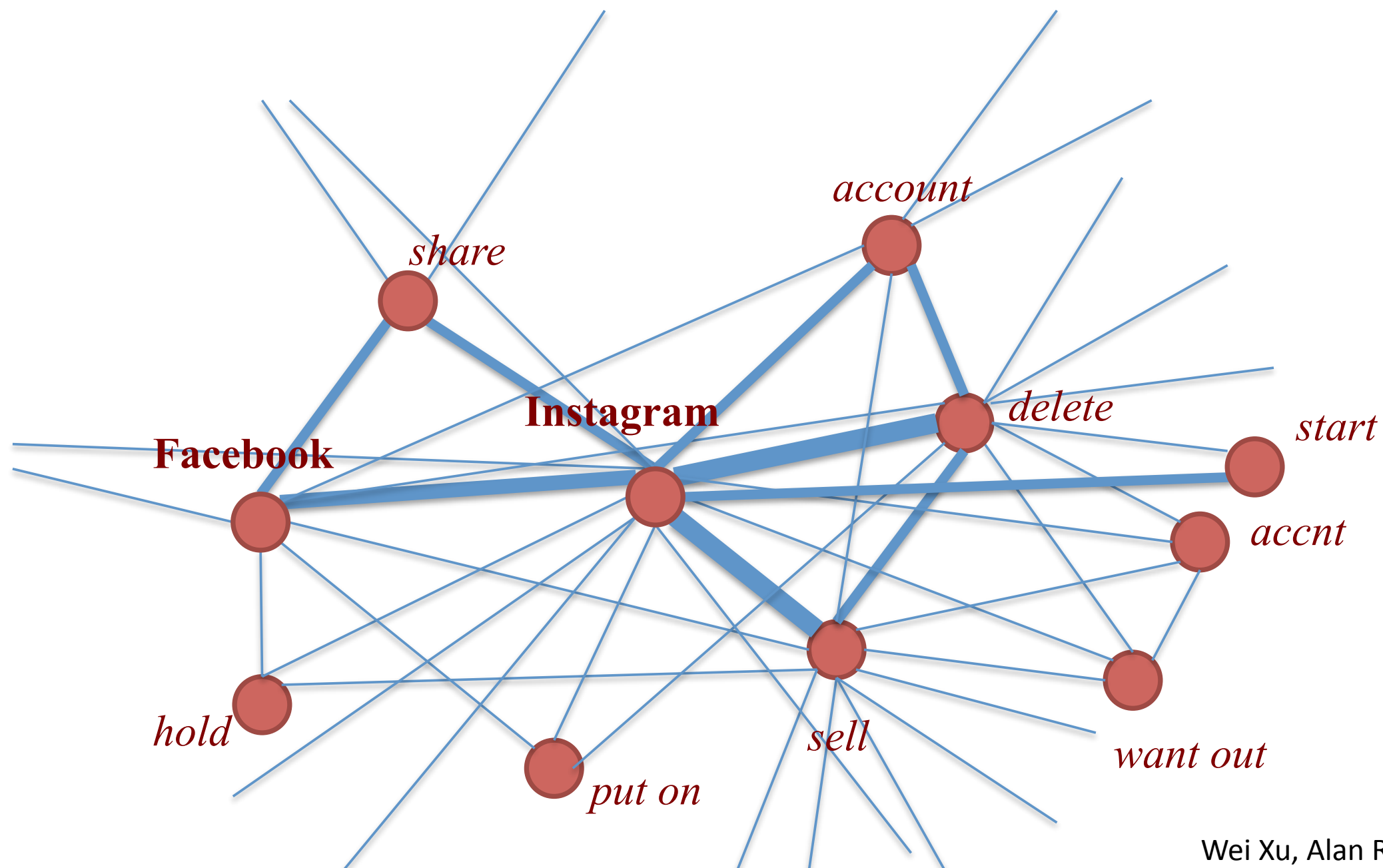
- Saliency score of nodes:

$$Score(u) = (1 - d) + d \times \sum_{v \in Adj(u)} \frac{e_{uv} \times Score(v)}{\sum_{w \in Adj(v)} e_{vw}}$$

adjacent nodes

- undirected graph
- iterate towards converge
- initial rank of node does not matter
- only edges **and their weights** matter
- total weight of the graph stays the same

