

# Microsoft Research

Each year Microsoft Research hosts hundreds of influential speakers from around the world including leading scientists, renowned experts in technology, book authors, and leading academics, and makes videos of these lectures freely available.

2013 © Microsoft Corporation. All rights reserved.

# *Small, $n=me$ , data*

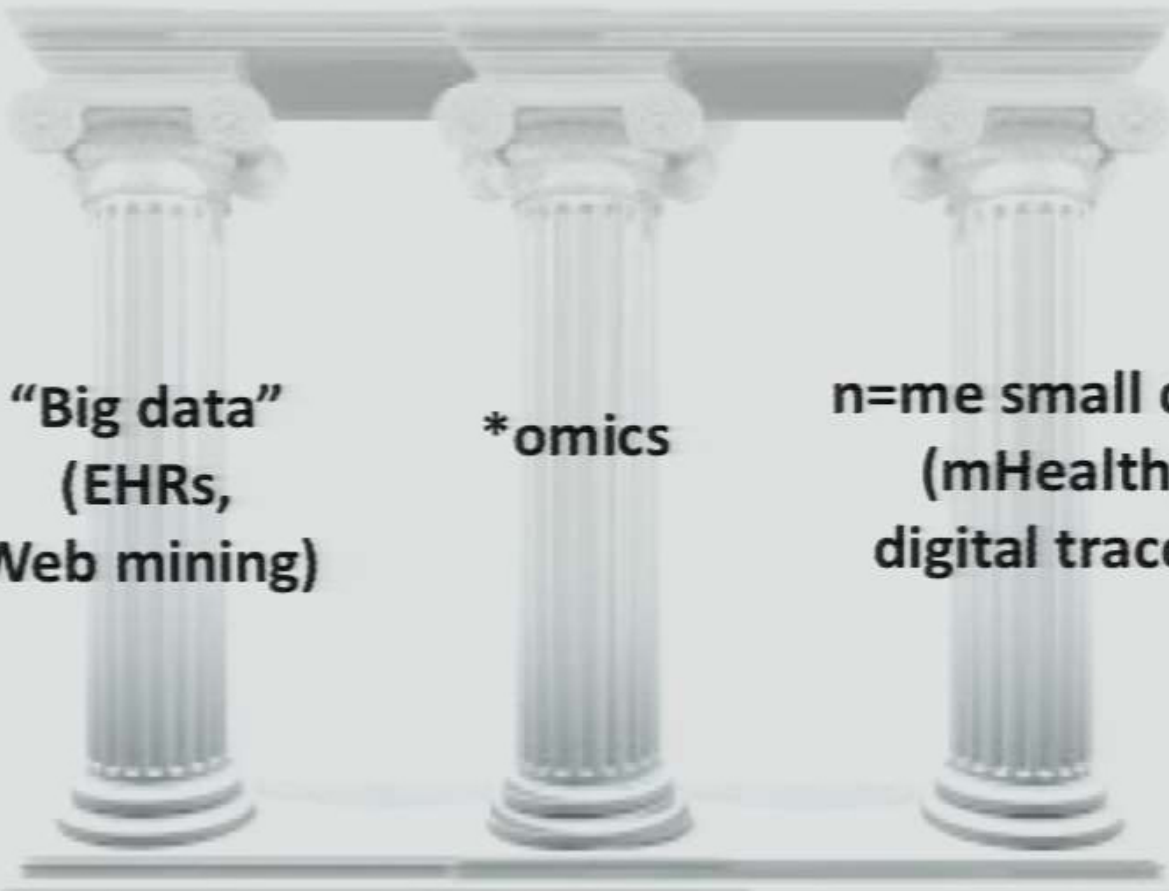
Deborah Estrin

Professor, Computer Science, Cornell NYC Tech  
Professor, Public Health, Weill Cornell Medical College  
Co-founder, Open mHealth

[destrin@cs.cornell.edu](mailto:destrin@cs.cornell.edu)

work done with collaborators from  
Cornell, UCLA, [openmhealth.org](http://openmhealth.org), ...

## Personalized, precision, medicine



**“Big data”  
(EHRs,  
Web mining)**

**\*omics**

**n=me small data  
(mHealth,  
digital traces)**

*harness previously-unmeasured function and behavior  
to fuel personalized and evidence-producing care*

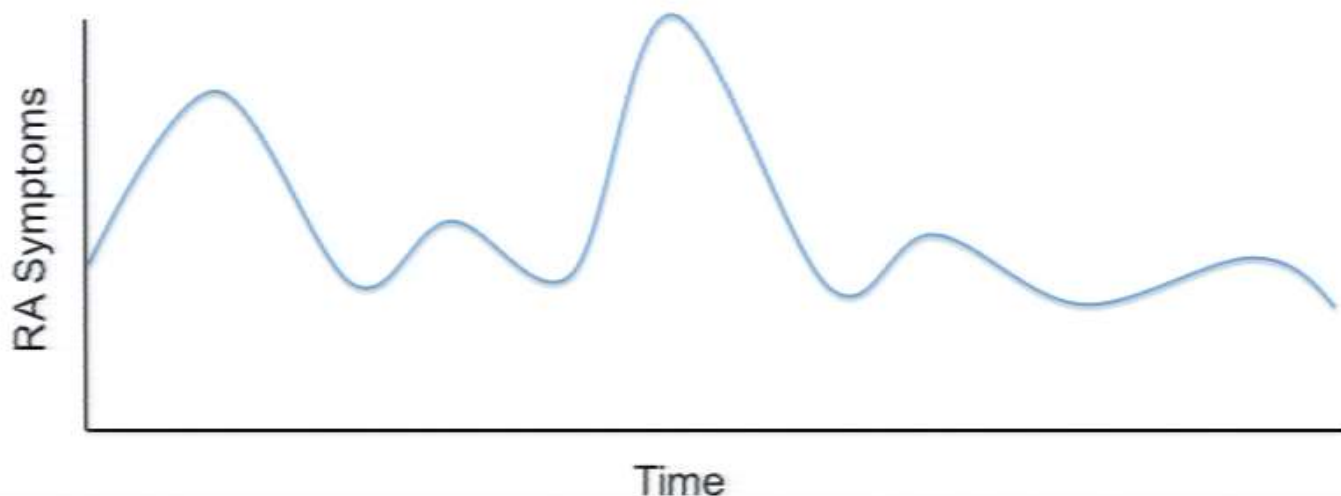
## A patient with arthritis



# Understanding fluctuations in RA disease activity

(D. Orange, R. Darnell, et al)

- 80% of RA patients experience relapsing-remitting course
  - unpredictable flares limit patient function, work productivity
- Treatment often has unwanted side effects (short courses of steroids: insomnia, hypertension, glucose intolerance, ..)
- “biologic agents” (TNF inhibitors) only make 60% of patients >20% better; 40%, >50% better; 15%, >70% better
- Early diagnosis and treatment leads to improved outcomes





### **Participant self-care**

*How is this new medication working for me?*



### **Clinical care**

*How is the patient responding to new care plan?*



### **Research evidence**

*What works best in different contexts?*



## Profound potential to rephrase 'does it work?'

(Complexes of)  
Exposures

strength of association?

Outcome

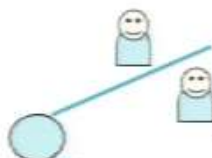
individual



population

'does it work on average?' (RCT)

100 people



50 people



Outcome measure

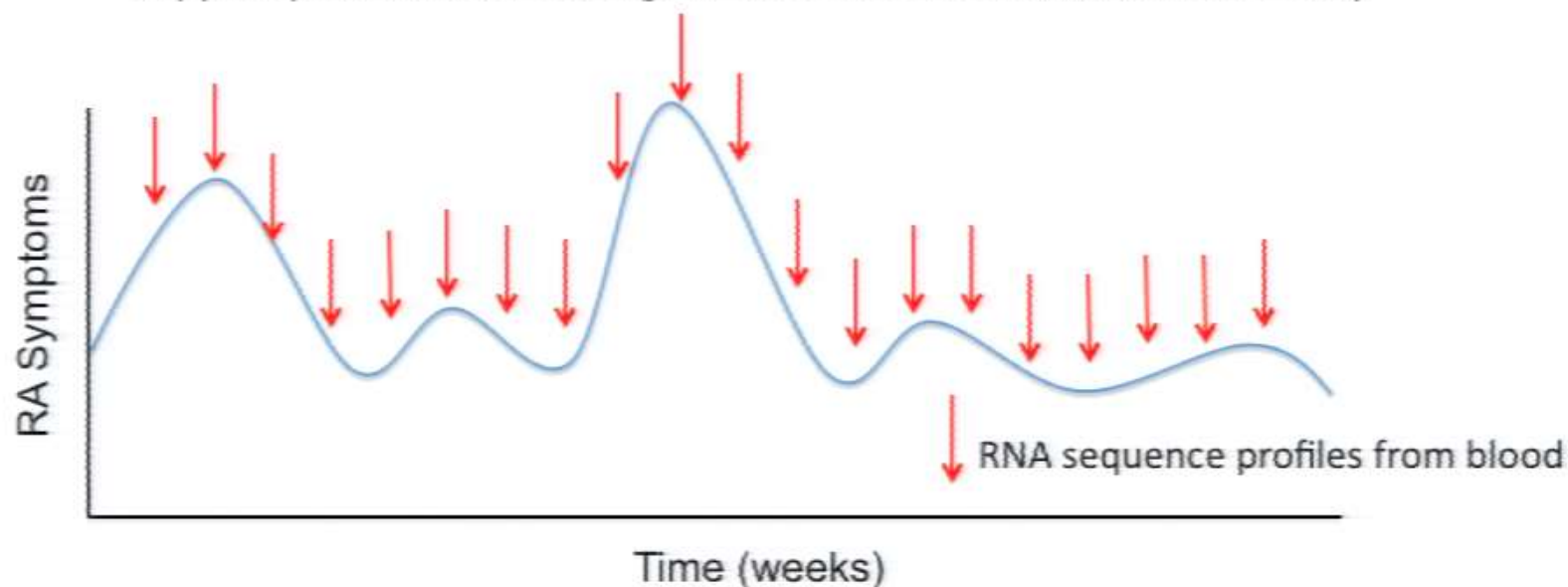
50 people



Outcome measure  
population

## Using bioinformatics to understand fluctuations in RA disease activity (Orange, Darnell, et al)

- Frequent blood samples taken and mailed in to support analysis pre, during and after flare
- mHealth:
  - make such studies practical
  - support adaptive sampling (sample more when signs of flare begin)?
  - support personalized management based on mechanism discovery





## Data Processing

(transform, cluster, infer, viz)



