

Uncovering News-Twitter Reciprocity via Interaction Patterns

Yue Ning¹ Sathappan Muthiah¹
Ravi Tandon² Naren Ramakrishnan¹

¹Discovery Analytics Center, Department of Computer Science, Virginia Tech

²Now with Department of Electrical and Computer Engineering, The University of Arizona



Outline

Introduction

- Problem Definition

Methodology

- Story Chaining

- Retrieval of Tweets

- Identify Interaction Patterns

- Clustering

- Topic Modeling

Experiments and Results

- Dataset

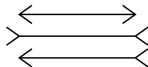
- Results

Conclusion

Problem Introduction



Social Media

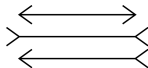


News Media

Problem Introduction



Social Media



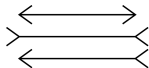
News Media



MIPTV: Ellen DeGeneres' Oscar Selfie Worth as Much as \$1 Billion



Problem Introduction



Social Media

News Media

Hollywood.com TV REVIEWS LIVE FEED FEINBERG FORECAST

MIPTV: Ellen DeGeneres' Oscar Selfie Worth as Much as \$1 Billion

The Guardian

Look at the newest colour photo of Pluto from Nasa on 13/2014

Pluto is stunning in latest color close-up from Nasa
Close-up images taken by the New Horizons spacecraft are combined with color data to paint a new and surprising portrait of the dwarf planet, Nasa said.
[View on web](#)

Motivation

- ▶ **News -> Twitter**
- ▶ Twitter -> News Media
- ▶ Explosion of information to comment/feed upon
- ▶ Cause for variations in such interdependencies
 - ▶ Temporal popularity of a "topic"
 - ▶ Geo-location (Africa vs Asia)

Motivation

- ▶ News -> Twitter
- ▶ Twitter -> News Media
- ▶ Explosion of information to comment/feed upon
- ▶ Cause for variations in such interdependencies
 - ▶ Temporal popularity of a "topic"
 - ▶ Geo-location (Africa vs Asia)

Motivation

- ▶ News -> Twitter
- ▶ Twitter -> News Media
- ▶ Explosion of information to comment/feed upon
- ▶ Cause for variations in such interdependencies
 - ▶ Temporal popularity of a "topic"
 - ▶ Geo-location (Africa vs Asia)

Motivation

- ▶ News -> Twitter
- ▶ Twitter -> News Media
- ▶ Explosion of information to comment/feed upon
- ▶ Cause for variations in such interdependencies
 - ▶ Temporal popularity of a "topic"
 - ▶ Geo-location (Africa vs Asia)

Motivation

- ▶ News -> Twitter
- ▶ Twitter -> News Media
- ▶ Explosion of information to comment/feed upon
- ▶ Cause for variations in such interdependencies
 - ▶ Temporal popularity of a "topic"
 - ▶ Geo-location (Africa vs Asia)

Motivation

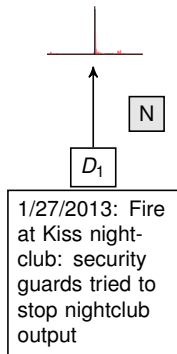
- ▶ News -> Twitter
- ▶ Twitter -> News Media
- ▶ Explosion of information to comment/feed upon
- ▶ Cause for variations in such interdependencies
 - ▶ Temporal popularity of a "topic"
 - ▶ Geo-location (Africa vs Asia)

One Example

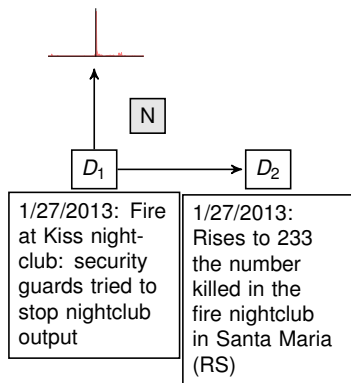
D_1

1/27/2013: Fire
at Kiss night-
club: security
guards tried to
stop nightclub
output