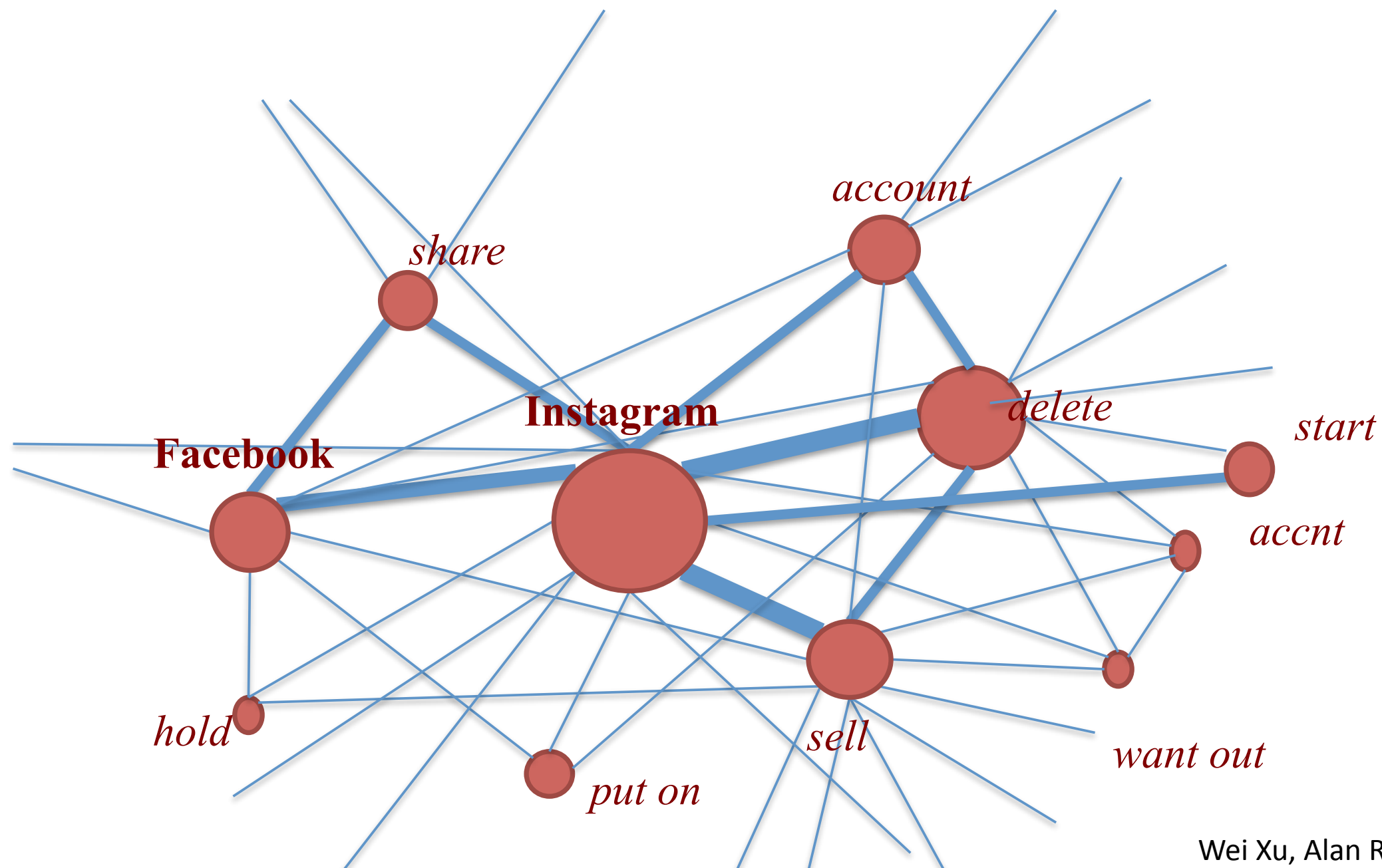
# Graph Ranking



Wei Xu, Alan Ritter, Ralph Grishman.