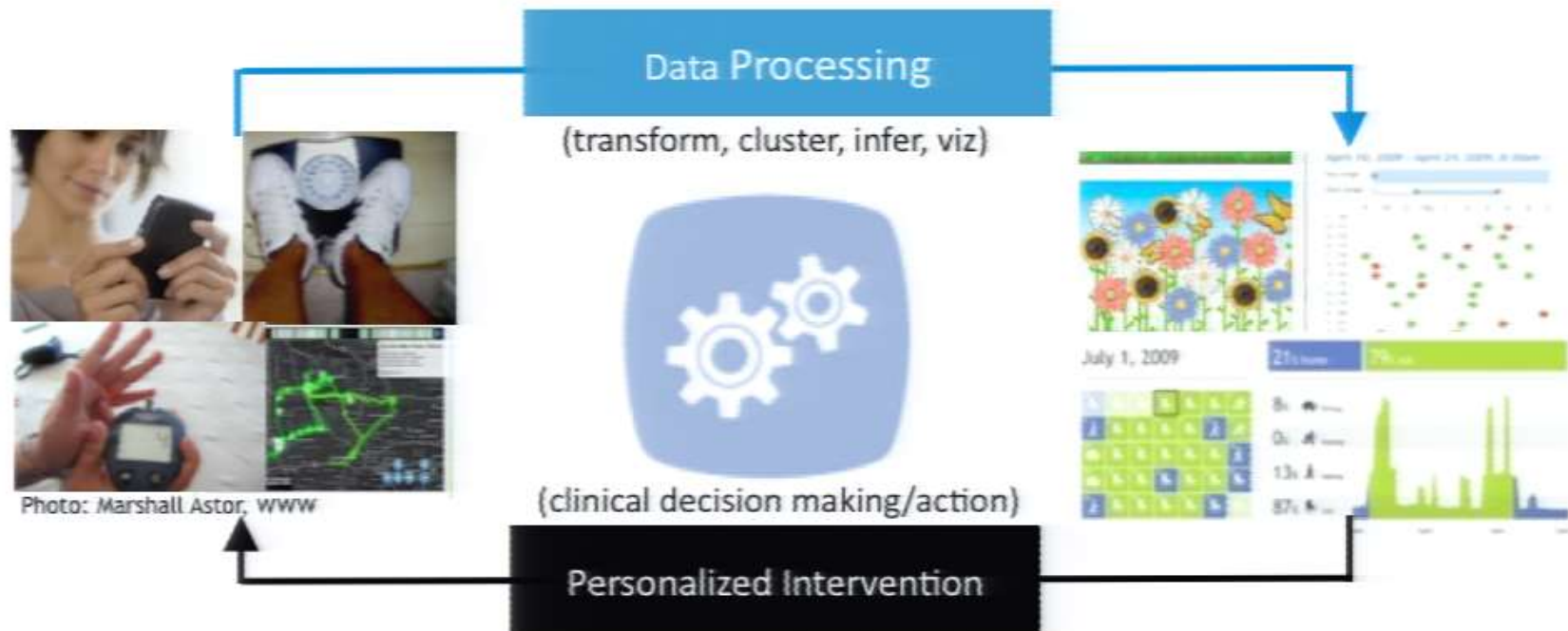
(clinical decision making/action)

## Personalized Intervention



Photo: Marshall Astor, [www](http://www)

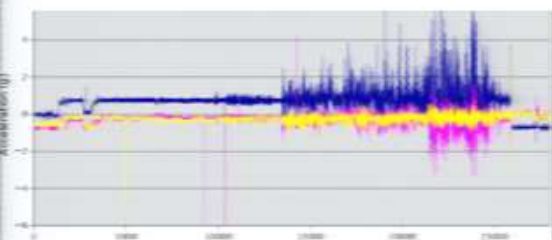
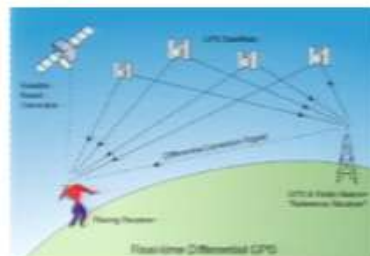
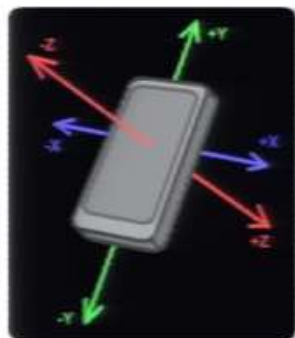




***ultimate goal: create scalable tools for PCP and clinic based care that supports precision, personalization, and continuity of care in:***

- RA, Lupus, Crohns, Asthma, MS
- Hospital discharge, Surgical recovery
- Pain management, Chronic fatigue, Migraines
- Depression, ADHD, insomnia, post-traumatic stress disorder
- Integrative medicine effectiveness
- Behavior change (individual, family , community)

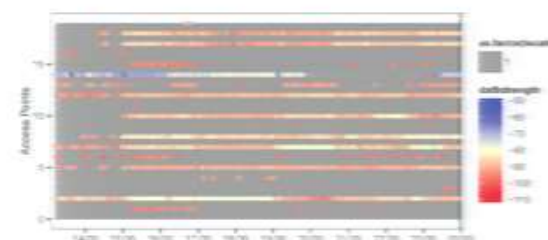
# Passively-recorded activity and location traces



Accelerometry



GPS Data



Ambient Wi-Fi Signals

## Data Processing

(transform, cluster, infer, viz)



(clinical decision making/action)

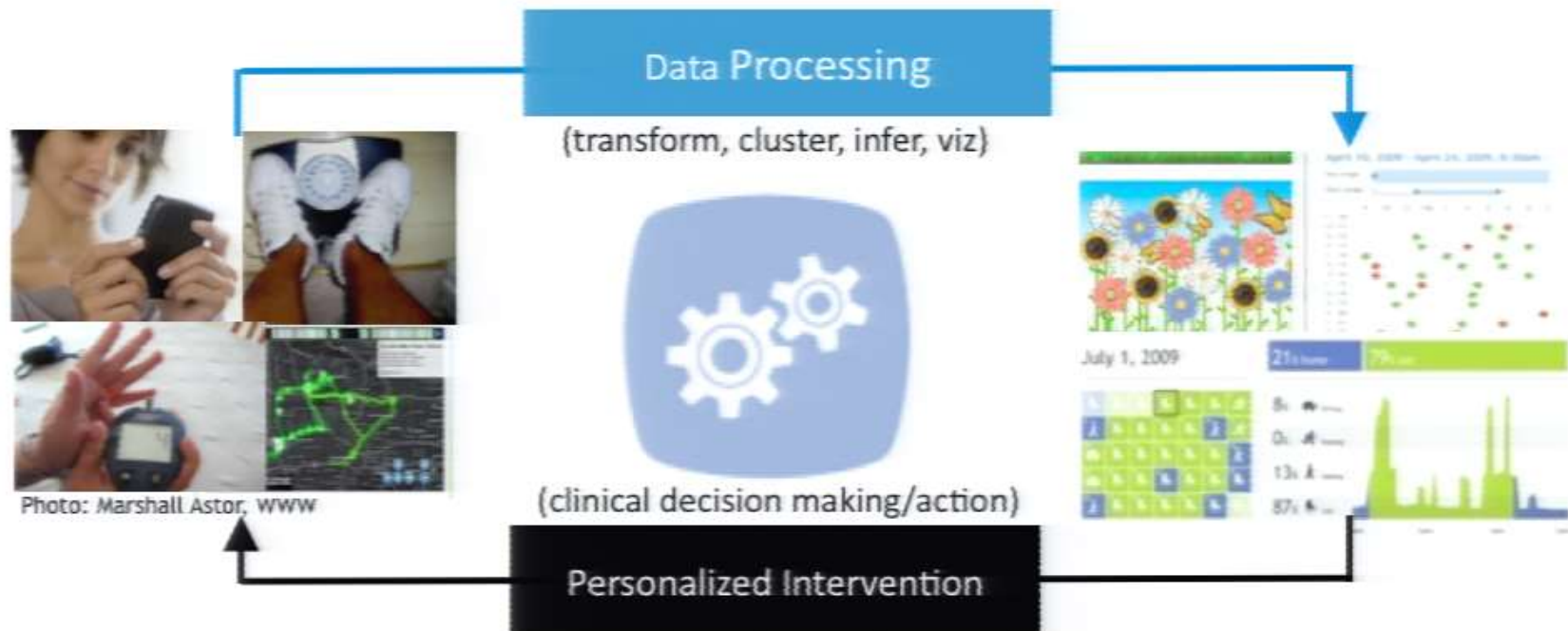
## Personalized Intervention



Photo: Marshall Astor, [www](http://www)







***ultimate goal: create scalable tools for PCP and clinic based care that supports precision, personalization, and continuity of care in:***

- RA, Lupus, Crohns, Asthma, MS
- Hospital discharge, Surgical recovery
- Pain management, Chronic fatigue, Migraines
- Depression, ADHD, insomnia, post-traumatic stress disorder
- Integrative medicine effectiveness
- Behavior change (individual, family, community)

KaurKremerMullainathan  
SelfControl.pdf

registry.openmealth.org.d  
ocumentation.wheel.graffle

KaurKremerMullainathan  
AERP&P2010.pdf

mPire One  
Page-092813-2.doc

AnnualReviewEdits4.pdf

mHealth-  
Greenhouse...813V2.docx  
148 KB

Precommitment and  
Flexibility 1992.pdf

mPire One  
Page-092813.doc

99034103.pdf

20130913\_173401  
copy.pdf

13312.pdf

ML01-Introduction.pdf

PDF14978959.pdf

347.pdf

AJMC June 2013 UHG  
format.ppt.pptx  
680 KB

raymerieeee2003.pdf

Alshaban and Taher.pdf

clockwise-  
aug-29-2013.pdf



## Data Processing

(transform, cluster, infer, viz)



