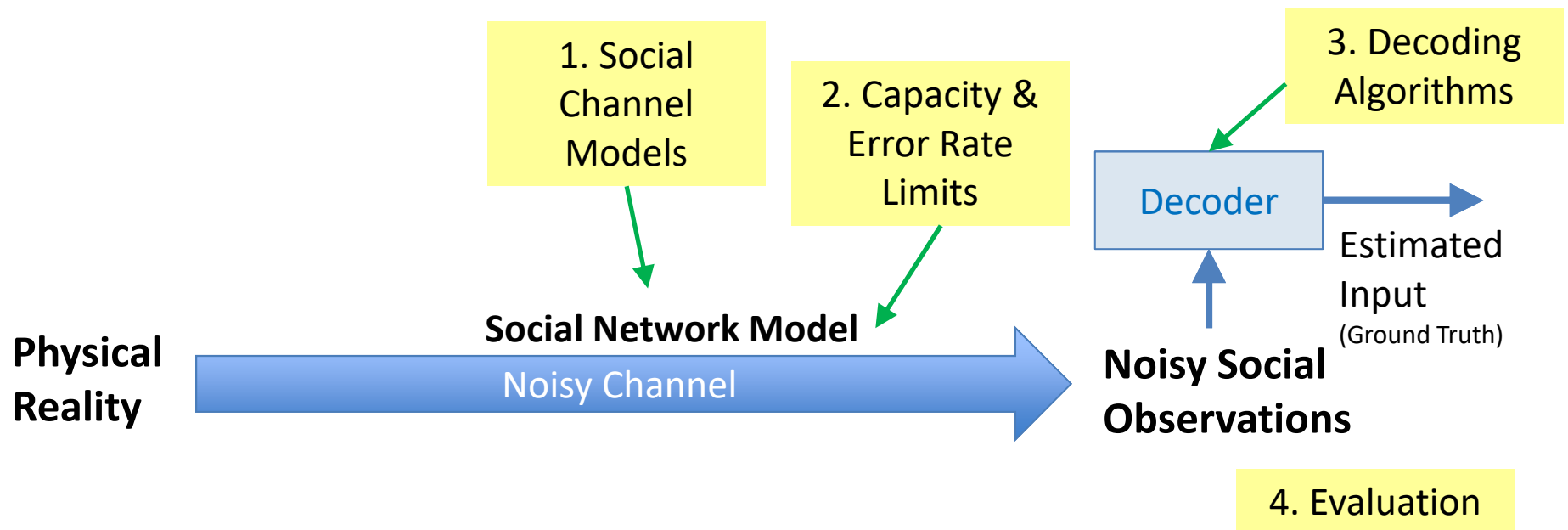# Veracity Analysis

# Fact-finding Research Motivation and Approach

**Goal:** Develop a mathematical foundation for "*social sensing*" – the exploitation of noisy social network data to attain reliable situation awareness.

1. Construct *models of "social channels"*

2. Establish the *fundamental feasibility/accuracy limits* on truth recovery from noisy social network data

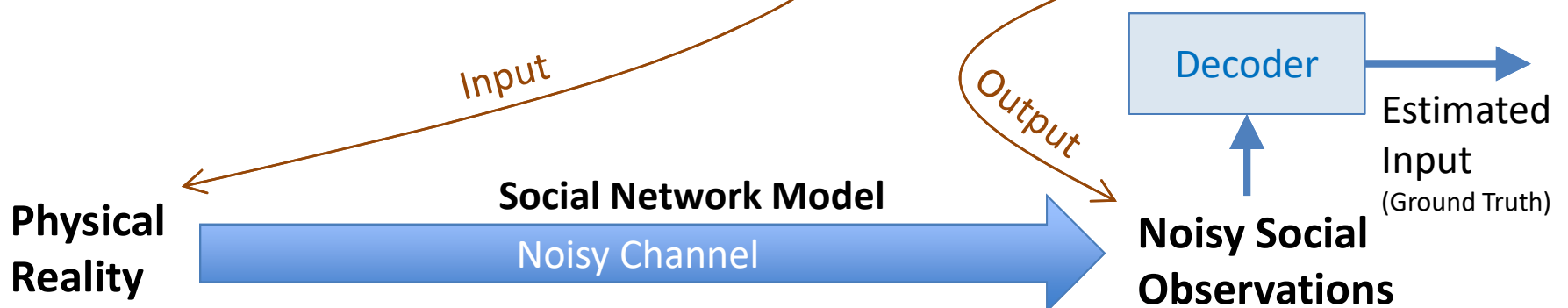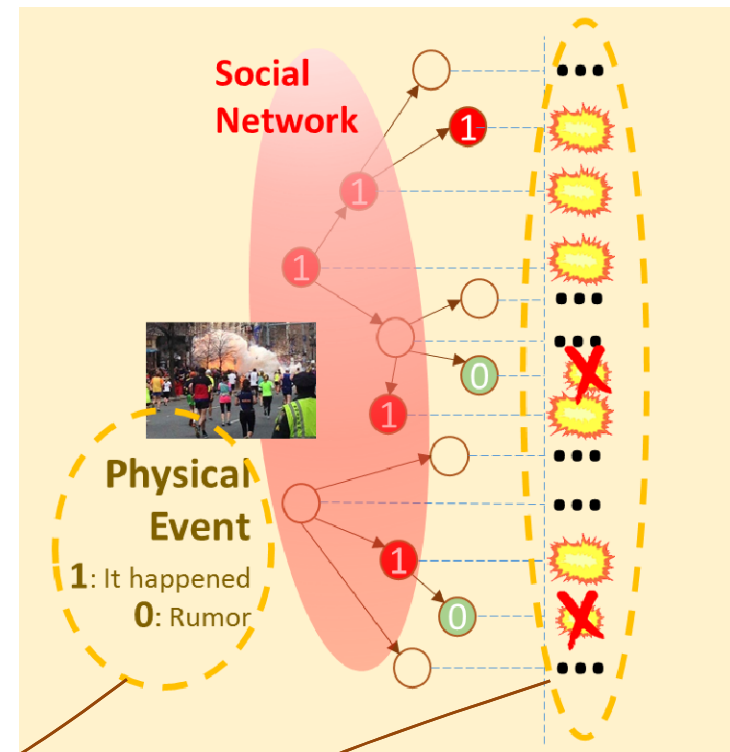3. Construct social-influence-aware fact-finding *algorithms* that approach these limits

1. Social Channel Models

2. Capacity & Error Rate Limits

3. Decoding Algorithms

Decoder

Estimated Input
(Ground Truth)