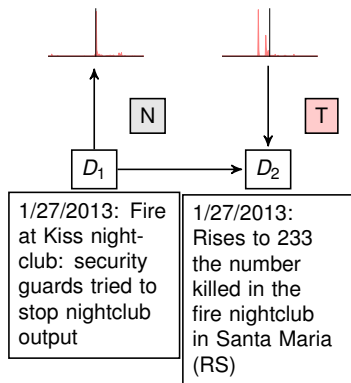
One Example



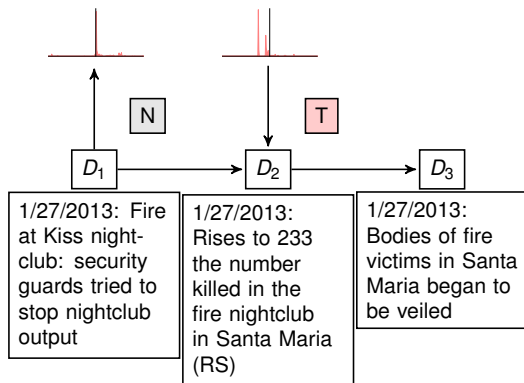
One Example



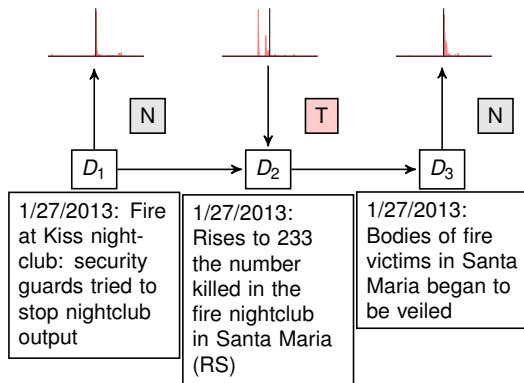
One Example



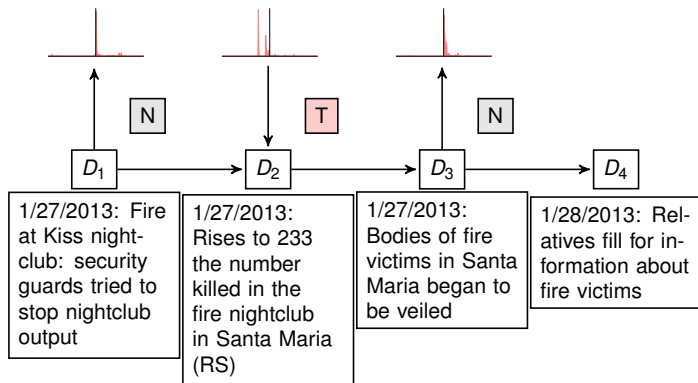
One Example



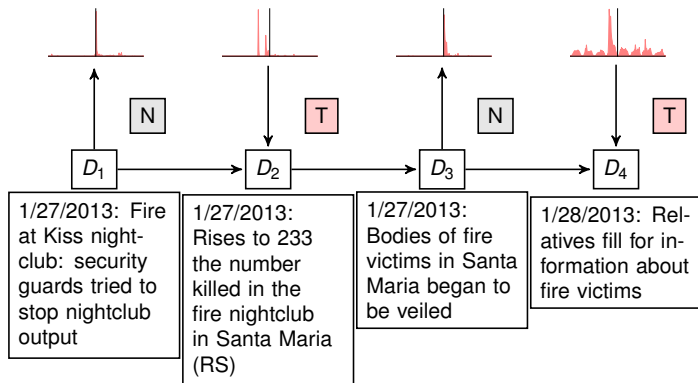
One Example



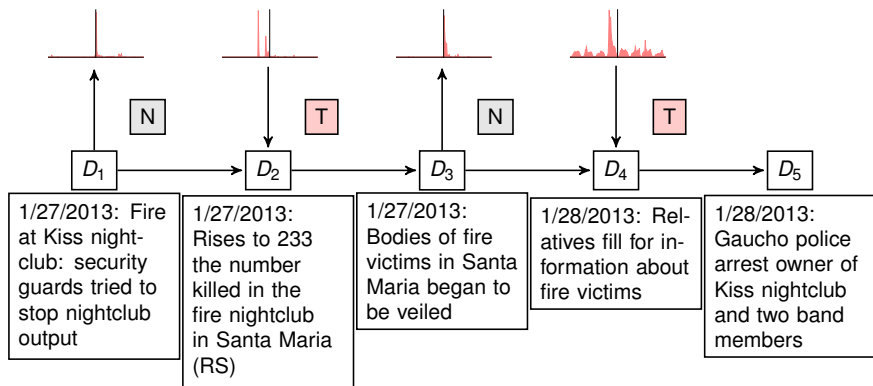
One Example



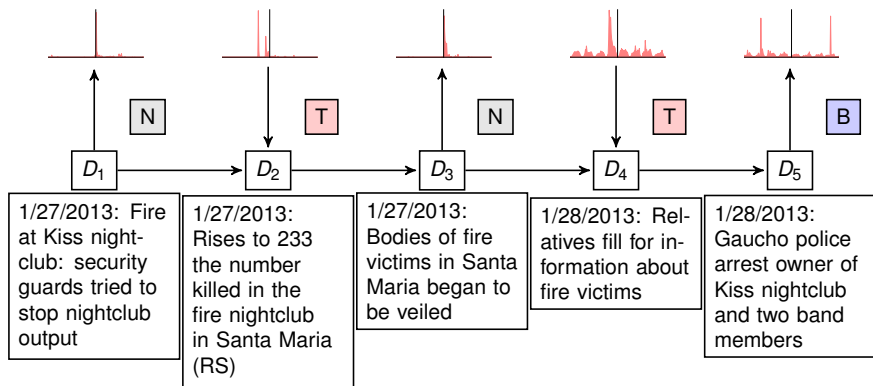
One Example



One Example



One Example



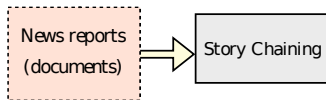
Goals

1. Understanding **the type of information flow** between news and Twitter.
2. **Chaining** similar news articles together.