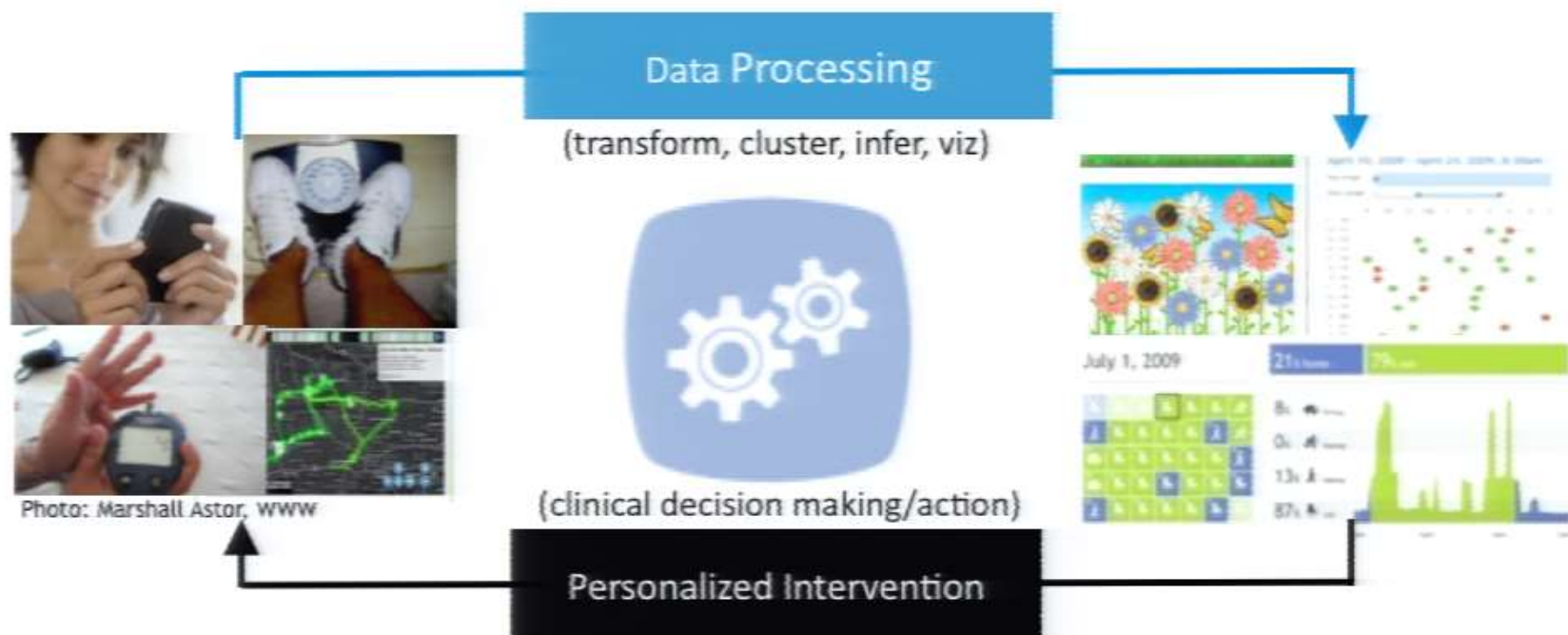
(clinical decision making/action)

## Personalized Intervention



Photo: Marshall Astor, [www](http://www)





***ultimate goal: create scalable tools for PCP and clinic based care that supports precision, personalization, and continuity of care in:***

- RA, Lupus, Crohns, Asthma, MS
- Hospital discharge, Surgical recovery
- Pain management, Chronic fatigue, Migraines
- Depression, ADHD, insomnia, post-traumatic stress disorder
- Integrative medicine effectiveness
- Behavior change (individual, family, community)



KaurKremerMullainathan  
SelfControl.pdf

registry.openmealth.org.d  
ocumentation.wheel.graffle

KaurKremerMullainathan  
AERP&P2010.pdf

mPire One  
Page-092813-2.doc

AnnualReviewEdits4.pdf

mHealth-  
Greenhouse...813V2.docx  
148 KB

Precommitment and  
Flexibility 1992.pdf

mPire One  
Page-092813.doc

99034103.pdf

20130913\_173401  
copy.pdf

13312.pdf

ML01-Introduction.pdf

PDF14978959.pdf

347.pdf

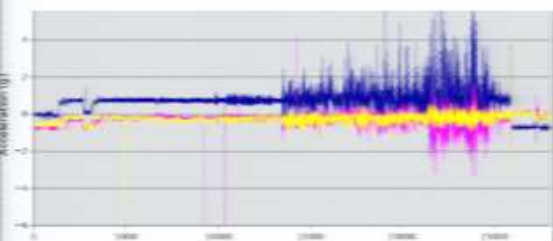
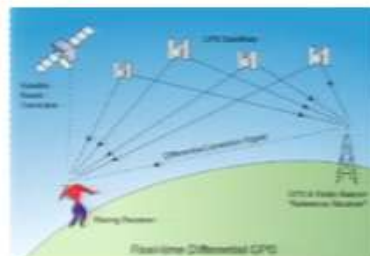
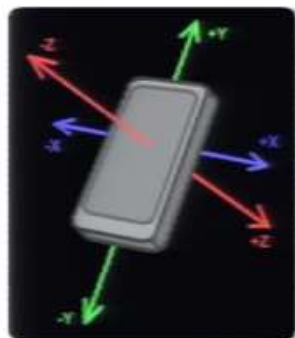
AJMC June 2013 UHG  
format.ppt.pptx  
680 KB

raymerieeee2003.pdf

Alshaban and Taher.pdf

clockwise-  
aug-29-2013.pdf

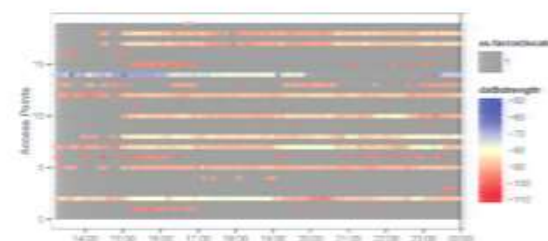
# Passively-recorded activity and location traces