**Social Network Model**

**Physical Reality**

Noisy Channel

**Noisy Social Observations**

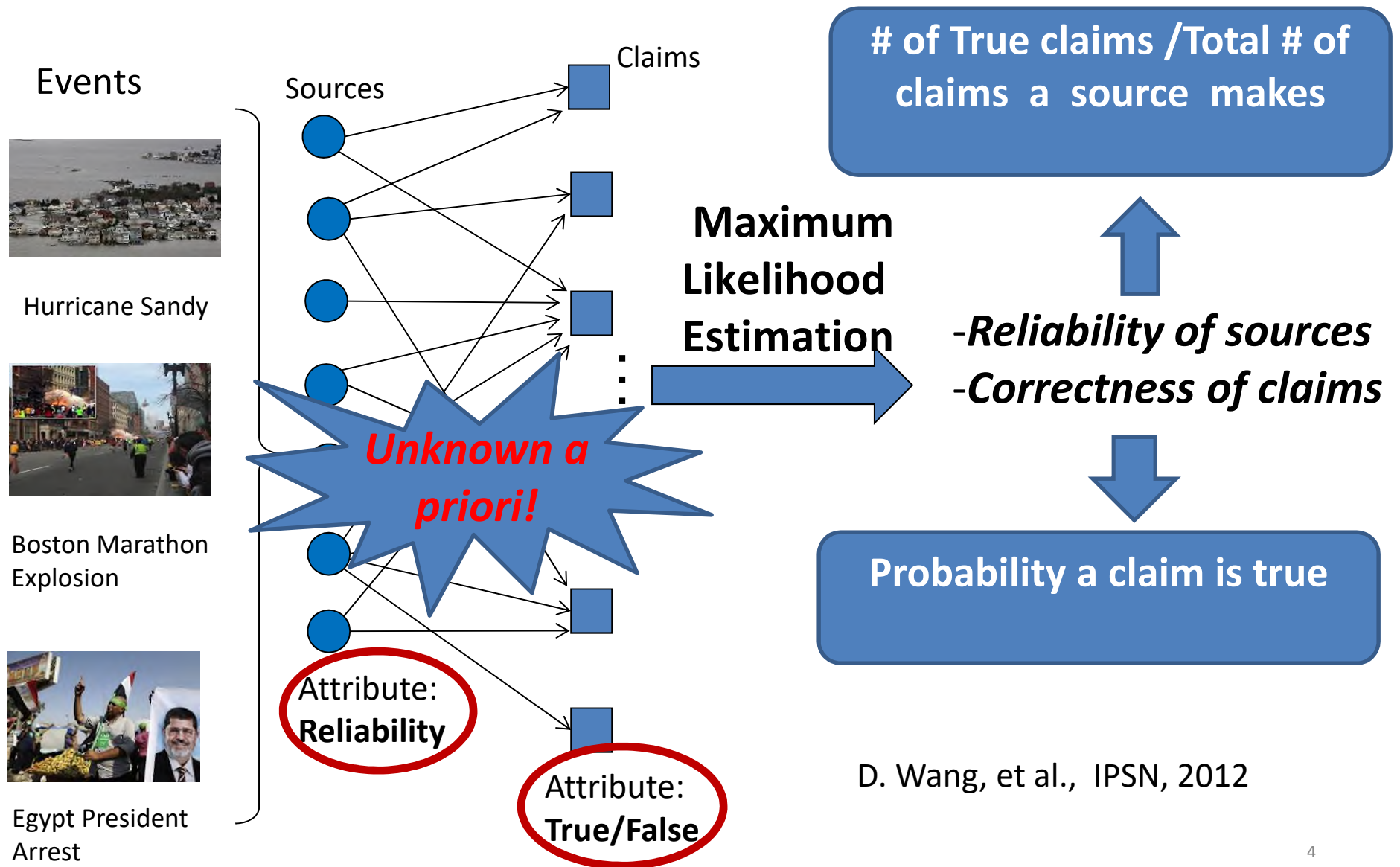4. Evaluation

# Motivation and Approach



**Approach:** Model the social network as a noisy channel that transforms "ground truth" into noisy observations

- Use information-theoretic results to understand its fundamental performance limits.

- Use estimation theory to build optimal fact-finders ("channel decoders" that approach these limits)

Social Network

Physical Event
**1**: It happened
**0**: Rumor

Input

Output

Decoder → Estimated Input (Ground Truth)

**Physical Reality**

**Social Network Model**
Noisy Channel

**Noisy Social Observations**

# Maximum Likelihood Estimation

Events

Sources

Claims

Hurricane Sandy

Boston Marathon Explosion

Egypt President Arrest

*Unknown a priori!*

Attribute: **Reliability**

Attribute: **True/False**

**Maximum Likelihood Estimation**

**# of True claims /Total # of claims a source makes**

-*Reliability of sources*
-*Correctness of claims*

**Probability a claim is true**

D. Wang, et al., IPSN, 2012

4

# Uncertain Data Provenance

**Events**

Hurricane Sandy

Boston Marathon Explosion

Egypt President Arrest

Sources

Claims

Attribute: **Reliability**

Attribute: **True/False**

Sources are not independent !
Consider the social network and source forwarding behaviors

twitter Follow me!

please retweet

facebook

D. Wang, et al., IPSN, 2014

# Formulate the Likelihood Function

Events

Sources

Claims

**SC:  Source Claim Graph**

**SD:  Social Dissemination Graph**

$$P(SC|SD, \theta) = \sum_z P(SC, z|SD, \theta)$$

Hurricane Sandy

Boston Marathon Explosion

Egypt President Arrest

Attribute: **Reliability**

Attribute: **True/False**

# Basis Definition

True Claim

False Claim

$a_i$

$b_i$

$$a_i = P(S_i C_j \mid C_j = true)$$

Using Bayesian Theorem: $a_i = \dfrac{t_i \times s_i}{d}$

where $d$ is the overal prior that a randomnly chozen claim is true

$$b_i = P(S_i C_j \mid C_j = false)$$

Using Bayesian Theorem: $b_i = \dfrac{(1 - t_i) \times s_i}{1 - d}$

where $d$ is the overal prior that a randomnly chozen claim is true

$$p_{i,k} = \dfrac{\text{number of time } S_i \text{ and } S_k \text{ make the same claim}}{\text{number of claims made by } S_k}$$
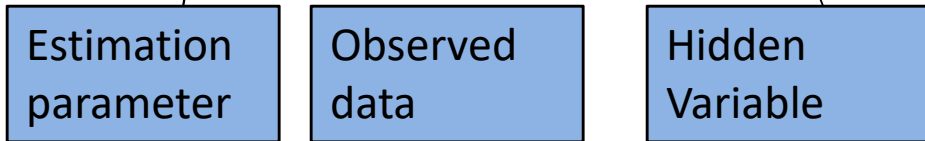
source $S_k$ is the parent node of source $S_i$ in social network
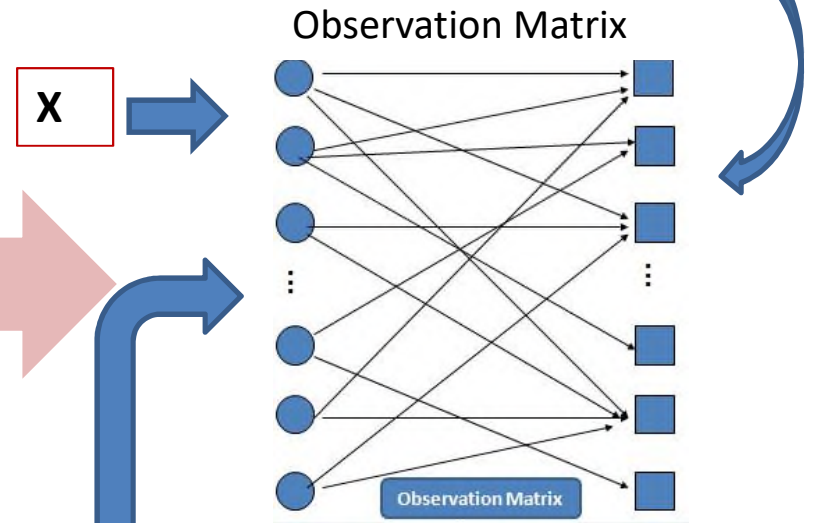
# Expectation Maximization

**Expectation Maximization**
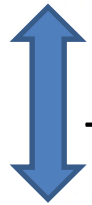
$$L(\theta; X) = p(X|\theta) = \sum_Z p(X, Z|\theta)$$

$Z = \{z_1, z_2, \ldots z_N\}$ where $z_j = 1$ when claim $C_j$ is true and 0 otherwise

| Estimation parameter | Observed data | Hidden Variable |
|---|---|---|

**X**

Observation Matrix



**Apply EM**

- Expectation Step (E-step)

$$Q\left(\theta|\theta^{(t)}\right) = E_{Z|X,\theta^{(t)}}\left[\log L(\theta; X, Z)\right]$$

- Maximization Step (M-step)

$$\theta^{(t+1)} = \underset{\theta}{\mathrm{argmax}}\, Q\left(\theta|\theta^{(t)}\right)$$

$$\theta = (a_1, a_2, \ldots a_M; b_1, b_2, \ldots b_M; d)$$

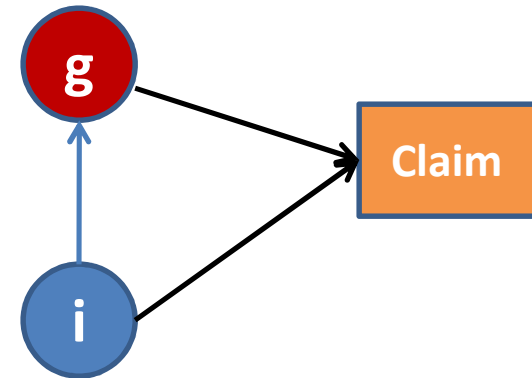**Find MLE of estimation parameter and values of hidden variables**