“A Preliminary Study of Tweet Summarization using Information Extraction” in LASM (2014)

# Graph Ranking

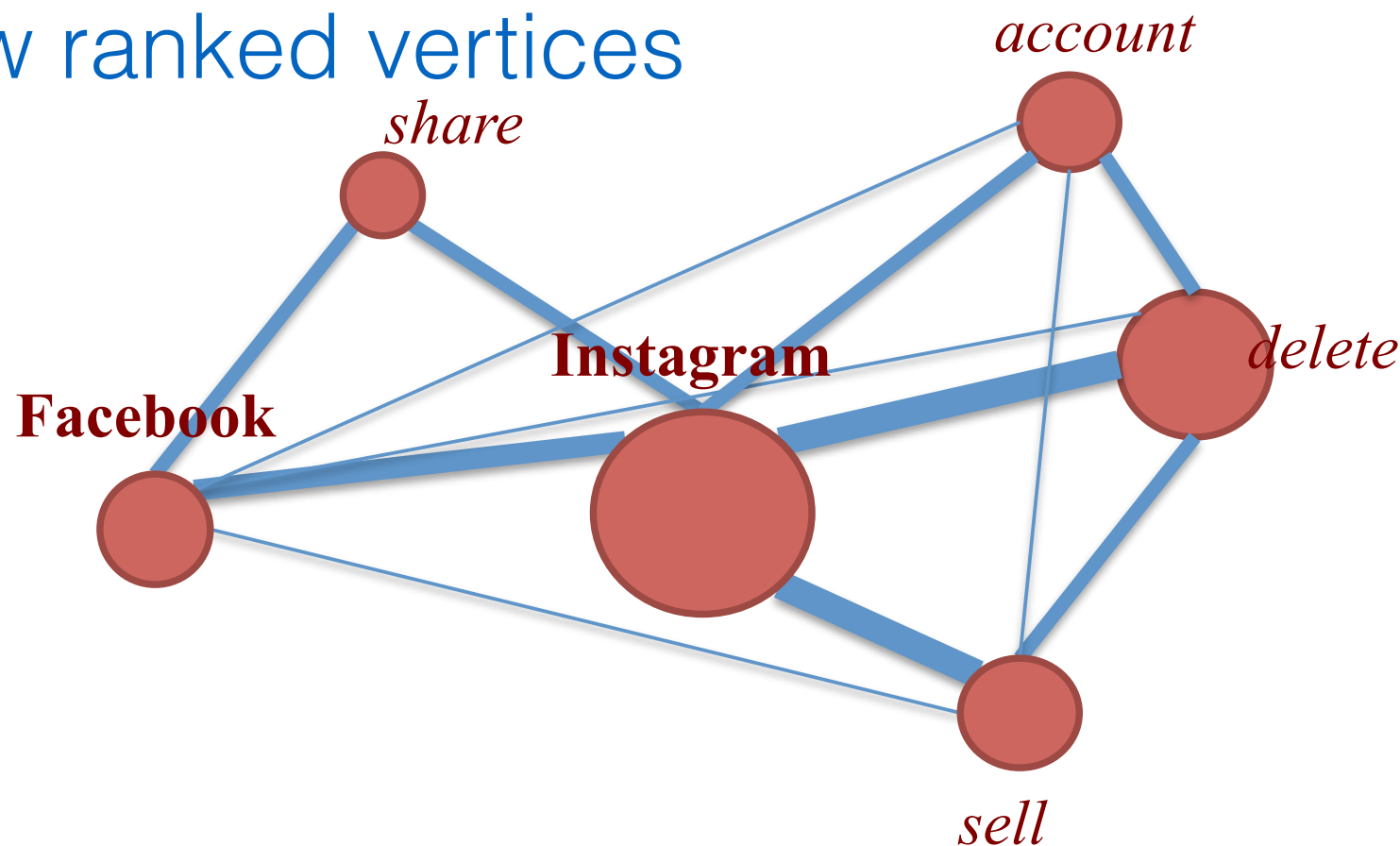