Accelerometry



GPS Data



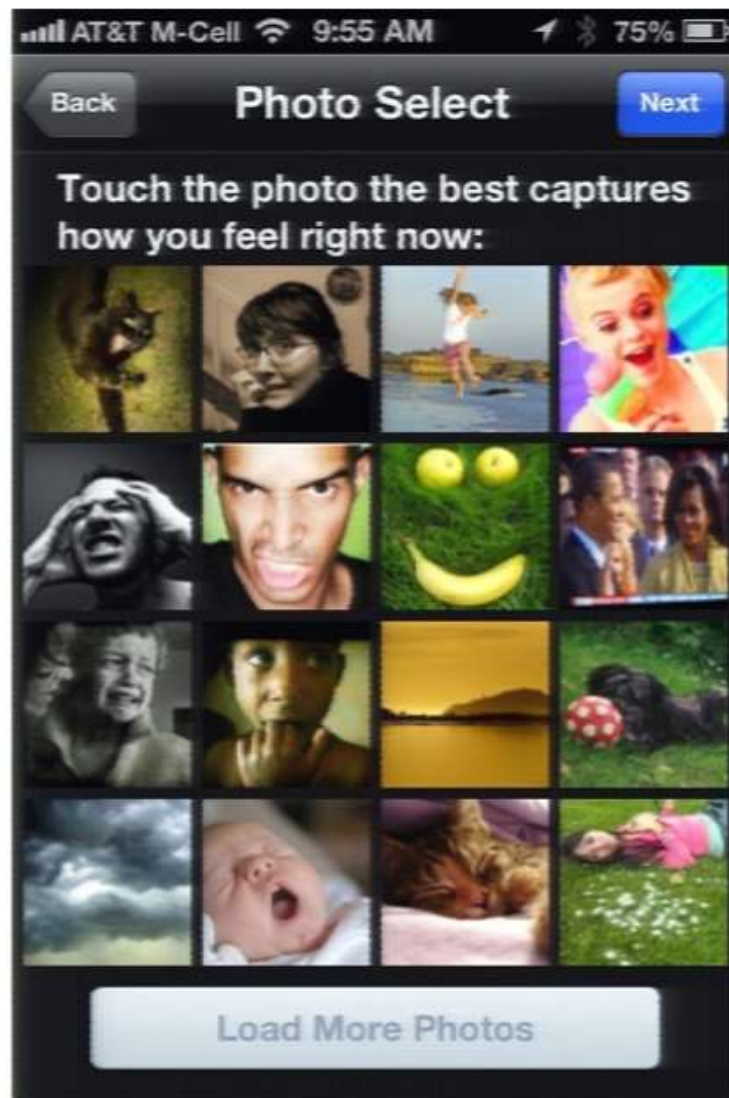
Ambient Wi-Fi Signals



# mobile apps data

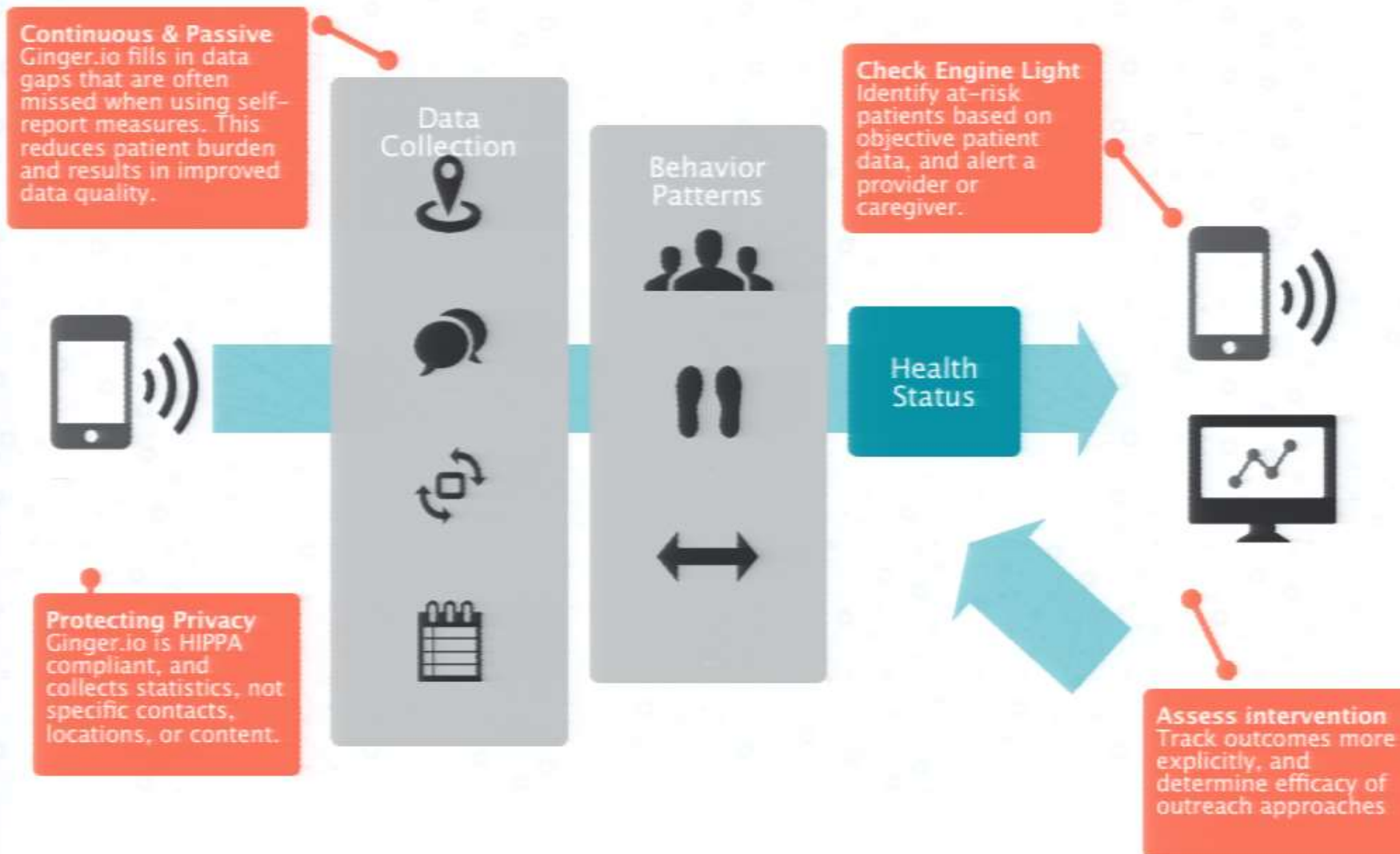


# Smart self report/EMA: Photographic Affect Meter, PAM (Pollak et al)

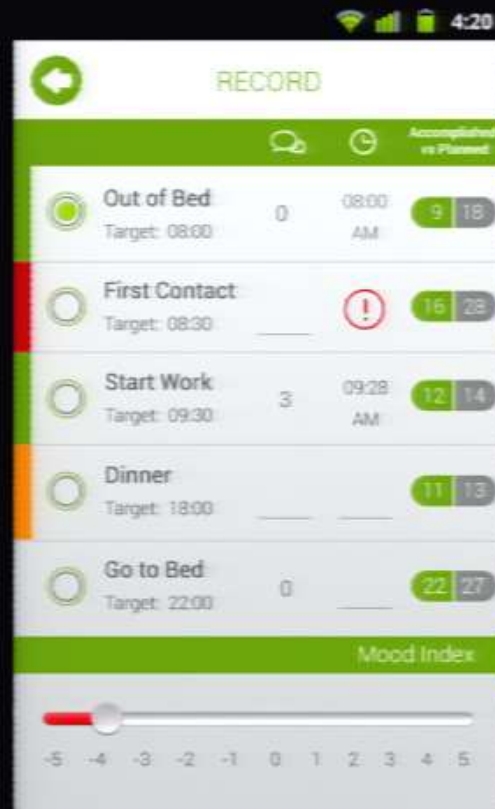




# Rich communication and activity data: The Ginger.io Platform



# Data driven tools for patients: MoodRhythm(TM) (Choudhury et al)



And yes...“Real Sensor” streams too



# Beyond mobile...

Leveraging digital traces from diverse consumer services

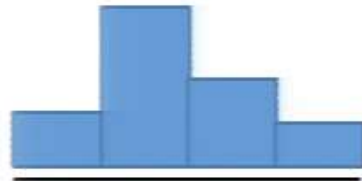
Your rows of their matrices...



## Example: consumer transaction patterns as small data



The goal is to transform purchasing patterns...



...into a consumption model (a categorical distribution of restaurant, fast food, drug store, etc. purchases)...



...on which descriptive statistics can be computed.



Sept Oct Nov

Since we have a progression of spending patterns over time...



Sept Oct Nov

...we can examine how the resulting distributions change...



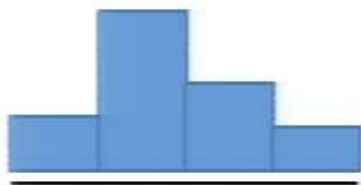
...and plot the statistics for each time frame, resulting in a time series which can then be correlated against other time series.



## Example: communication language patterns as small data



The goal is to transform text communication...



...into a “bag of words” model (a categorical distribution)...



...on which descriptive statistics can be computed.



Sept Oct Nov

Since we have a progression of emails over time...



Sept Oct Nov

...we can examine how the resulting distributions change...



...and plot the statistics for each time frame, resulting in a time series which can then be correlated against other time series.



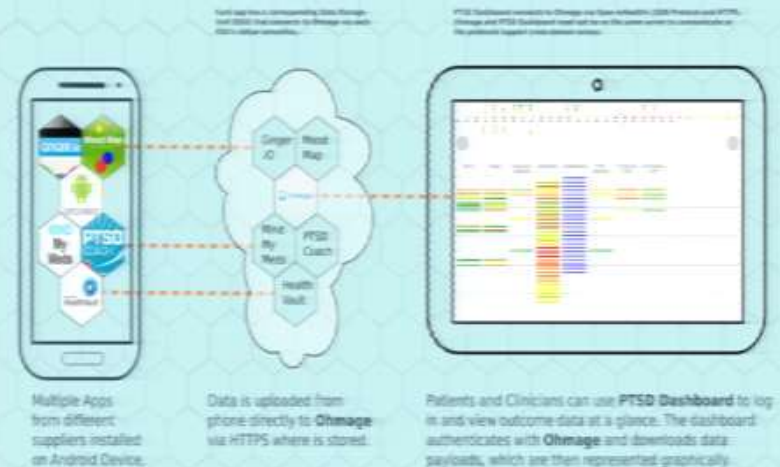
no one data stream tells the story

its about integration, fusion, and sense-making

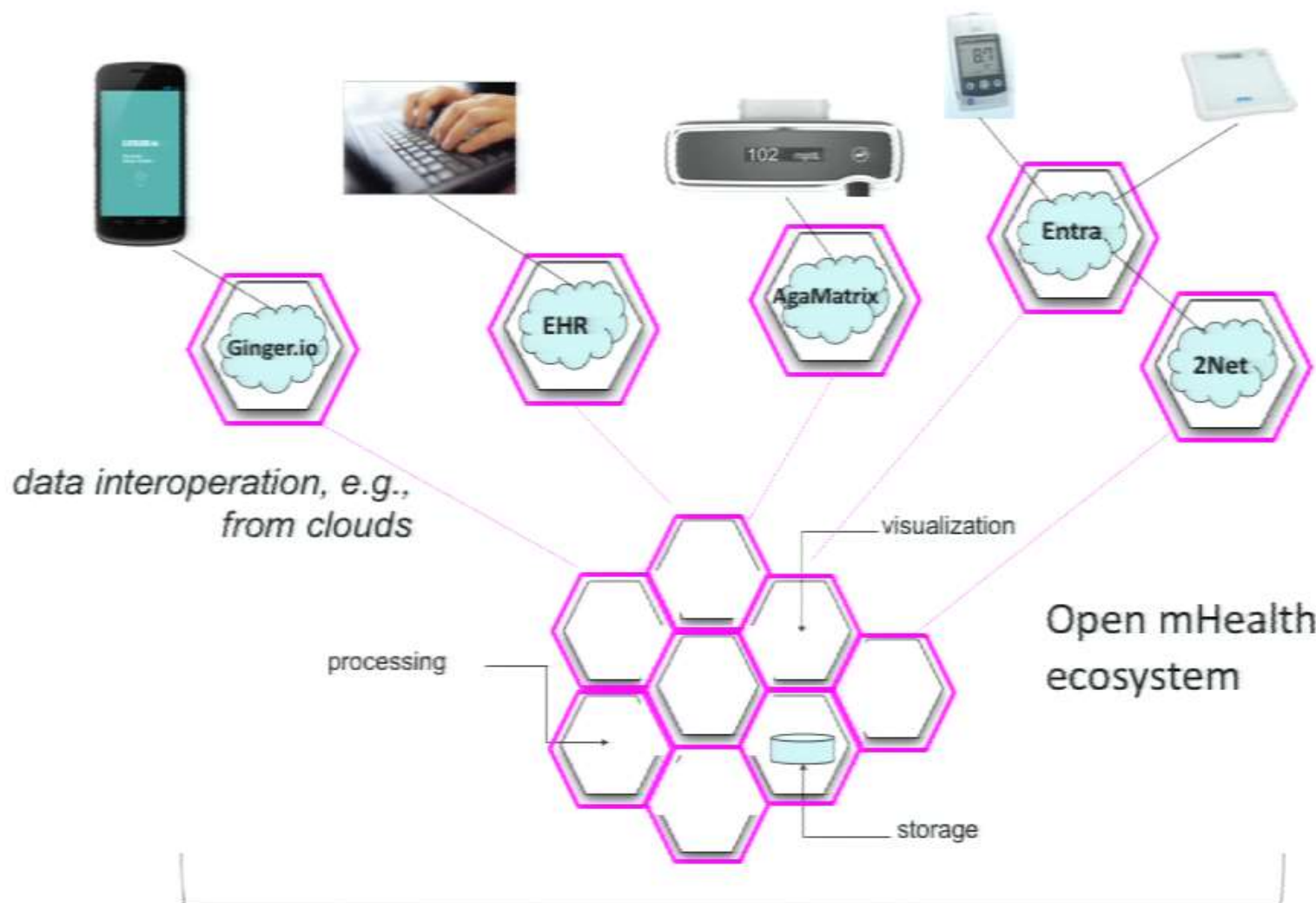
## Open mHeath Case Study Diabetes Scenario



## Open mHeath Case Study PTSD Scenario

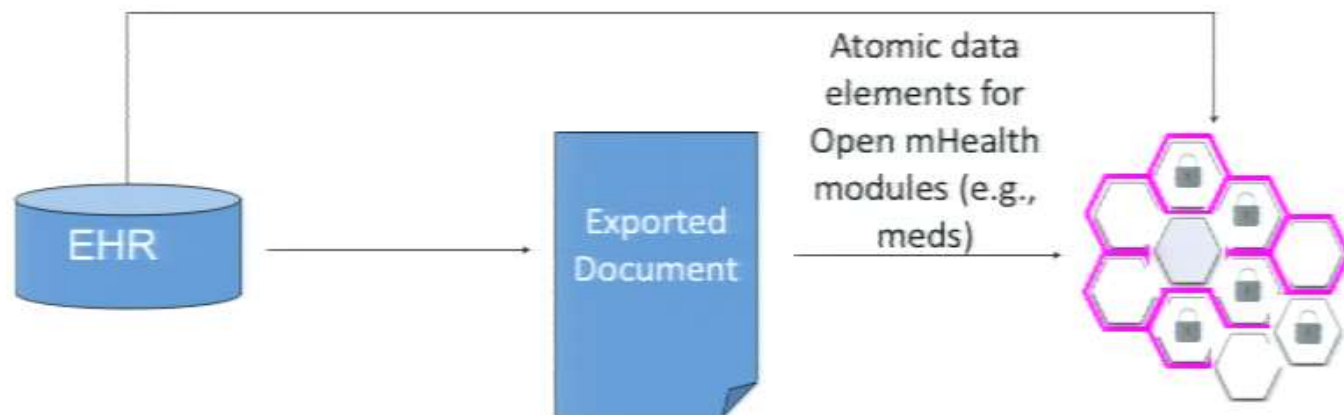


**promote modularity, data integration, and software reusability**



# Eventual integration w/ clinical data sources and workflows via EHR

*data from the EHR*



*data into the EHR*



## Key challenge 1: making “clinical sense” out of raw n=me data



???



Transform raw data streams into ***behavioral biomarkers***:  
specific behavioral traits to measure progress of disease and treatment

**state classification**

- sedentary/ambulatory
- at home/work
- app analytics (games, media...)
- communication



Transform raw data streams into **behavioral biomarkers**:  
specific behavioral traits to measure progress of disease and treatment

#### summarization

- ambulatory/sedentary cumulative and durations, walking speed
- sleep times, meal times
- time spent key locations, diameter of day
- social interaction



#### state classification

- sedentary/ambulatory
- at home/work
- app analytics (games, media...)
- communication





# Transform raw data streams into **behavioral biomarkers**: specific behavioral traits to measure progress of disease and treatment

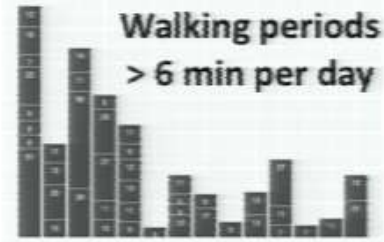
## behavioral biomarker

- individual's patterns; relevance is symptom and condition dependent
- 'function, fatigue, pain, depression, insomnia, cognition, self-medication...