**Likelihood Function Incorporating Source Dependency**

$$P(SC, z|SD, \theta) = \prod_{j=1}^{N} P(z_j) \times$$

$$\left\{ \prod_{g \in M_j} P(S_g C_j | \theta, z_j) \prod_{i \in c_g} P(S_i C_j | S_g C_j) \right\}$$

$$P(z_j) = \begin{cases} d & z_j = 1 \\ (1-d) & z_j = 0 \end{cases}$$

$$P(S_g C_j | \theta, z_j) = \begin{cases} a_g & z_j = 1, S_g C_j = 1 \\ (1-a_g) & z_j = 1, S_g C_j = 0 \\ b_g & z_j = 0, S_g C_j = 1 \\ (1-b_g) & z_j = 0, S_g C_j = 0 \end{cases}$$

$$P(S_i C_j | S_g C_j) = \begin{cases} p_{ig} & S_g C_j = 1, S_i C_j = 1 \\ 1 - p_{ig} & S_g C_j = 1, S_i C_j = 0 \end{cases}$$

g

Claim

i

**Dependent Sources**

**E-Step**

$$Q\left(\theta|\theta^{(n)}\right) = \sum_{j=1}^{N}\left\{Z(n,j) \times \left[\left\{\sum_{g \in M_j}\left(\log P(S_gC_j|\theta, z_j)\right.\right.\right.$$

$$+ \sum_{i \in c_g}\log P(S_iC_j|S_gC_j)\right)\right\} + \log d\right]$$

$$+ (1 - Z(n,j)) \times \left[\left\{\sum_{g \in M_j}\left(\log P(S_gC_j|\theta, z_j)\right.\right.\right.$$

$$+ \sum_{i \in c_g}\log P(S_iC_j|S_gC_j)\right)\right\} + \log(1-d)\right]\right\} \qquad (10)$$

**M-Step**

$$a_g^{(n+1)} = a_g^* = \frac{\sum_{j \in SJ_g}Z(n,j)}{\sum_{j=1}^{N}Z(n,j)}$$

$$a_i^{(n+1)} = a_i^* = \frac{\sum_{j \in \overline{SJ}_g \cap SJ_i}Z(n,j)}{\sum_{j \in \overline{SJ}_g}Z(n,j)} \qquad \text{for } i \in c_g$$

$$b_g^{(n+1)} = b_g^* = \frac{\sum_{j \in SJ_g}(1 - Z(n,j))}{\sum_{j=1}^{N}(1 - Z(n,j))}$$

$$b_i^{(n+1)} = b_i^* = \frac{\sum_{j \in \overline{SJ}_g \cap SJ_i}(1 - Z(n,j))}{\sum_{j \in \overline{SJ}_g}(1 - Z(n,j))} \qquad \text{for } i \in c_g$$

$$d^{(n+1)} = d^* = \frac{\sum_{j=1}^{N}Z(n,j)}{N}$$

# Simple Illustrative Examples

Example 1

Example 2

**C** True Claim

**SD** Links

**SC** Links

**SD** Links that are ignored

Missing **SC** Links

**S** Source that made claim

# The Apollo Fact-finder

**http://apollo.cs.illinois.edu/**

# EM is Integrated with Apollo
## A Real World Application



**Data Collection Frontend**

**Information Analysis Frontend**

# Evaluation using Real Twitter Traces

| Trace | Hurricane Sandy | Hurricane Irene | Egypt Unrest |
|---|---|---|---|
| Time duration | 14 days (Nov.2-15, 2012) | 8 days (Aug.26-Sept.2, 2011) | 18 days (Feb.2-Feb.19,2011) |
| Locations | 16 cities in East Coasts | New York | Cairo, Egypt |
| # of users tweeted | 7,583 | 207,562 | 13,836 |
| # of tweets | 12,931 | 387,827 | 93,208 |
| # of users crawled in social network | 704,941 | 2,510,316 | 5,285,160 |
| # of follower-followee links | 37,597 | 3,902,713 | 10,490,098 |

# Estimate Latent Social Dissemination (SD) Network



**Estimate Latent Social Dissemination Network**

i → j

i  Follow  j
**FF SD**

i  retweet  j
**RT SD**

i  Epidemic Cascade*  j
**EC SD**

* P. Netrapalli and S. Sanghavi. Learning the graph of epidemic cascades. SIGMETRICS '12, pages 211–222, New York, NY, USA, 2012. ACM.