Wei Xu, Alan Ritter, Ralph Grishman.

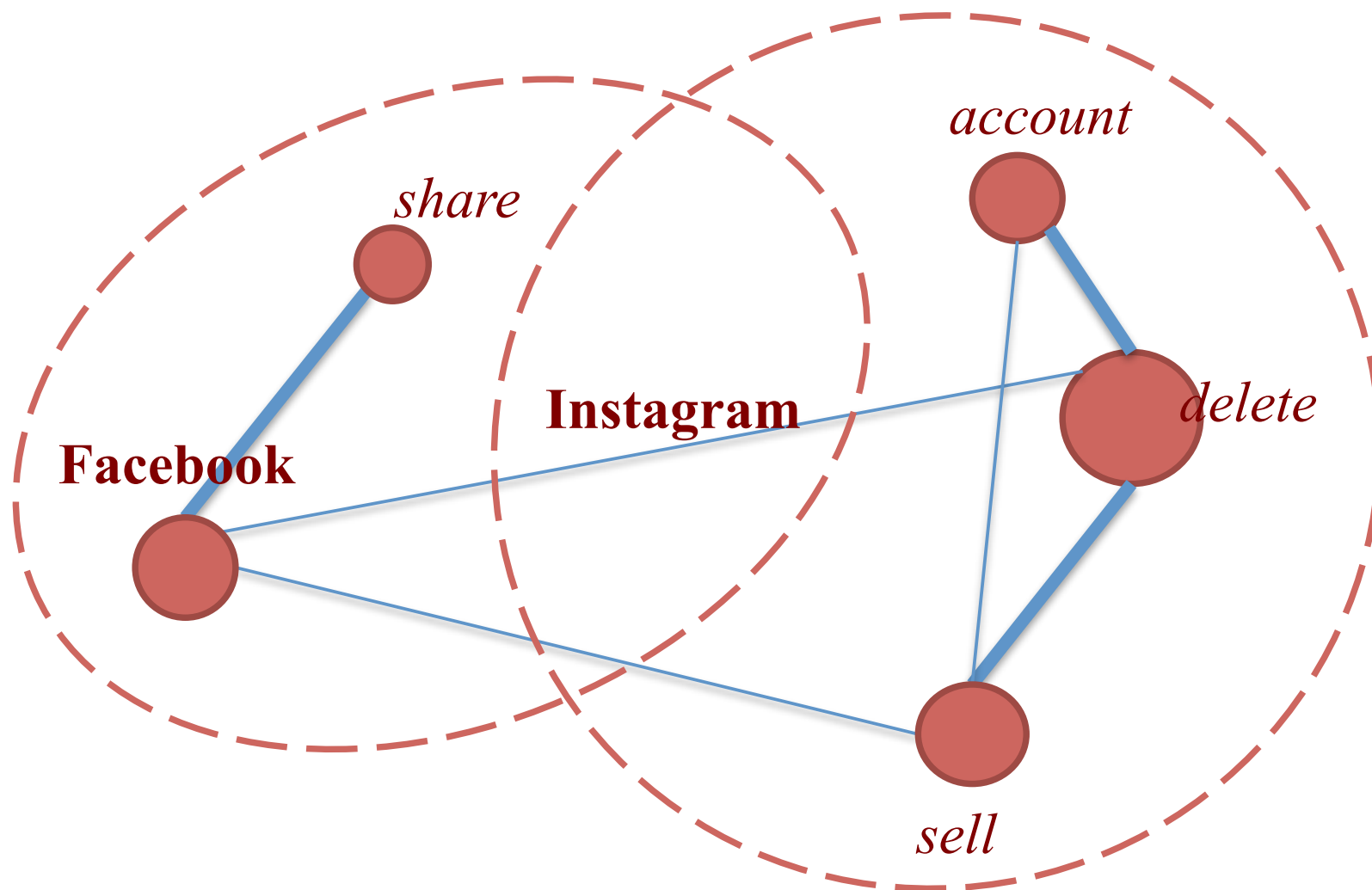
“A Preliminary Study of Tweet Summarization using Information Extraction” in LASM (2014)