## Hours at home per day



## Walking periods > 6 min per day



## summarization

- ambulatory/sedentary cumulative and durations, walking speed
- sleep times, meal times
- time spent key locations, diameter of day
- social interaction

Today

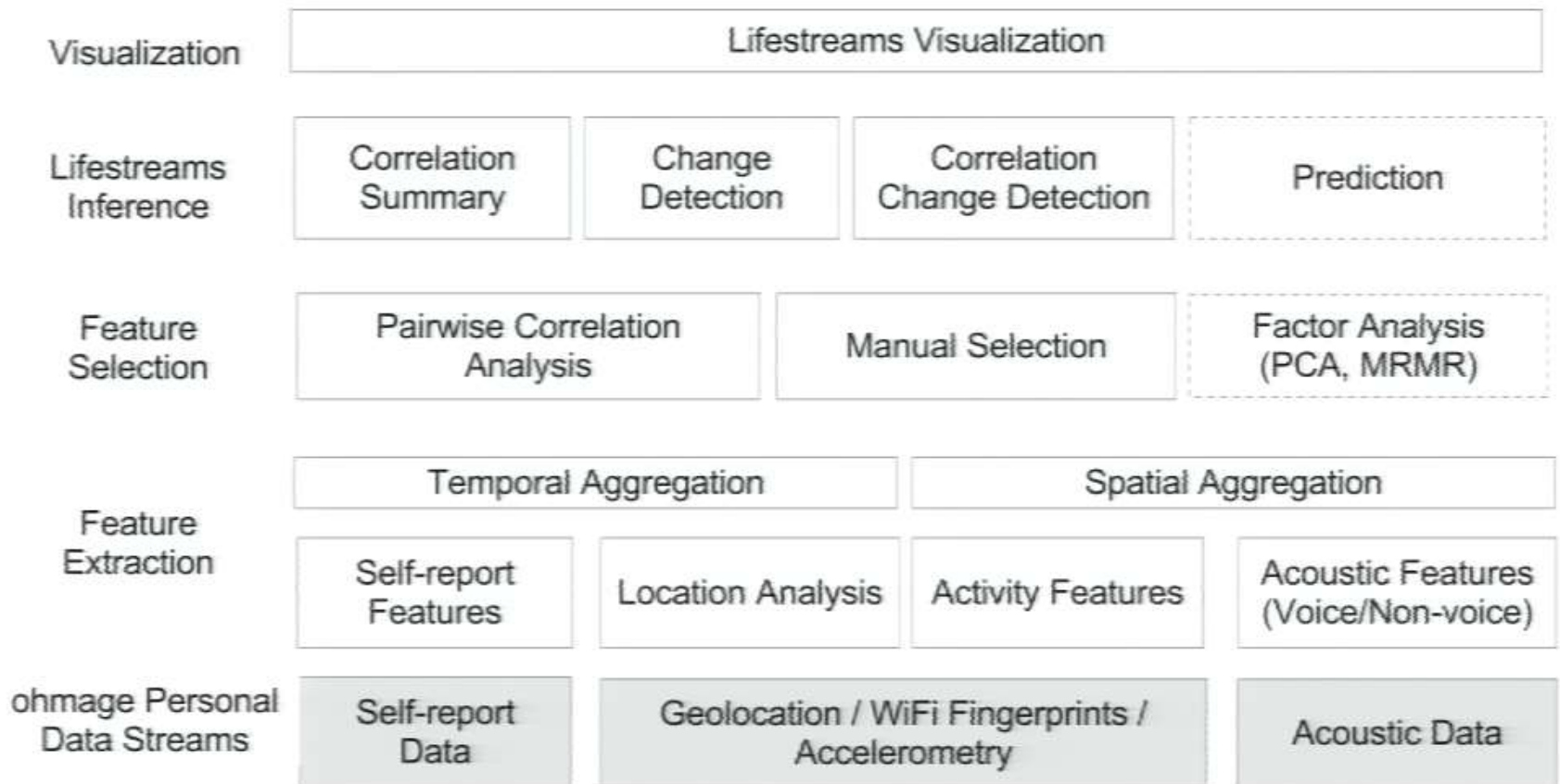


## state classification

- sedentary/ambulatory
- at home/work
- app analytics (games, media...)
- communication



# Modular tools to identify, iterate, share building blocks for behavioral biomarkers, sensemaking



(Lifestreams, Hsieh et al, Sensys 2013)

## Key challenge 2: small data governance

- Each data source has shared/other origins
- Individual has control over their corpus of data streams to correlate, fuse
- App/service utility derives from lack of anonymity
- Selective sharing embodied in apps--TMI works both ways in clinical domains



**Socio-technical architecture for small data:**  
individual as nexus of control



**Beyond 'health'**

**Life**

**quantified  
consumer**

**Health**

Independent living transitions

Auto-immune/inflammatory disease

Chronic pain: treatment and evidence

Mental health: depression, stress, ...

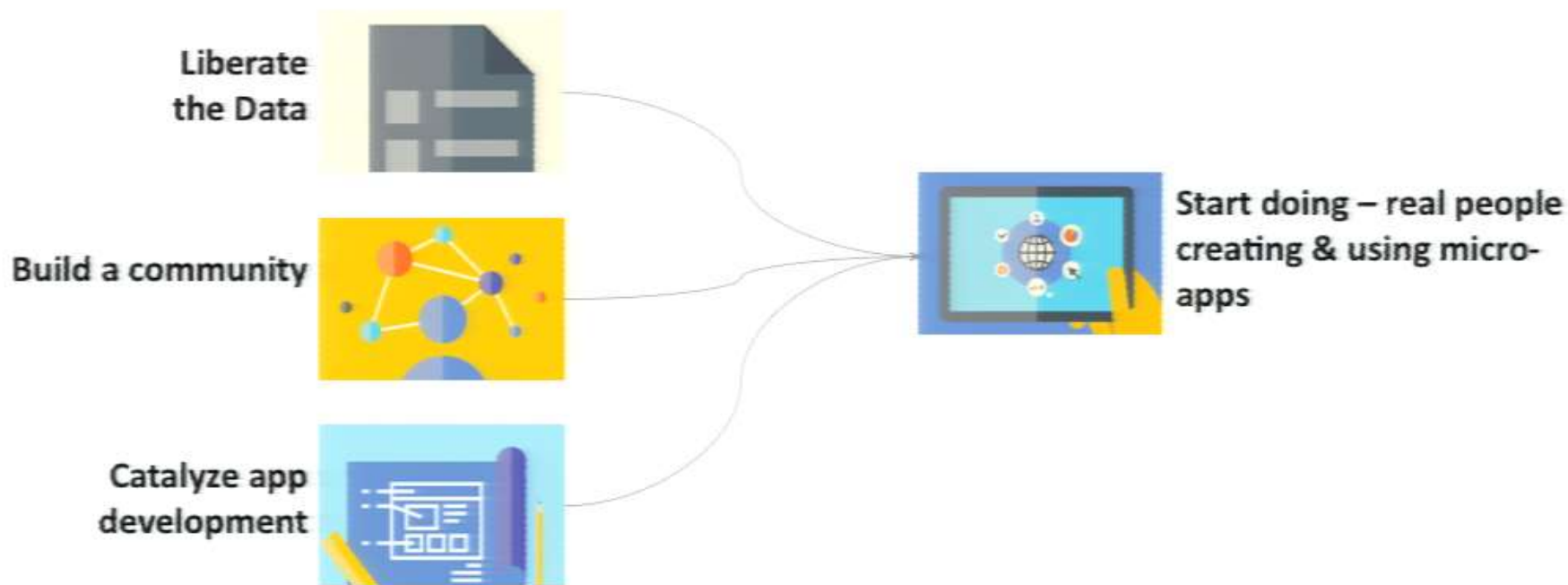
*\*Omics research*

**games,  
media,  
dating**

**family  
apps**

**quantified  
student**

## **mpire:** test-bed for small data and personal informatics



### **Partners**



**It's not rocket science....**

**It's pocket science**



## A Tech campus for the 21st century

### A campus that embraces and embeds external engagement

- Between technology user and technology creator: co-innovation
  - Connective media, Healthier life, Built environment
- Between the academic and the non-academic worlds
  - Large and small businesses, Startups, Government, Non-profits





# SCALABLE INFLUENCE ESTIMATION IN CONTINUOUS-TIME DIFFUSION NETWORKS

Nan Du <sup>1</sup>

Joint work with Le Song <sup>1</sup>, Manuel Gomez Rodriguez<sup>2</sup> and Hongyuan Zha<sup>1</sup>

<sup>1</sup>Georgia Institute of Technology

<sup>2</sup>Max Planck Institute for Intelligent Systems



- Diffusions of news, events, and virus, take place over social and information networks.



- Diffusions of news, events, and virus, take place over social and information networks.
- **Influence Estimation** : how to predict how many people will follow the fashion lead by the influential users ?





- Diffusions of news, events, and virus, take place over social and information networks.
- **Influence Estimation** : how to predict how many people will follow the fashion lead by the influential users ?
- **Influence Maximization** : how to identify such influential users to trigger the largest expected number of follow-ups ?



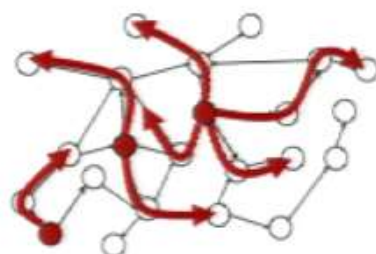
- Diffusions of news, events, and virus, take place over social and information networks.
- **Influence Estimation** : how to predict how many people will follow the fashion lead by the influential users ?
- **Influence Maximization** : how to identify such influential users to trigger the largest expected number of follow-ups ?
- **Time-Sensitive** : influence most users before time  $T$ .

- Diffusions of news, events, and virus, take place over social and information networks.
- **Influence Estimation** : how to predict how many people will follow the fashion lead by the influential users ?
- **Influence Maximization** : how to identify such influential users to trigger the largest expected number of follow-ups ?
- **Time-Sensitive** : influence most users before time  $T$ .
- **Scalability** : deal with large networks (millions of nodes) in practice.