15

# Evaluation on Sandy Trace



**Understanding Source Dependency Helps !**

# Evaluation on Irene Trace

# Evaluation on Egypt Trace



**Understanding Source Dependency Helps !**

RT: Retweet

FF: Follower-Followee

EC: Epidemic Cascade

# Ground Truth Events Found by Social EM vs Regular EM
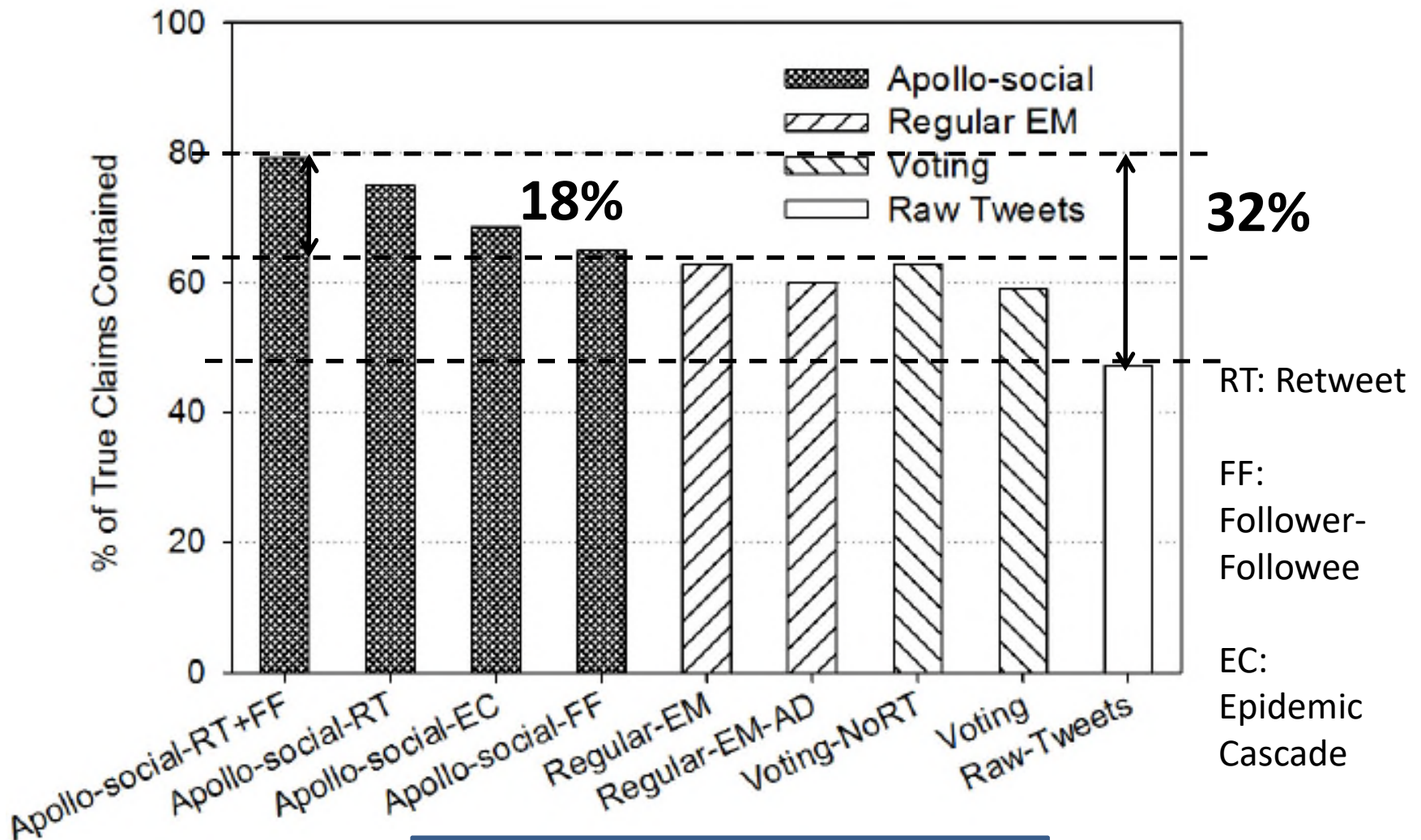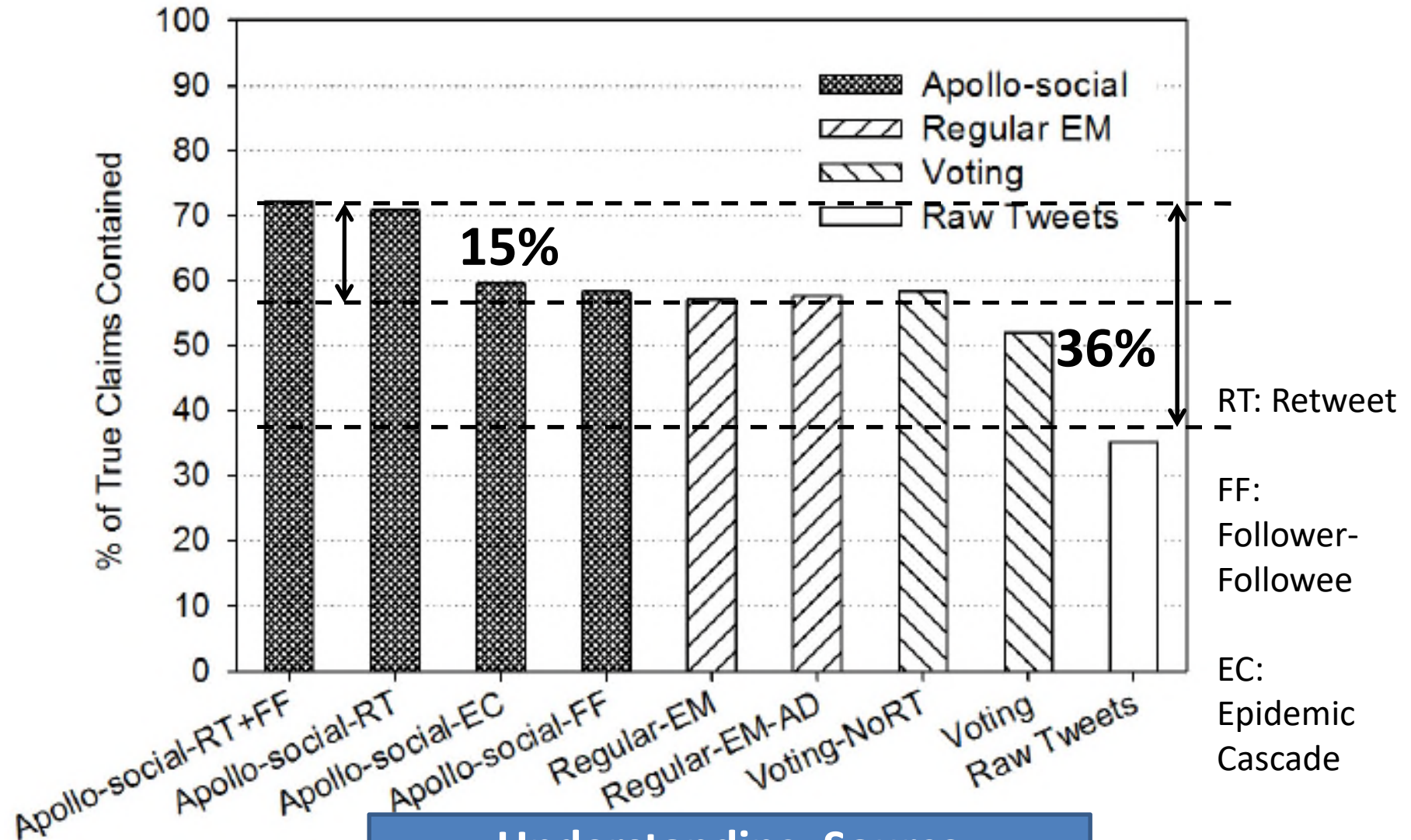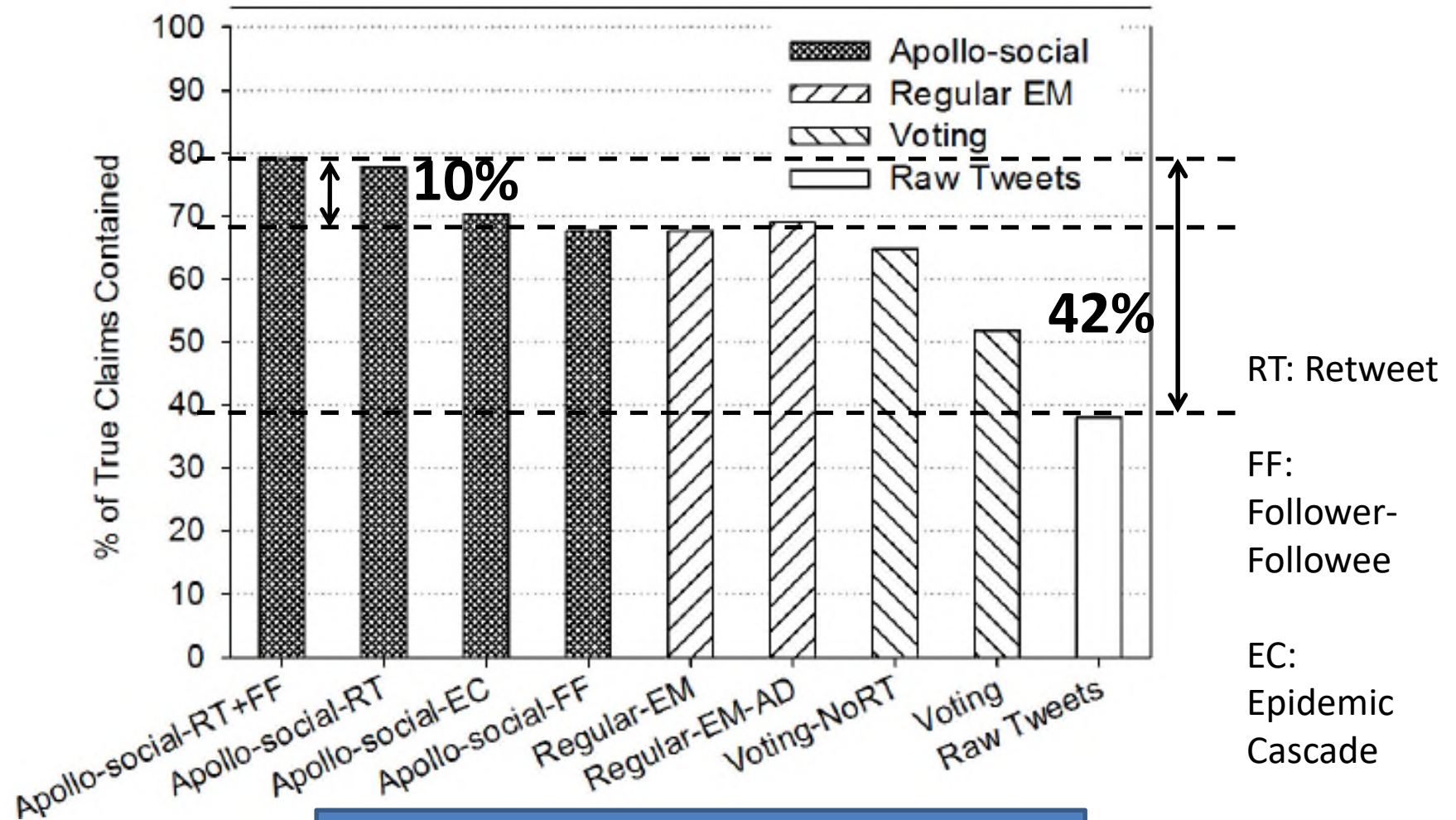
| # | Media | Tweet found by Apollo-social | Tweet found by Regular EM |
|---|-------|------------------------------|---------------------------|
| 1 | Rockland County Executive C. Scott Vanderhoef is announcing a Local Emergency Order restricting the amount of fuel that an individual can purchase at a gas station. | Rockland County Orders Restrictions on Gas Sales - Nyack-Piermont, NY Patch http://t.co/cDSrqpa2 | **MISSING** |
| 2 | New York City Mayor Michael Bloomberg has announced that the city will impose an indefinite program of gas rationing after fuel shortages led to long lines and frustration at the pump in the wake of superstorm Sandy. | Gas rationing plan set for New York City: The move follows a similar announcement last week in New Jersey to eas... http://t.co/nkmF7U9I | RT @nytimes: Breaking News: Mayor Bloomberg Imposes Odd-Even Gas Rationing Starting Friday, as Does Long Island http://t.co/eax7KMVi |
| 3 | New Jersey authorities filed civil suits Friday accusing seven gas stations and one hotel of price gouging in the wake of Hurricane Sandy. | RT @MarketJane: NJ plans price gouging suits against 8 businesses. They include gas stations and a lodging provider. | **MISSING** |
| 4 | The rationing system: restricting gas sales to cars with even-numbered license plates on even days, and odd-numbered on odd days will be discontinued at 6 a.m. Tuesday, Gov. Chris Christie announced on Monday. | # masdirin City Room: Gas Rationing in New Jersey to End Tuesday # news | RT @nytimes: City Room: Gas Rationing in New Jersey to End Tuesday http://t.co/pYIVOmPo |
| 5 | New Yorkers can expect gas rationing for at least five more days: Bloomberg. | Mayor Bloomberg: Gas rationing in NYC will continue for at least 5 more days. @eyewitnessnyc #SandyABC7 | Bloomberg: Gas Rationing To Stay In Place At Least Through The Weekend http://t.co/mmqqjYRx |