# Graph Partitioning

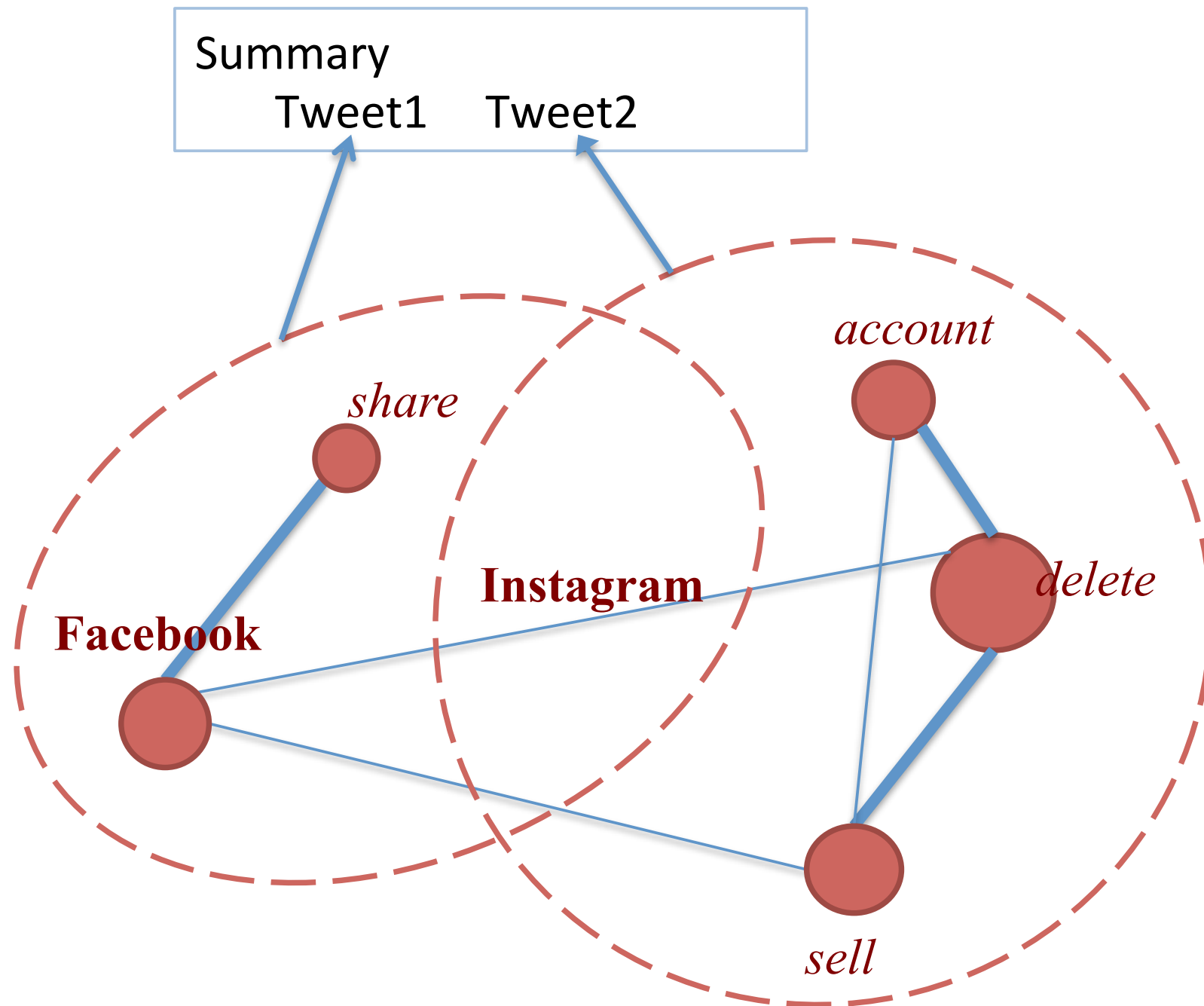
- local graph partitioning by PageRank [Andersen et al., 2006] : a good partition of the graph can be obtained by separating high ranked vertices from low ranked vertices