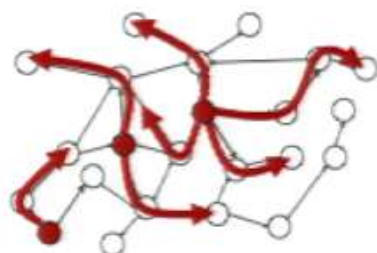
time-sensitive viral marketing

## 1 Continuous-time diffusion process.

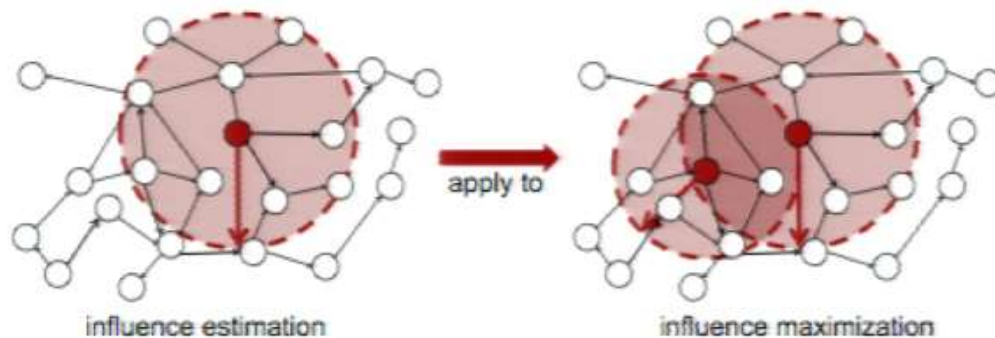




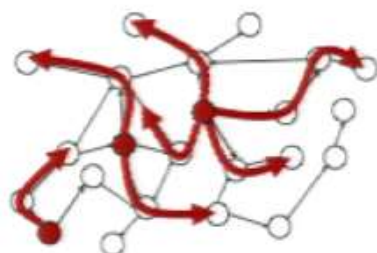
## 1 Continuous-time diffusion process.



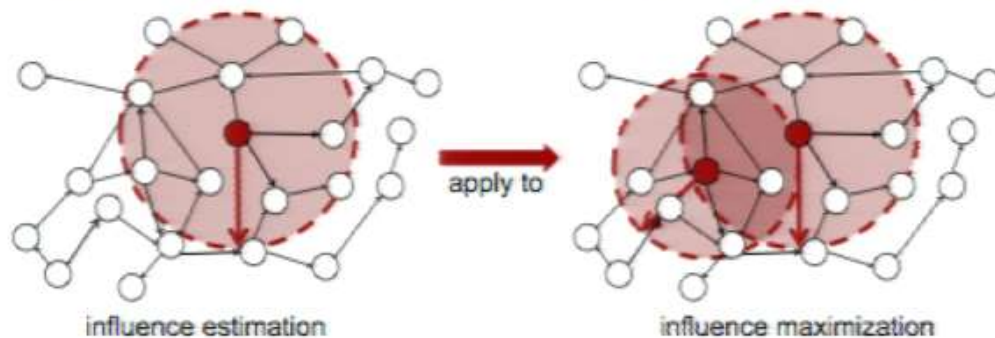
## 2 Efficient influence estimation and maximization.



## 1 Continuous-time diffusion process.



## 2 Efficient influence estimation and maximization.

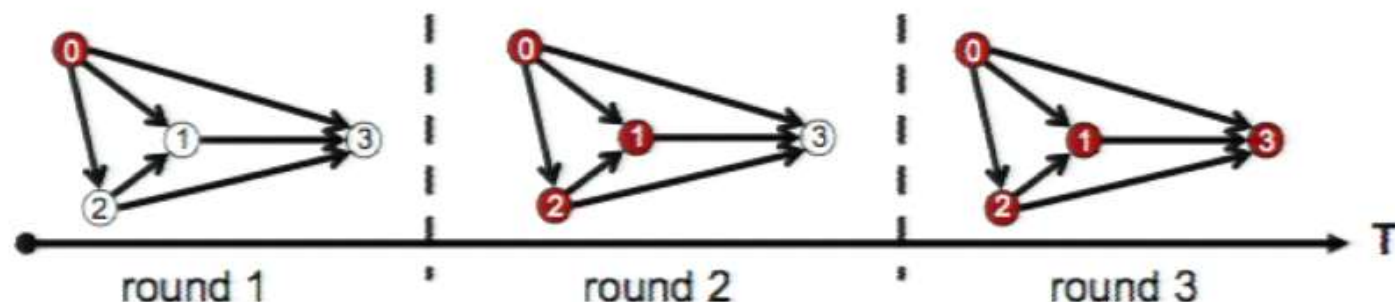


## 3 Experimental evaluation with synthetic and real diffusion data.

# CONTINUOUS VS. DISCRETE TIME DIFFUSION MODEL

- Traditionally, diffusion has been modeled as discrete steps (or rounds).

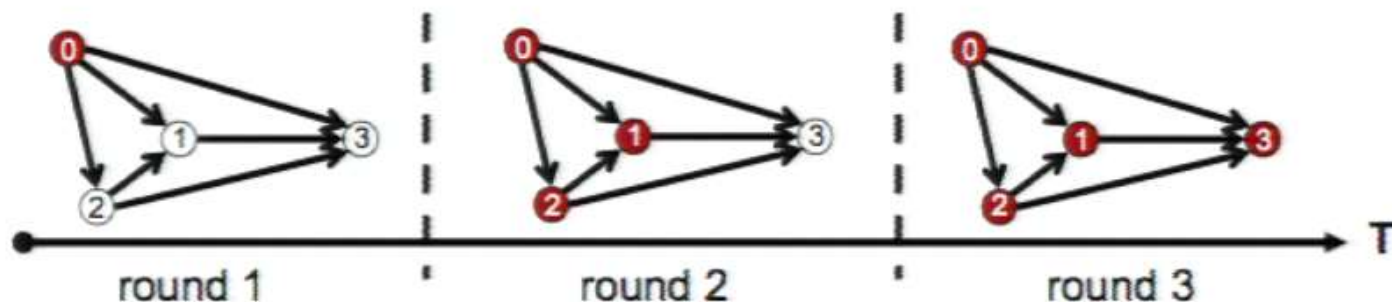
● infected    ○ uninfected



# CONTINUOUS VS. DISCRETE TIME DIFFUSION MODEL

- Traditionally, diffusion has been modeled as discrete steps (or rounds).

● infected ○ uninfected



- In reality, propagation does not go in rounds !

- how long is each round ?
- how many rounds do we need ?





# CONTINUOUS-TIME INDEPENDENT CASCADE MODEL

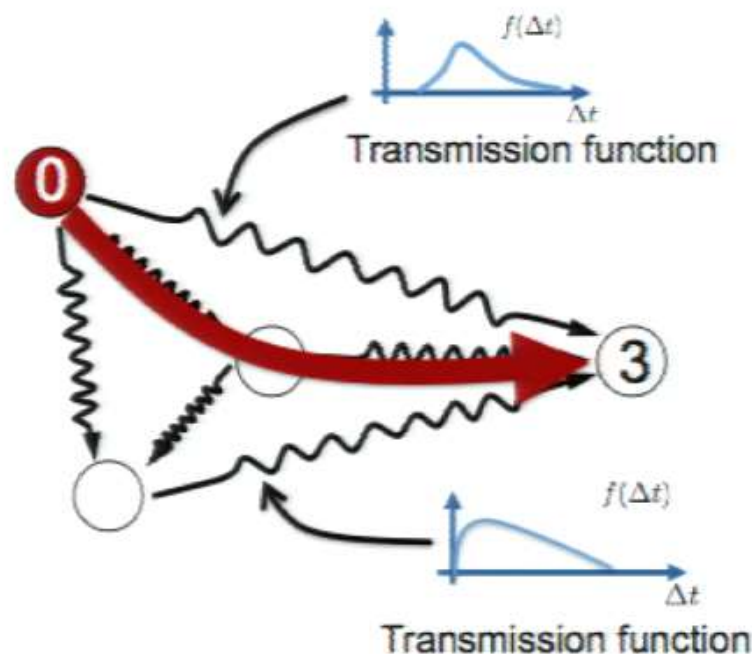
- Model mutually independent transmission time

$$\tau_{ji} = t_i - t_j.$$

- Pairwise conditional density ( transmission function )

$$f_{ji}(t_i|t_j) = f_{ji}(t_i - t_j).$$

- A network with stochastic edge weights.
- Infection time  
 $t_i$  = length of shortest path.



# ABSOLUTE INFECTION TIME VIEW

- The influence of sources  $\mathcal{A}$  by time  $T$  is

$$\sigma(\mathcal{A}, T) = \mathbb{E} \left[ \sum_{i \in \mathcal{V}} \mathbb{I}\{t_i \leq T\} \right] = \sum_{i \in \mathcal{V}} \Pr\{t_i \leq T\}$$

I

# ABSOLUTE INFECTION TIME VIEW

- The influence of sources  $\mathcal{A}$  by time  $T$  is

$$\sigma(\mathcal{A}, T) = \mathbb{E} \left[ \sum_{i \in \mathcal{V}} \mathbb{I}\{t_i \leq T\} \right] = \sum_{i \in \mathcal{V}} \Pr\{t_i \leq T\}$$

- Infection probability

$$\Pr\{t_i \leq T\} = \int_0^\infty \cdots \int_{t_i=0}^T \cdots \int_0^\infty \left( \prod_{j \in \mathcal{V}} p(t_j | \{t_l\}_{l \in \pi_j}) \right) \left( \prod_{j \in \mathcal{V}} dt_j \right)$$

# ABSOLUTE INFECTION TIME VIEW

- The influence of sources  $\mathcal{A}$  by time  $T$  is

$$\sigma(\mathcal{A}, T) = \mathbb{E} \left[ \sum_{i \in \mathcal{V}} \mathbb{I} \{t_i \leq T\} \right] = \sum_{i \in \mathcal{V}} \Pr \{t_i \leq T\}$$

- Infection probability

$$\Pr \{t_i \leq T\} = \int_0^\infty \cdots \int_{t_i=0}^T \cdots \int_0^\infty \left( \prod_{j \in \mathcal{V}} p(t_j | \{t_l\}_{l \in \pi_j}) \right) \left( \prod_{j \in \mathcal{V}} dt_j \right)$$

- Need to integrate all possible configurations of cascades where  $t_i < T$ .
- No closed form solution for general heterogeneous transmission function.
- Hard to approximate.

# NEIGHBORHOOD SIZE ESTIMATION

- Influence function

$$\sigma(\mathcal{A}, T) = \sum_{i \in \mathcal{V}} \Pr\{t_i \leq T\}$$

- No need to calculate  $\Pr\{t_i \leq T\}$  individually.



# NEIGHBORHOOD SIZE ESTIMATION

- Influence function

$$\sigma(\mathcal{A}, T) = \sum_{i \in \mathcal{V}} \Pr\{t_i \leq T\}$$

- No need to calculate  $\Pr\{t_i \leq T\}$  individually.