TABLE III. GROUND TRUTH EVENTS AND RELATED CLAIMS FOUND BY APOLLO-SOCIAL VS REGULAR EM IN SANDY

# One Interesting Example



**Shark in the street!**

**Suppressed by Social EM**

The Washington Post

Posted at 08:53 AM ET, 08/26/2011

**Hurricane Irene: 'Photo' of shark swimming in street is fake**

By Sarah Anne Hughes

*FAKE!*

Holy moly! A (fake) picture of a shark swimming on a Puerto Rico street! (Reddit)

http://www.washingtonpost.com/blogs/blogpost/post/hurricane-irene-photo-of-shark-swimming-in-street-is-fake/2011/08/26/gIQABHAvfJ_blog.html
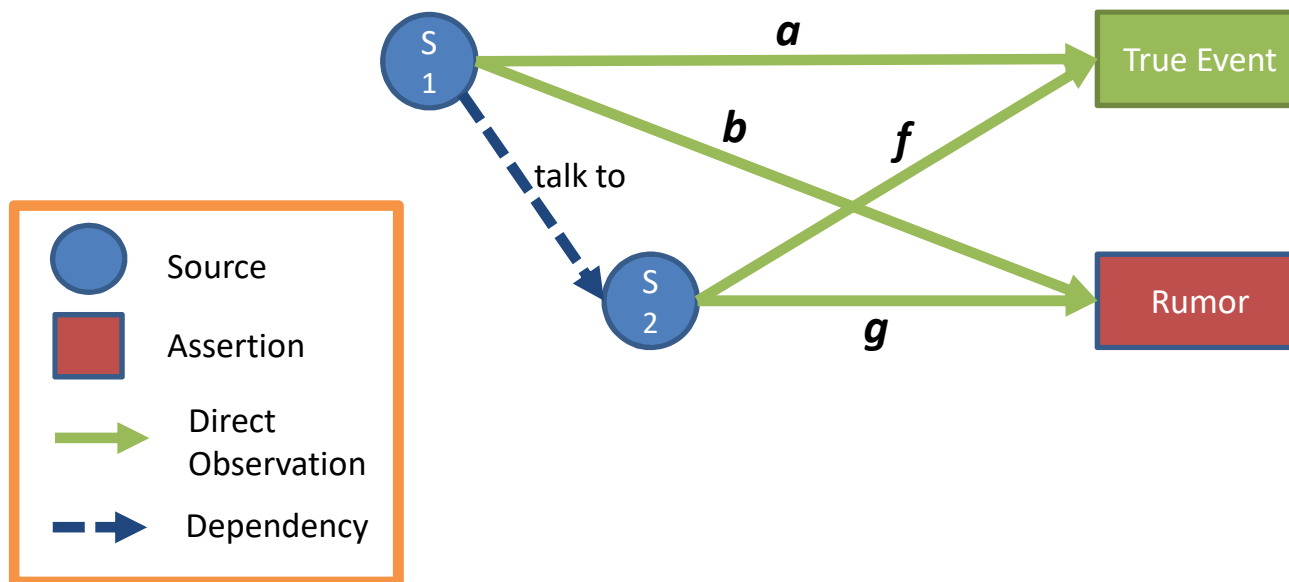
# Social Sensing: Source Dependencies

- Failure of physical sensor: independent

- Failure of social sensing sensor: dependent
  - People talk and influence each other
  - Correlated errors

- We need to formulate source dependency correctly!

# Estimator Parameters

# Expectation-Maximization Solution

E Step

$$\mathcal{Q}(\theta|\theta^{(t)}) = \sum_{j=1}^{m} P(C_j|SC_j; \theta^{(t)}) \sum_{C_j \in \{0,1\}} \ln(P(C_j; \theta))$$
$$\left( \sum_{i=1}^{n} \ln(P(S_iC_j|C_j; \theta, D_{ij})) \right)$$

M Step

$$a_i^{(t+1)} = \frac{\sum_{C_j \in S_iC_1^{D_0}} P(C_j = 1|S_iC_j; \theta^{(t)})}{\sum_{C_j \in S_iC_1^{D_0} \bigcup S_iC_0^{D_0}} P(C_j = 1|S_iC_j; \theta^{(t)})}$$

$$f_i^{(t+1)} = \frac{\sum_{C_j \in S_iC_1^{D_1}} P(C_j = 1|S_iC_j; \theta^{(t)})}{\sum_{C_j \in S_iC_1^{D_1} \bigcup S_iC_0^{D_1}} P(C_j = 1|S_iC_j; \theta^{(t)})}$$

$$b_i^{(t+1)} = \frac{\sum_{C_j \in S_iC_1^{D_0}} P(C_j = 0|S_iC_j; \theta^{(t)})}{\sum_{C_j \in S_iC_1^{D_0} \bigcup S_iC_0^{D_0}} P(C_j = 0|S_iC_j; \theta^{(t)})}$$
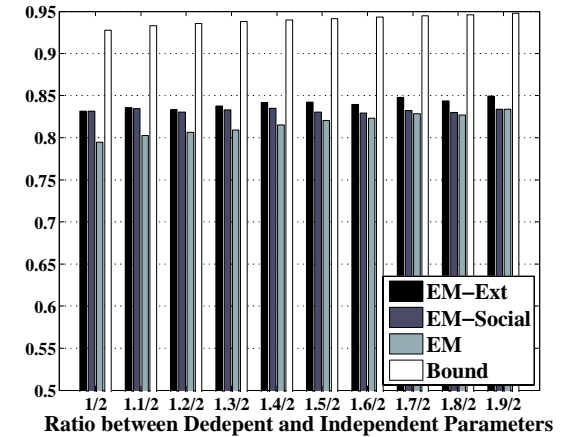
$$g_i^{(t+1)} = \frac{\sum_{C_j \in S_iC_1^{D_1}} P(C_j = 0|S_iC_j; \theta^{(t)})}{\sum_{C_j \in S_iC_1^{D_1} \bigcup S_iC_0^{D_1}} P(C_j = 0|S_iC_j; \theta^{(t)})}$$

$$z^{(t+1)} = \frac{\sum_{j=1}^{m} P(C_j = 0|S_iC_j; \theta^{(t)})}{m}$$

# *Simulation of Dependency-Aware Estimator*

# *Empirical Evaluation*

# Optimal Estimator

**Social Network**

**Physical Event**

**1**: True Event
**0**: False

$C_j$

**Prediction**

$P(C_j\ happens|Observations)$

$P(C_j\ is\ a\ rumor|Observations)$

**Error**

Source (silent)

Source (not silent)