# Graph Partitioning

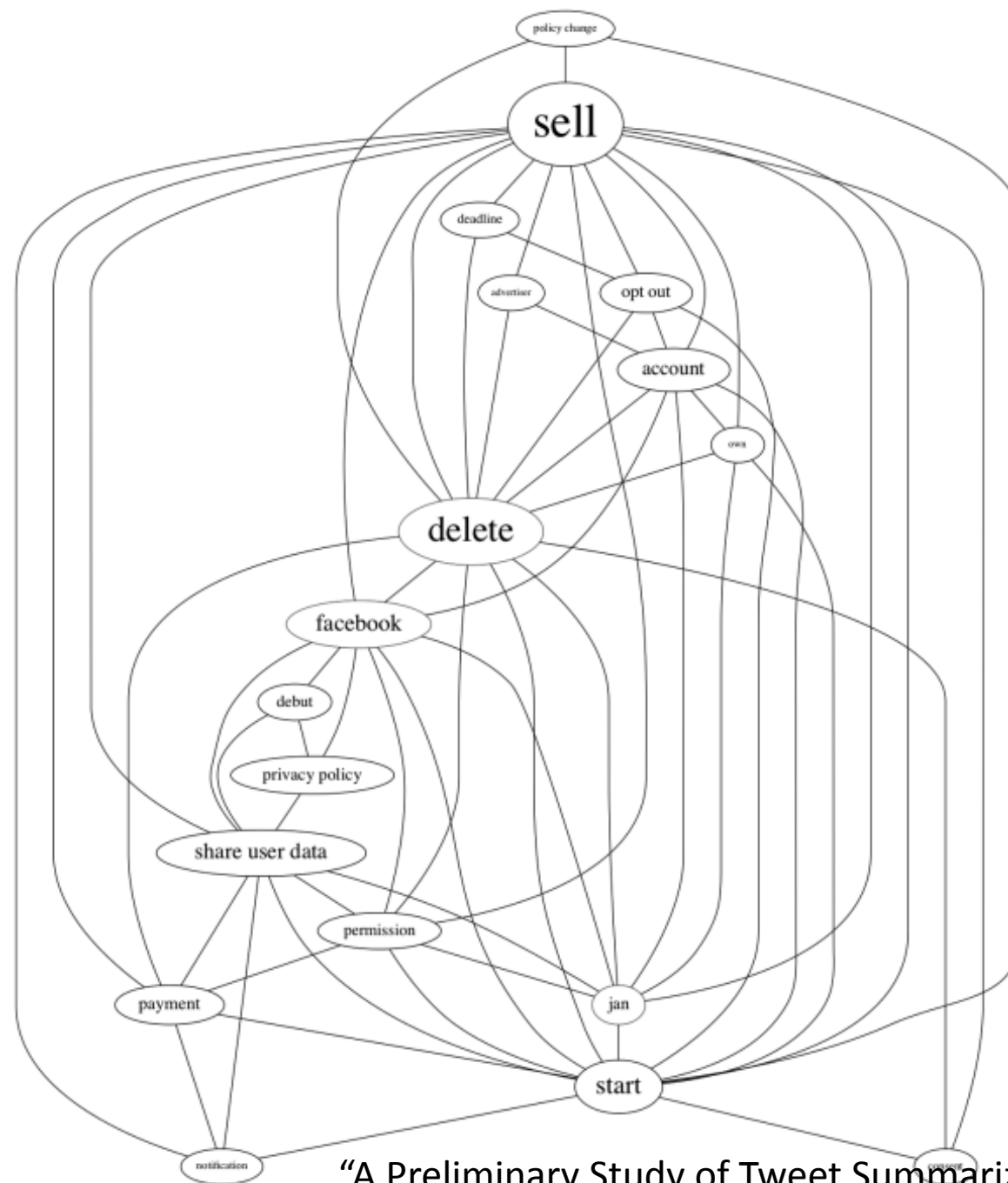


# Graph Partitioning





# Example Event Graph



Wei Xu, Alan Ritter, Ralph Grishman.

“A Preliminary Study of Tweet Summarization using Information Extraction” in LASM (2014)

# Example Summary

Instagram 1/16/2013	EventRank (Flexible)	<ul style="list-style-type: none"><li>- So Instagram can sell your pictures to advertisers without u knowing starting January 16th I'm bout to delete my instagram !</li><li>- Instagram debuts new privacy policy , set to share user data with Facebook beginning January 16</li></ul>
	SumBasic	<ul style="list-style-type: none"><li>- Instagram will have the rights to sell your photos to Advertisers as of jan 16</li><li>- Over for Instagram on January 16th</li><li>- Instagram says it now has the right to sell your photos unless you delete your account by January 16th <a href="http://t.co/tsjic6yA">http://t.co/tsjic6yA</a></li></ul>