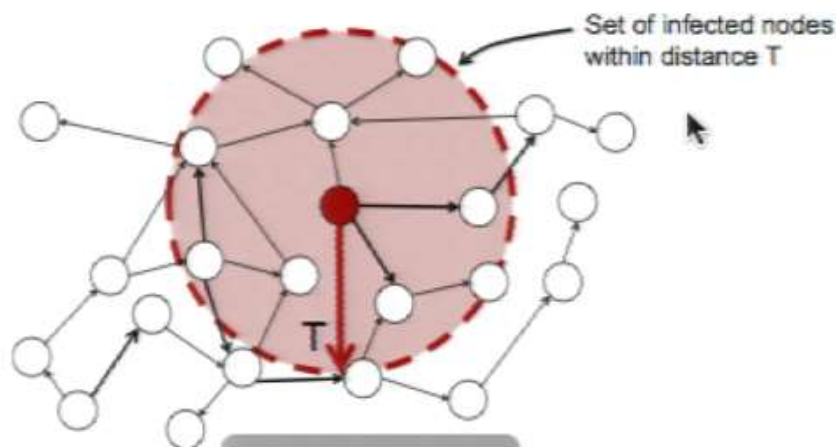
# NEIGHBORHOOD SIZE ESTIMATION

- Influence function

$$\sigma(\mathcal{A}, T) = \sum_{i \in \mathcal{V}} \Pr \{t_i \leq T\}$$

- No need to calculate  $\Pr \{t_i \leq T\}$  individually.
- Given a set of  $\{\tau_{ji}^l\}_{(j,i) \in \mathcal{E}}$ , only care about

$$\sum_{i \in \mathcal{V}} \mathbb{I} \{t_i \leq T\} = |\mathcal{N}(\{j\}, T)| = |\{i : t_i \leq T\}|$$



## NEIGHBORHOOD SIZE ESTIMATION

- Sample  $n$  sets of  $G_l := \{\tau_{ji}^l\}_{(j,i) \in \mathcal{E}} \sim \prod_{(j,i) \in \mathcal{E}} f_{ji}(\tau_{ji})$

# NEIGHBORHOOD SIZE ESTIMATION

- Sample  $n$  sets of  $G_l := \{\tau_{ji}^l\}_{(j,i) \in \mathcal{E}} \sim \prod_{(j,i) \in \mathcal{E}} f_{ji}(\tau_{ji})$
- Average the counts across  $n$  samples.

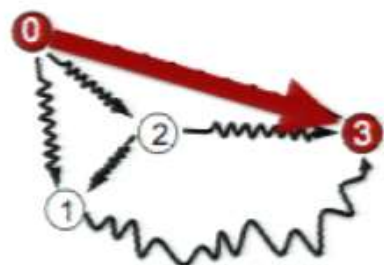
$$\begin{aligned}\sigma(\mathcal{A}, T) &= \mathbb{E} \left[ \sum_{i \in \mathcal{V}} \mathbb{I}\{t_i \leq T\} \right] \\ &\approx \frac{1}{n} \left( \sum_{i \in \mathcal{V}} \mathbb{I}\{t_i \leq T | G_1\} + \dots + \sum_{i \in \mathcal{V}} \mathbb{I}\{t_i \leq T | G_n\} \right)\end{aligned}$$

# NEIGHBORHOOD SIZE ESTIMATION

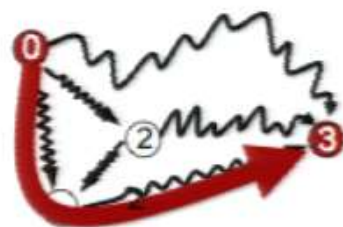
- Sample  $n$  sets of  $G_l := \{\tau_{ji}^l\}_{(j,i) \in \mathcal{E}} \sim \prod_{(j,i) \in \mathcal{E}} f_{ji}(\tau_{ji})$
- Average the counts across  $n$  samples.

$$\sigma(\mathcal{A}, T) = \mathbb{E} \left[ \sum_{i \in \mathcal{V}} \mathbb{I}\{t_i \leq T\} \right]$$
$$\approx \frac{1}{n} \left( \sum_{i \in \mathcal{V}} \mathbb{I}\{t_i \leq T | G_1\} + \dots + \sum_{i \in \mathcal{V}} \mathbb{I}\{t_i \leq T | G_n\} \right)$$

- To calculate  $\mathbb{I}\{t_i \leq T | G_l\}$ , check whether **length of shortest path**  $\leq T$  on each sampled network.



✓  $\mathbb{I}(t_3 \leq T)$



Page 20 of 49



✗  $\mathbb{I}(t_3 \leq T)$

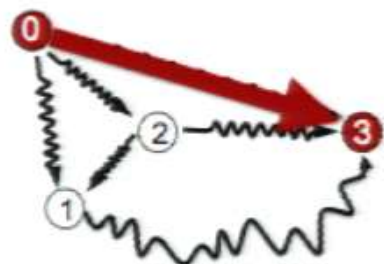


# NEIGHBORHOOD SIZE ESTIMATION

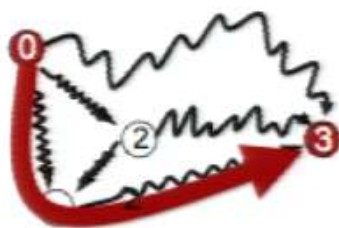
- Sample  $n$  sets of  $G_l := \{\tau_{ji}^l\}_{(j,i) \in \mathcal{E}} \sim \prod_{(j,i) \in \mathcal{E}} f_{ji}(\tau_{ji})$
- Average the counts across  $n$  samples.

$$\sigma(\mathcal{A}, T) = \mathbb{E} \left[ \sum_{i \in \mathcal{V}} \mathbb{I}\{t_i \leq T\} \right]$$
$$\approx \frac{1}{n} \left( \sum_{i \in \mathcal{V}} \mathbb{I}\{t_i \leq T | G_1\} + \dots + \sum_{i \in \mathcal{V}} \mathbb{I}\{t_i \leq T | G_n\} \right)$$

- To calculate  $\mathbb{I}\{t_i \leq T | G_l\}$ , check whether **length of shortest path**  $\leq T$  on each sampled network.



✓  $\mathbb{I}(t_3 \leq T)$



✓  $\mathbb{I}(t_3 \leq T)$



✗  $\mathbb{I}(t_3 \leq T)$

# NEIGHBORHOOD SIZE ESTIMATION

- Using shortest path is not scalable.

I

# NEIGHBORHOOD SIZE ESTIMATION

- Using shortest path is not scalable.
- Influence Estimation of a single source  $j$ 
  - $\sigma(\{j\}, T)$
  - Compute all shortest paths from  $j$  to the other nodes.

I

# NEIGHBORHOOD SIZE ESTIMATION

- Using shortest path is not scalable.
- Influence Estimation of a single source  $j$ 
  - $\sigma(\{j\}, T)$
  - Compute all shortest paths from  $j$  to the other nodes.
- Which source is the best ?
  - Chose  $j$  with the largest  $\sigma(\{j\}, T)$
  - Try source  $j = 0, \dots, |\mathcal{V}| - 1, O(|\mathcal{V}|^2)$
- Quadratic in network size  
Can not deal with large networks !

I

## NEIGHBORHOOD SIZE ESTIMATION

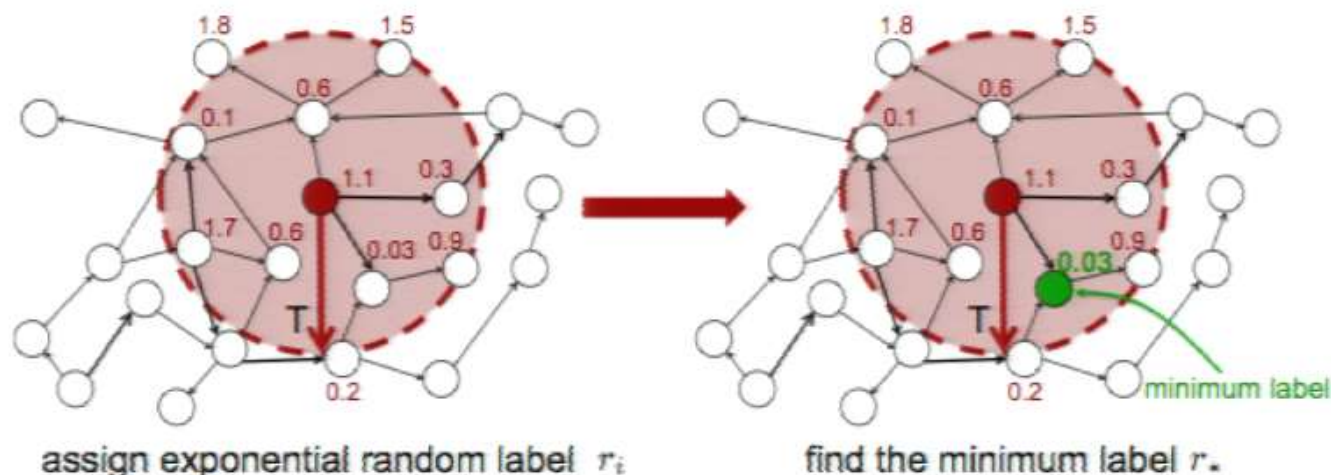
- Directly estimate the neighborhood size by Cohen's algorithm !

I



# NEIGHBORHOOD SIZE ESTIMATION

- Directly estimate the neighborhood size by Cohen's algorithm !

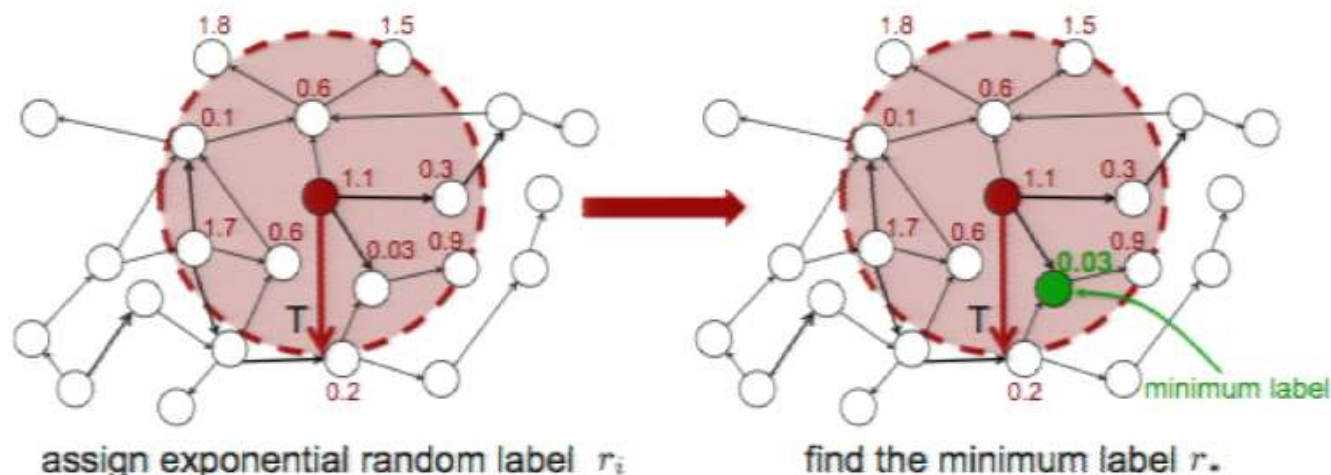


- Draw  $m$  sets of i.i.d random labels  $\{r_i^u\}_{u=1}^m \sim e^{-r_i}$ .