3. Identifying **major interaction patterns**
 - ▶ **Cluster** story chains and understanding their differences
 - ▶ **Identify main topics** of interest within such clusters.

System Framework

News reports
(documents)

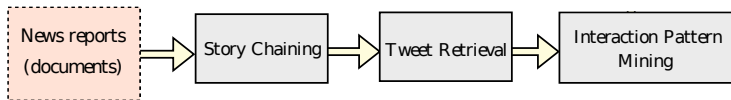
System Framework



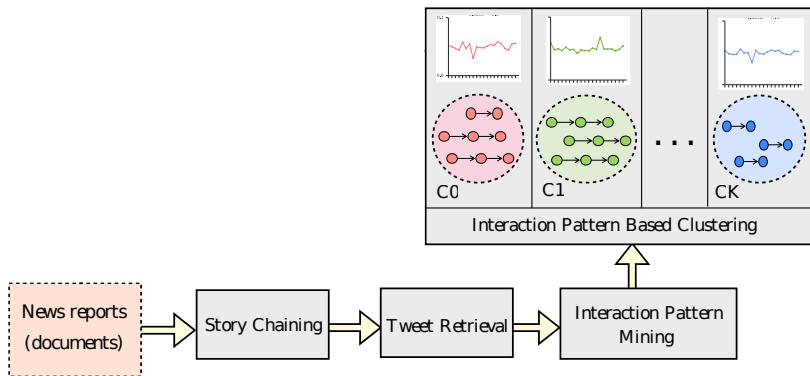
System Framework



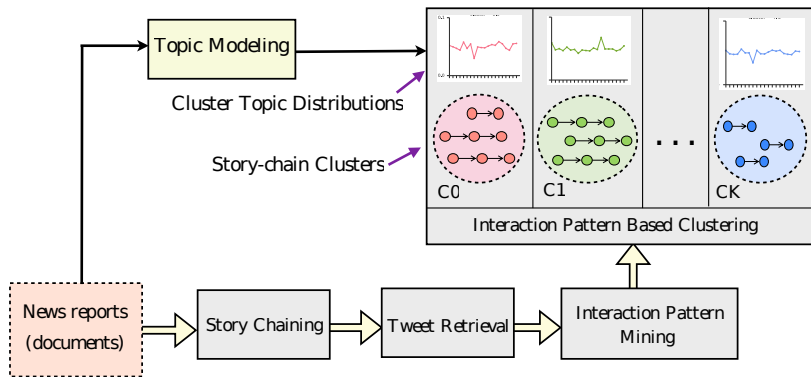
System Framework



System Framework



System Framework



Story Chaining Algorithm

¹ **Goal:** identifying all documents related to a news story and to keep track of the news story as new documents arrive.