Assertion

Source Dependency

Claim

# Error Bounds

- Alleged event: "France bombs Iraq"

| Pravda | Jazeera | Ahram | Falsehood Probability | Truth Probability |
|--------|---------|-------|-----------------------|-------------------|
| Silent | Silent | Silent | 99% | 1% |
| Silent | Silent | Report | 80% | 20% |
| Silent | Report | Silent | 90% | 10% |
| Silent | Report | Report | 40% | 60% |
| Report | Silent | Silent | 95% | 5% |
| Report | Silent | Report | 60% | 40% |
| Report | Report | Silent | 70% | 30% |
| Report | Report | Report | 5% | 95% |

# Error Bounds

- Alleged event: "France bombs Iraq"

| The odds | Pravda | Jazeera | Ahram | Falsehood Probability | Truth Probability |
|---|---|---|---|---|---|
| 4% | Silent | Silent | Silent | 99% | 1% |
| 10% | Silent | Silent | Report | 80% | 20% |
| 10% | Silent | Report | Silent | 90% | 10% |
| 20% | Silent | Report | Report | 40% | 60% |
| 20% | Report | Silent | Silent | 95% | 5% |
| 13% | Report | Silent | Report | 60% | 40% |
| 13% | Report | Report | Silent | 70% | 30% |
| 10% | Report | Report | Report | 5% | 95% |

# Error Bounds

- Alleged event: "France bombs Iraq"

| The odds | Pravda | Jazeera | Ahram | Falsehood Probability | Truth Probability |
|----------|--------|---------|-------|-----------------------|-------------------|
| 4% | Silent | Silent | Silent | 99% | 1% |
| 10% | Silent | Silent | Report | 80% | 20% |
| 10% | Silent | Report | Silent | 90% | 10% |
| 20% | Silent | Report | Report | 40% | 60% |
| 20% | Report | Silent | Silent | 95% | 5% |
| 13% | Report | Silent | Report | 60% | 40% |
| 13% | Report | Report | Silent | 70% | 30% |
| 10% | Report | Report | Report | 5% | 95% |

Odds of omission = 4%

# Error Bounds

- Alleged event: "France bombs Iraq"

| The odds | Pravda | Jazeera | Ahram | Falsehood Probability | Truth Probability |
|----------|--------|---------|-------|----------------------|-------------------|
| 4% | Silent | Silent | Silent | 99% | 1% |
| 10% | Silent | Silent | Report | 80% | 20% |
| 10% | Silent | Report | Silent | 90% | 10% |
| 20% | Silent | Report | Report | 40% | 60% |
| 20% | Report | Silent | Silent | 95% | 5% |
| 13% | Report | Silent | Report | 60% | 40% |
| 13% | Report | Report | Silent | 70% | 30% |
| 10% | Report | Report | Report | 5% | 95% |

Odds of omission = 4%

Odds of error = 0.1 * 0.2 + …

# Error Bounds

- Alleged event: "France bombs Iraq"

| The odds | Pravda | Jazeera | Ahram | Falsehood Probability | Truth Probability |
|----------|--------|---------|-------|-----------------------|-------------------|
| 4% | Silent | Silent | Silent | 99% | 1% |
| 10% | Silent | Silent | Report | 80% | 20% |
| 10% | Silent | Report | Silent | 90% | 10% |
| 20% | Silent | Report | Report | 40% | 60% |
| 20% | Report | Silent | Silent | 95% | 5% |
| 13% | Report | Silent | Report | 60% | 40% |
| 13% | Report | Report | Silent | 70% | 30% |
| 10% | Report | Report | Report | 5% | 95% |

Odds of omission = 4%

Odds of error = 0.1 * 0.2 + 0.1 * 0.1 + …

# Error Bounds

- Alleged event: "France bombs Iraq"

| The odds | Pravda | Jazeera | Ahram | Falsehood Probability | Truth Probability |
|---|---|---|---|---|---|
| 4% | Silent | Silent | Silent | 99% | 1% |
| 10% | Silent | Silent | Report | 80% | 20% |
| 10% | Silent | Report | Silent | 90% | 10% |
| 20% | Silent | Report | Report | 40% | 60% |
| 20% | Report | Silent | Silent | 95% | 5% |
| 13% | Report | Silent | Report | 60% | 40% |
| 13% | Report | Report | Silent | 70% | 30% |
| 10% | Report | Report | Report | 5% | 95% |

Odds of omission = 4%

Odds of error = 0.1 * 0.2 + 0.1 * 0.1 + 0.2 * 0.4 + …
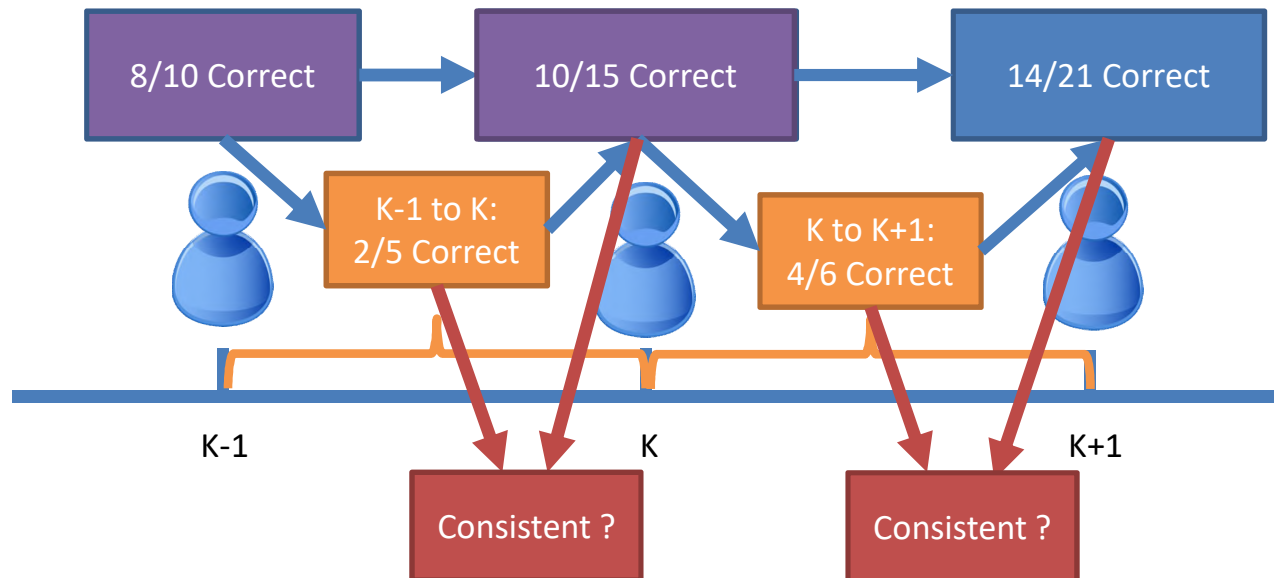
# Error Bounds

- Alleged event: "France bombs Iraq"

| The odds | Pravda | Jazeera | Ahram | Falsehood Probability | Truth Probability |
|----------|--------|---------|-------|----------------------|-------------------|
| 4% | Silent | Silent | Silent | 99% | 1% |
| 10% | Silent | Silent | Report | 80% | 20% |
| 10% | Silent | Report | Silent | 90% | 10% |
| 20% | Silent | Report | Report | 40% | 60% |
| 20% | Report | Silent | Silent | 95% | 5% |
| 13% | Report | Silent | Report | 60% | 40% |
| 13% | Report | Report | Silent | 70% | 30% |
| 10% | Report | Report | Report | 5% | 95% |