# Example Event Graph

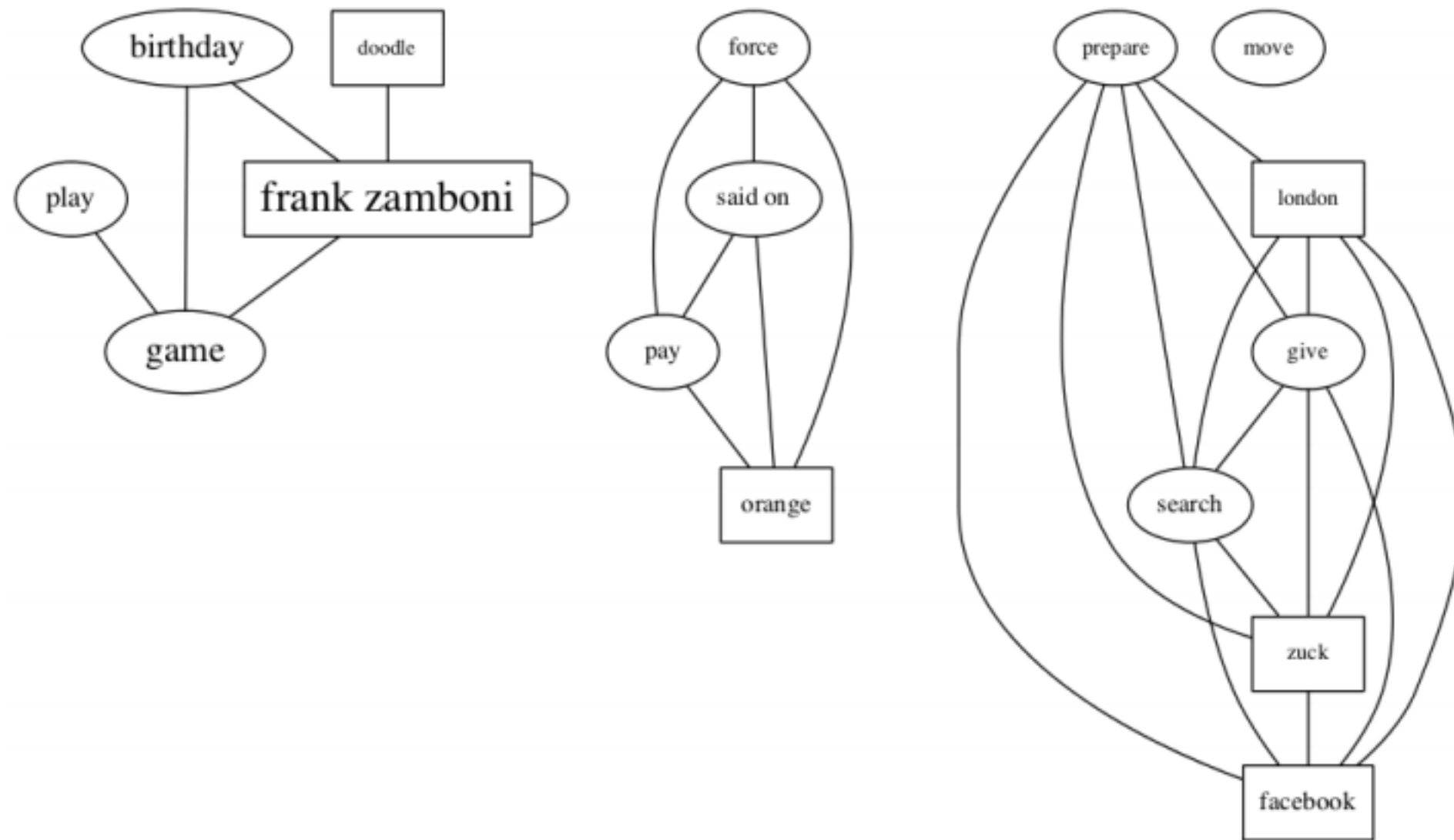


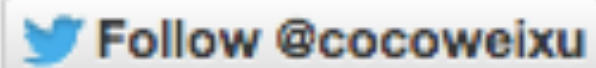
Figure 2: Event graph of 'Google - 1/16/2013', an example of event cluster with multiple focuses

# Example Summary

Google 1/16/2013	EventRank (Flexible)	<ul style="list-style-type: none"><li>- Google 's home page is a Zamboni game in celebration of Frank Zamboni 's birthday January 16 #GameOn</li><li>- Today social , Tomorrow Google ! Facebook Has Publicly Redefined Itself As A Search Company <a href="http://t.co/dAevB2V0">http://t.co/dAevB2V0</a> via @sai</li><li>- Orange says has it has forced Google to pay for traffic . The Head of the Orange said on Wednesday it had ... <a href="http://t.co/dOqAHhWi">http://t.co/dOqAHhWi</a></li></ul>
	SumBasic	<ul style="list-style-type: none"><li>- Tomorrow's Google doodle is going to be a Zamboni! I may have to take a vacation day.</li><li>- the game on google today reminds me of hockey #tooexcited #saturday</li><li>- The fact that I was soooo involved in that google doodle game says something about this Wednesday #TGIW You should try it!</li></ul>

# Research Questions

- What is the perfect length of multi-tweet summary?  
variable length
- Will IE help summarization on Twitter?
  - noisy text: performance of IE?  
summary is more readable and newsworthy
  - short context: still need in-depth event analysis?  
self-contained (no coref.) → better event graph
  - redundant: is word enough?  
unbalanced event graph → easier partitioning



**Instructor: Wei Xu**  
**[www.cis.upenn.edu/~xwe/](http://www.cis.upenn.edu/~xwe/)**

**Course Website: [socialmedia-class.org](http://socialmedia-class.org)**