Method: To assess if two documents are referring to the same underlying context, we calculate their similarity scores with respect to three features:

- ▶ - textual features, denoted by $T(D_i)$
- ▶ - spatial features, denoted by $L(D_i)$, e.g. city, state, country
- ▶ - actors, denoted by $A(D_i)$, e.g. Hillary Clinton.

¹J. Schlachter, A. Ruvinsky, L. Asencios Reynoso, S. Muthiah, and N. Ramakrishnan, "Leveraging topic models to develop metrics for evaluating the quality of narrative threads extracted from news stories", in *Proc. of the 6th International Conference on Applied Human Factors and Ergonomics, AHFE*, Elsevier, 2015.

Story Chaining Algorithm (Cont.)

The total weighted similarity measure between two documents, D_i and D_j , is then defined as follows:

$$\begin{aligned} \text{sim}(D_i, D_j) \triangleq & \underbrace{\alpha f(\mathcal{T}(D_i), \mathcal{T}(D_j))}_{\text{textul features}} + \underbrace{\beta f(\mathcal{L}(D_i), \mathcal{L}(D_j))}_{\text{spatial features}} \\ & + \underbrace{\eta f(\mathcal{A}(D_i), \mathcal{A}(D_j))}_{\text{actor features}} \end{aligned}$$

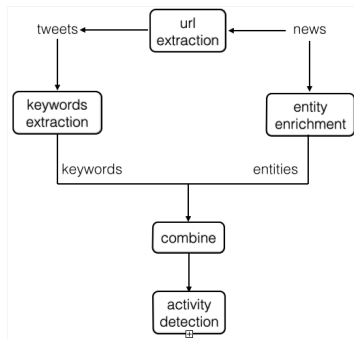
The coherence between a chain C_j and document D_i is defined as

$$\text{coh}(D_i, C_j) = \theta g(\mathcal{L}(D_i), \mathcal{L}(C_j)) + \phi g(\mathcal{A}(D_i), \mathcal{A}(C_j))$$

where g is any similarity measure and the coefficients θ, ϕ are chosen such that $\theta + \phi = 1$.

Twitter Profile for News

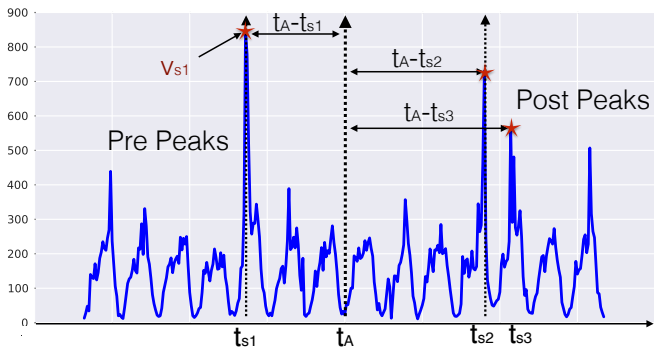
1. Collect tweets based on URL.
2. Extract entity keywords from news.
3. Filter keywords together.
4. Download hourly count metrics.