Odds of omission = 4%
Odds of error = 0.1 * 0.2 + 0.1 * 0.1 + 0.2 * 0.4 + 0.2 * 0.05 + 0.4 * 0.13 + 0.3 * 0.13 + 0.05 * 0.1
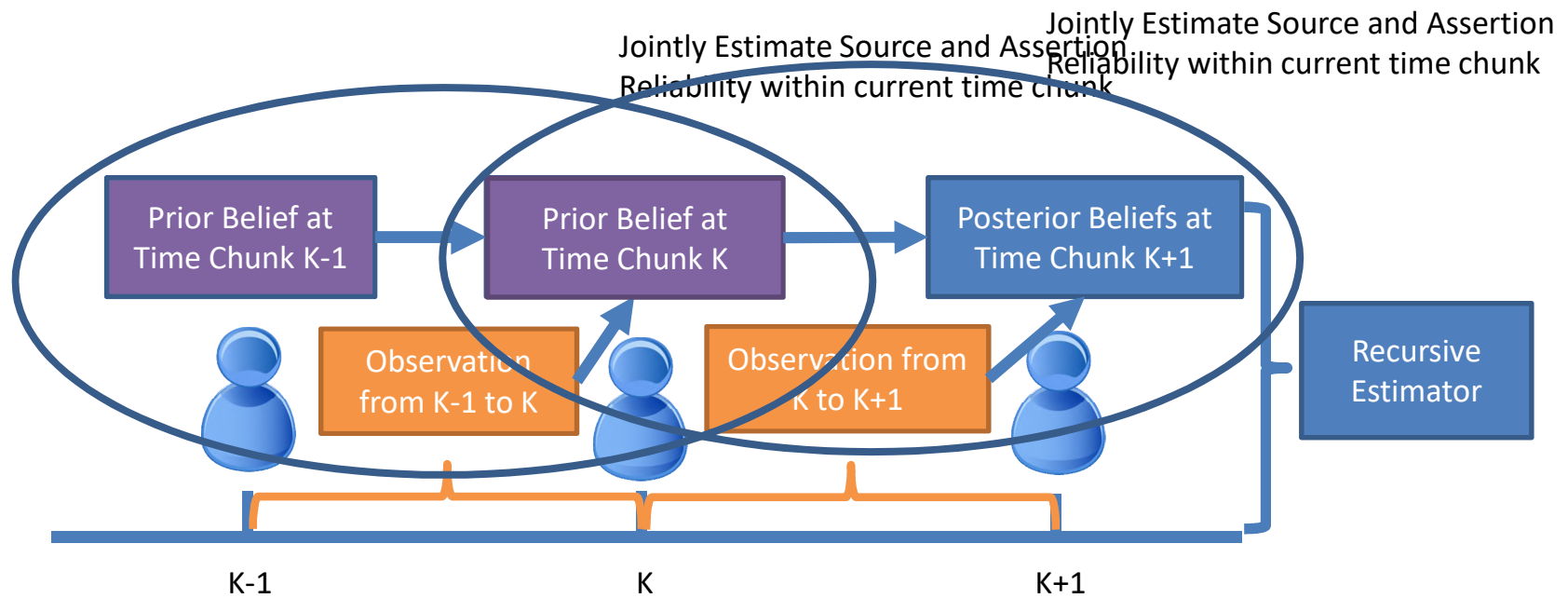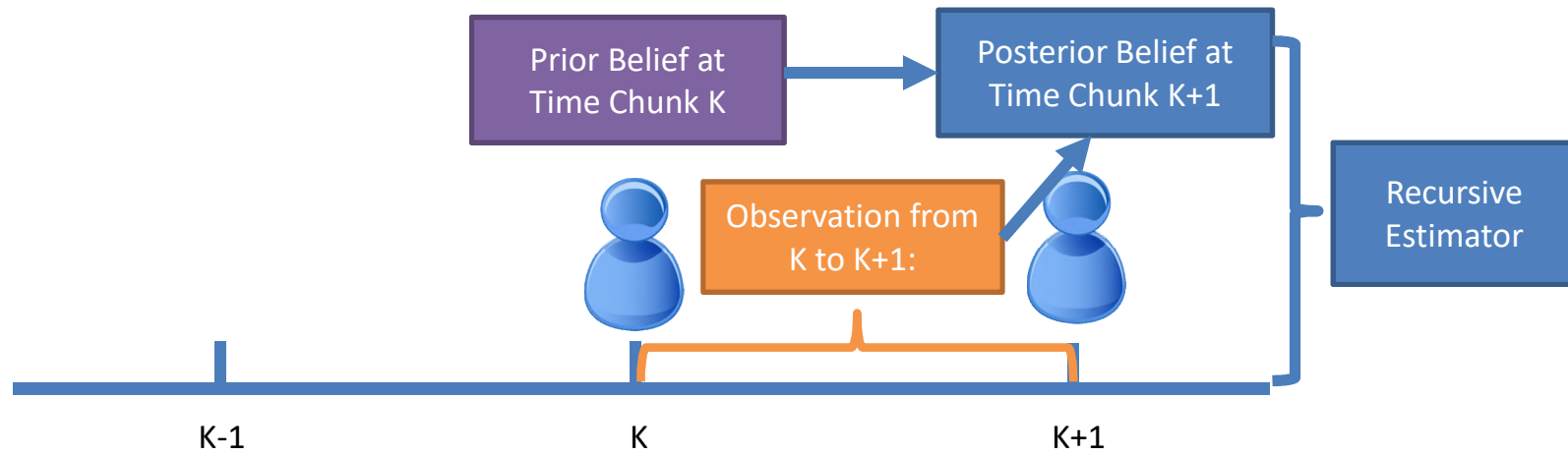       = 23.6%