Interaction Patterns

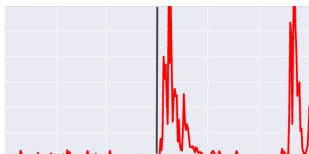
- ▶ Peak detection ²
- ▶ Incoming influence (\mathcal{W}^{pre}) and outgoing influence ($\mathcal{W}^{\text{post}}$):

$$\mathcal{W}^{\text{pre}} = \sum_{s \in S_{\text{pre}}} \frac{v_s}{t_A - t_s}, \quad \mathcal{W}^{\text{post}} = \sum_{s \in S_{\text{post}}} \frac{v_s}{t_s - t_A} \quad (1)$$

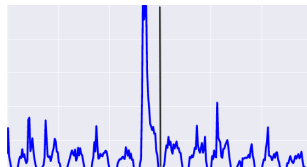


²M. Duarte, “Notes on scientific computing for biomechanics and motor control”, 2015.

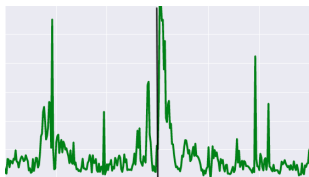
Interaction States



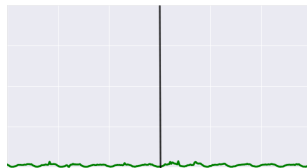
N



T



B



E

Interaction States (Cont.)

$$\text{State}(D_i) = \begin{cases} N, & \text{if } \mathcal{W}^{\text{pre}} < \rho, \mathcal{W}^{\text{post}} \geq (1 + \lambda)\mathcal{W}^{\text{pre}} \\ E, & \text{if } \mathcal{W}^{\text{pre}} < \rho, \mathcal{W}^{\text{post}} < (1 + \lambda)\mathcal{W}^{\text{pre}} \\ T, & \text{if } \mathcal{W}^{\text{pre}} \geq \rho, \mathcal{W}^{\text{post}} < (1 + \lambda)\mathcal{W}^{\text{pre}} \\ B, & \text{if } \mathcal{W}^{\text{pre}} \geq \rho, \mathcal{W}^{\text{post}} \geq (1 + \lambda)\mathcal{W}^{\text{pre}} \end{cases}$$

Interaction States (Cont.)

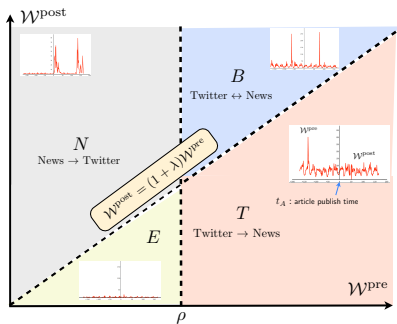


Figure: Geometric Interpretation of States

Interpretation from a Different Dimension

- ▶ Are sports events always related with Bi-directional Interactions?
- ▶ Do Twitter users focus more on sports and entertainment?
- ▶ Latent Dirichlet Allocation (LDA) for hidden topic analysis on clusters. The topic distributions for one cluster is defined by:

$$\mathbf{C}_{j,k} = \frac{\sum_{d_{ij} \in c_j} n_{d_{ij}} \theta(d_{ij}, k)}{\sum_{d_{ij}} n_{d_{ij}}}, \quad (2)$$

where

- ▶ $n_{d_{ij}}$ refers to the frequency of d_i in cluster C_j .
- ▶ $\theta(d_{ij}, k)$ refers to the topic proportions for this document.
- ▶ k is the topic index.

Interpretation from a Different Dimension

- ▶ Are sports events always related with Bi-directional Interactions?
- ▶ Do Twitter users focus more on sports and entertainment?
- ▶ Latent Dirichlet Allocation (LDA) for hidden topic analysis on clusters. The topic distributions for one cluster is defined by:

$$\mathbf{C}_{j,k} = \frac{\sum_{d_{ij} \in c_j} n_{d_{ij}} \theta(d_{ij}, k)}{\sum_{d_{ij}} n_{d_{ij}}}, \quad (2)$$

where

- ▶ $n_{d_{ij}}$ refers to the frequency of d_i in cluster C_j .
- ▶ $\theta(d_{ij}, k)$ refers to the topic proportions for this document.
- ▶ k is the topic index.

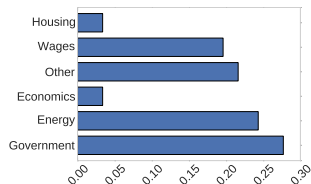
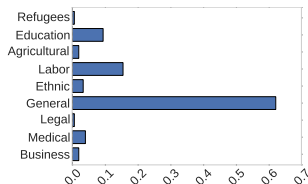
Dataset

Real data from Brazil during the period from Nov. 2012 to Sep. 2013:

- ▶ Protest related articles: GSR
- ▶ Other articles: NON-GSR

Table: Statistical properties of GSR and Non-GSR chains.

Category	% of Twitter starts	Avg-Time-Lag(hour)
GSR Chains	40%	10.95
Non-GSR Chains	73%	5.26



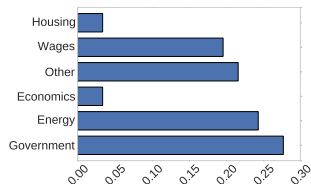
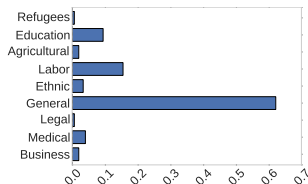
Dataset (Cont.)

Real data from Brazil during the period from Nov. 2012 to Sep. 2013:

- ▶ Protest related articles: GSR
- ▶ Other articles: NON-GSR

Table: Statistical properties of GSR and Non-GSR chains.

Category	% of Twitter starts	Avg-Time-Lag(hour)
GSR Chains	40%	10.95
Non-GSR Chains	73%	5.26



GSR Dataset

Category	% News starts	% of Twitter starts
Housing related protests	100%	0%
Agriculture	100%	0%
Medical	74%	26%
Other (religious & cultural)	60%	40%
General Population	30%	70%
Govt. Policies	23%	77%

Table: % of Twitter, News starts for GSR story-chains

Cluster Results

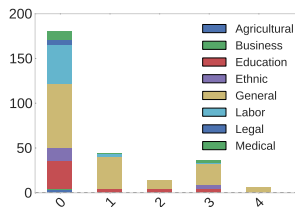


Figure: Population Distribution of Clusters (K-Medoids)

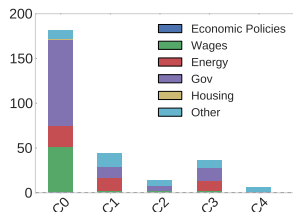


Figure: Event Type Distribution of Clusters (K-Medoids)

Topics in Clusters

ID	Frequent Sub-patterns	Top Topics
C0	“NBNBTNTN”, “NTNTN”	Local Events
C1	“NT”, “NTNT”	Local Events
C2	“TNT”	Local Events, Ads, Technology
C3	“T”, “TB”	Others, Protest, Sports
C4	“TNENT”, “TEB”	Protest, Government, Entertainment

Table: Top topics for clusters

Main Influencer

We define the *influence weight* of a story chain as the average of the difference of pre- and post- influence weights:

$$\frac{\sum_i (\mathcal{W}_i^{\text{pre}} - \mathcal{W}_i^{\text{post}})}{n}$$

where the summation is over n , the number of articles in a chain.

Main Influencer (Cont.)

Table: Story Chains with Interaction Patterns and Main Influencer

ID	IP	IW	MI	Story Summary
SC1	TT	0.514	Twitter	"Marco Feliciano protest at church door"
SC2	TN	0.48	Twitter	"25% Teachers are on strike."
SC3	NNNNBNTBN	-0.422	News	"Fire in Kiss Nightclub in Santa Maria "
SC4	NBNNTN	-0.405	News	"Governor decree official mourning"
SC5	TTTNN	5.0e-05	Both	"Nadal back to Brazil"
SC6	NNTNTTNTN	-1.7e-04	Both	"Nissan sells more than 100 thousand"

Conclusion

- ▶ A new framework for discovering the direction of information flow over time across **news and Twitter**.
- ▶ Uncover the **interaction patterns** over stories and test our proposed method on real data.
- ▶ Cluster on encoded story chains and discover topics

Observation 1

Twitter as a social network platform serves as a fast way to draw attention from public for many social events such as sports news.

Observation 2

News media is quicker to report events regarding political, economical and business issues.

Thank you!
Q&A