# Example
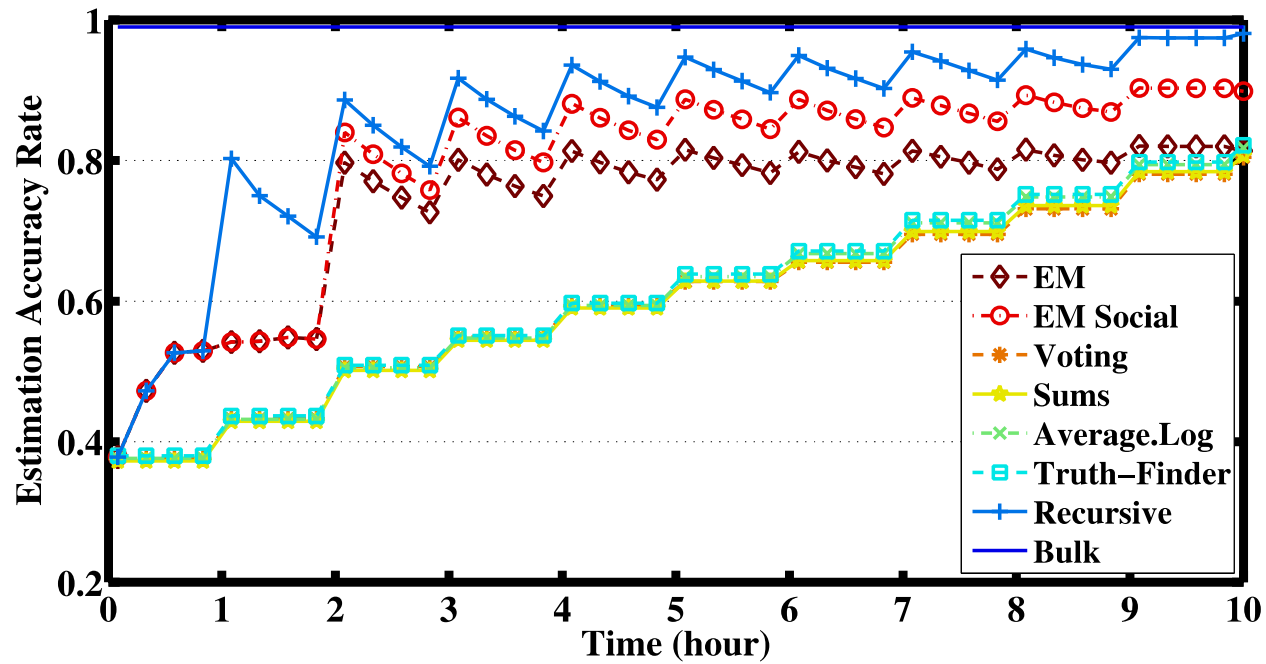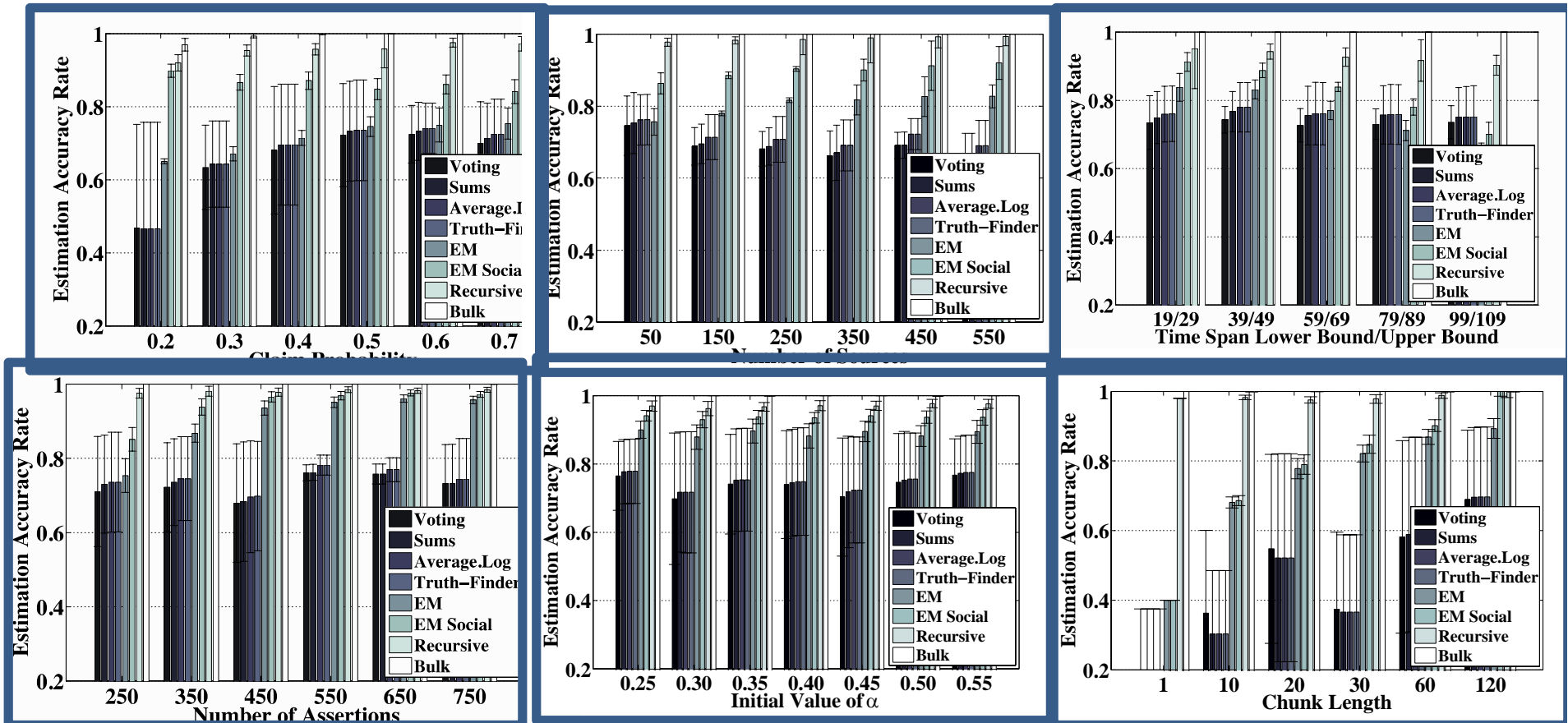
# Overview

# Recursive Estimator

# Recursive Estimator

- Recursively update the belief of reliability distribution:
    - Compute mean reliability (Compute $1^{st}$ Moment)
        - source reliability parameters, $\theta_i$
        - probability of correctness, $P(t(C) = 1 | SC_k, D, \theta)$
    - Computing the error variance (Compute $2^{nd}$ Moment)
        - error variance of source reliability parameters, $\theta_i$
    - Computing the posterior belief (Update Distribution with Moment Matching)
        - updated belief in source reliability

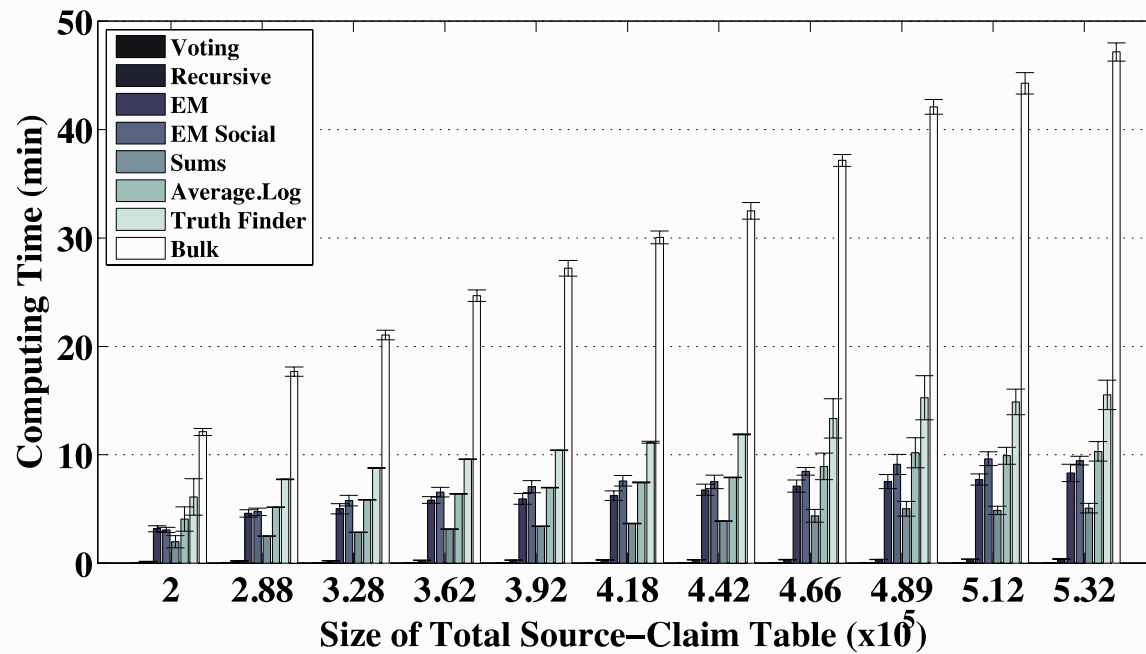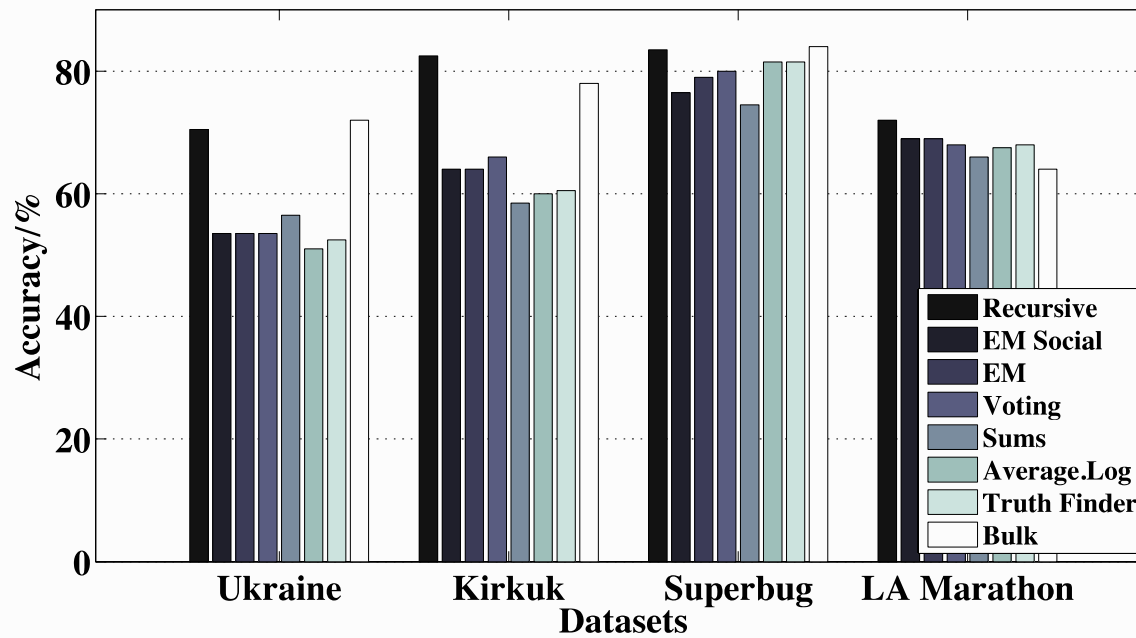# Synthetic Data: estimation accuracy of 10-hour trace

# Synthetic Data: computation time

# Empirical Evaluation: Empirical Accuracy Results

# Empirical Evaluation: empirical execution time