$\mathbb{I}$

# NEIGHBORHOOD SIZE ESTIMATION

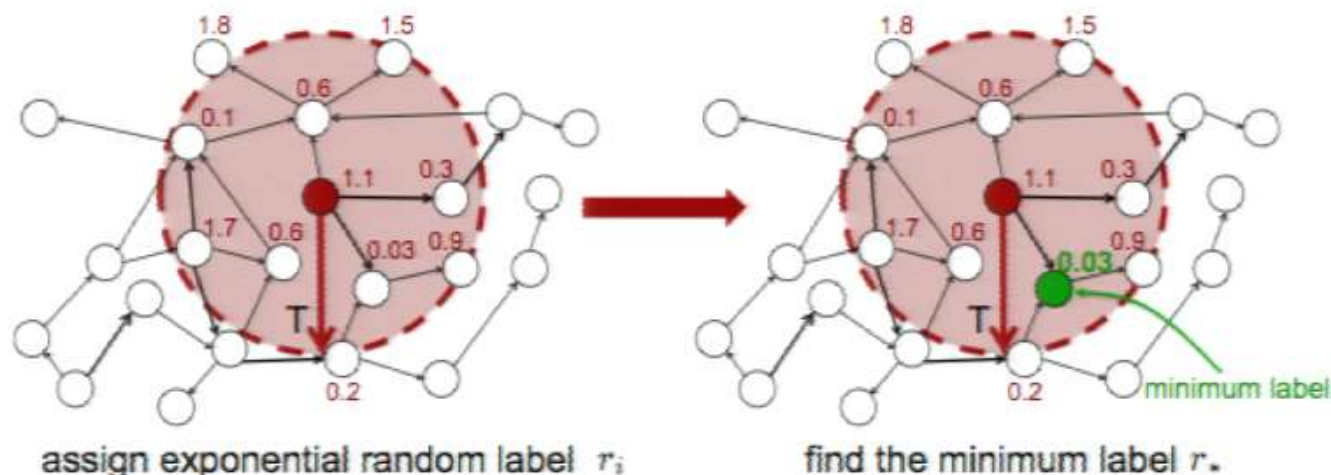
- Directly estimate the neighborhood size by Cohen's algorithm !



- Draw  $m$  sets of i.i.d random labels  $\{r_i^u\}_{u=1}^m \sim e^{-r_i}$ .
- Find the minimum label  $\{r_*^u\}_{u=1}^m$  within distance  $T$  by Cohen's algorithm in  $\tilde{O}(|\mathcal{E}|)$ .

# NEIGHBORHOOD SIZE ESTIMATION

- Directly estimate the neighborhood size by Cohen's algorithm !



- Draw  $m$  sets of i.i.d random labels  $\{r_i^u\}_{u=1}^m \sim e^{-r_i}$ .
- Find the minimum label  $\{r_*^u\}_{u=1}^m$  within distance  $T$  by Cohen's algorithm in  $\tilde{O}(|\mathcal{E}|)$ .
- Estimate  $|\mathcal{N}(\{j\}, T)| \approx \frac{m-1}{\sum_{u=1}^m r_*^u}$  using the property of exponential distribution.

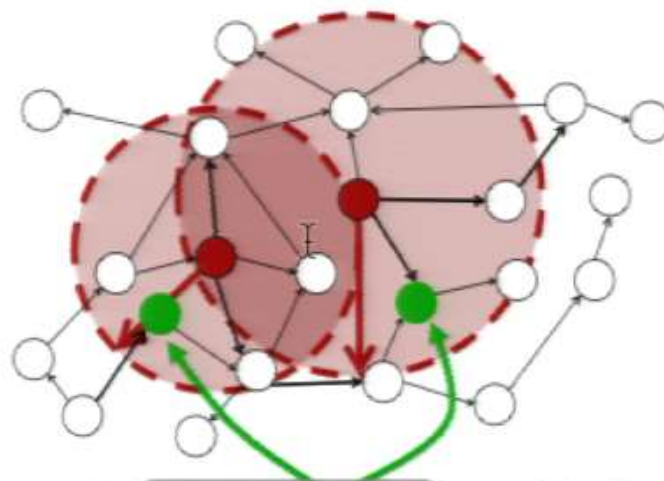
# MULTIPLE SOURCES

- Multiple sources  $\mathcal{A}$

$$\mathcal{N}(\mathcal{A}, T) = \bigcup_{s \in \mathcal{A}} \mathcal{N}(s, T).$$

- The overall least label

$$r_* = \min_{i \in \mathcal{A}} \min_{j \in \mathcal{N}(i, T)} r_j$$



mini Page 28 of 49 least labels

# OVERALL ALGORITHM CONTINEST

1. Sample  $n$  sets of random transmission times



$$\{\tau_{ji}^l\}_{(j,i) \in \mathcal{E}} \sim \prod_{(j,i) \in \mathcal{E}} f_{ji}(\tau_{ji})$$

I



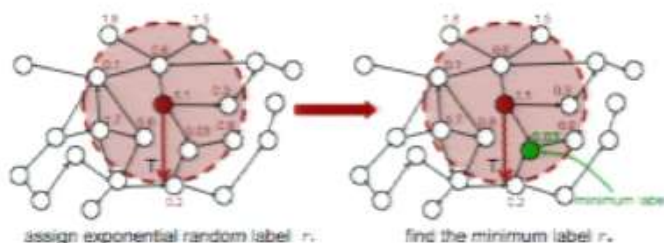
# OVERALL ALGORITHM CONTINEST

1. Sample  $n$  sets of random transmission times



$$\{\tau_{ji}^l\}_{(j,i) \in \mathcal{E}} \sim \prod_{(j,i) \in \mathcal{E}} f_{ji}(\tau_{ji})$$

2. Given a set of  $\{\tau_{ji}^l\}_{(j,i) \in \mathcal{E}}$ , sample  $m$  sets of random labels

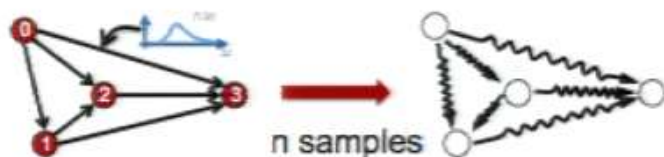


$$\{r_i^u\}_{i \in \mathcal{V}} \sim \prod_{i \in \mathcal{V}} \exp(-r_i)$$

I

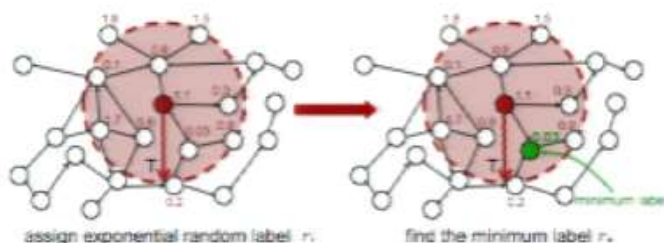
# OVERALL ALGORITHM CONTINEST

1. Sample  $n$  sets of random transmission times



$$\{\tau_{ji}^l\}_{(j,i) \in \mathcal{E}} \sim \prod_{(j,i) \in \mathcal{E}} f_{ji}(\tau_{ji})$$

2. Given a set of  $\{\tau_{ji}^l\}_{(j,i) \in \mathcal{E}}$ , sample  $m$  sets of random labels



$$\{r_i^u\}_{i \in \mathcal{V}} \sim \prod_{i \in \mathcal{V}} \exp(-r_i)$$

3. Find the minimum label  $\{r_*^u\}_{u=1}^m$  within  $T$  using Cohen's algorithm.

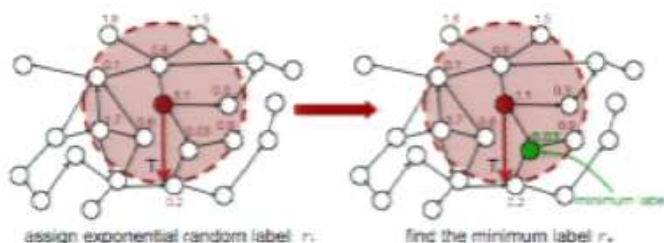
# OVERALL ALGORITHM CONTINEST

1. Sample  $n$  sets of random transmission times



$$\{\tau_{ji}^l\}_{(j,i) \in \mathcal{E}} \sim \prod_{(j,i) \in \mathcal{E}} f_{ji}(\tau_{ji})$$

2. Given a set of  $\{\tau_{ji}^l\}_{(j,i) \in \mathcal{E}}$ , sample  $m$  sets of random labels



$$\{r_i^u\}_{i \in \mathcal{V}} \sim \prod_{i \in \mathcal{V}} \exp(-r_i)$$

3. Find the minimum label  $\{r_*^u\}_{u=1}^m$  within  $T$  using Cohen's algorithm.
4. Estimate  $\sigma(\mathcal{A}, T)$  by sample averages

$$\sigma(\mathcal{A}, T) \approx \frac{1}{n} \sum_{l=1}^n \left( (m-1) / \sum_{u_l=1}^m r_*^{u_l} \right)$$

# OVERALL ALGORITHM CONTINEST

## THEOREM

*Draw the following number of samples for the set of random transmission times*

$$n \geq \frac{C\Lambda(T, 1/m)}{\epsilon^2} \log \left( \frac{2|\mathcal{V}|}{\delta} \right),$$

*and for each set of random transmission times, draw  $m$  set of random labels.  
Then  $|\hat{\sigma}(\mathcal{A}, T) - \sigma(\mathcal{A}, T)| \leq \epsilon$  uniformly for all  $\mathcal{A}$  with  $|\mathcal{A}| \leq C$ , with probability at least  $1 - \delta$ .*

I



# OVERALL ALGORITHM CONTINEST

## THEOREM

*Draw the following number of samples for the set of random transmission times*

$$n \geq \frac{C\Lambda(T, 1/m)}{\epsilon^2} \log \left( \frac{2|\mathcal{V}|}{\delta} \right),$$

*and for each set of random transmission times, draw  $m$  set of random labels.  
Then  $|\hat{\sigma}(\mathcal{A}, T) - \sigma(\mathcal{A}, T)| \leq \epsilon$  uniformly for all  $\mathcal{A}$  with  $|\mathcal{A}| \leq C$ , with probability at least  $1 - \delta$ .*

- Implications : influence at the longer time window  $T$  requires more samples.
- In practice : large  $n = 10K$  allows small  $m = 5$  to achieve good performance.



# INFLUENCE MAXIMIZATION

- We seek to solve

$$\mathcal{A}^* = \operatorname{argmax}_{|\mathcal{A}| \leq c} \sigma(\mathcal{A}, T)$$

which is NP-hard in general.

I

# INFLUENCE MAXIMIZATION

- We seek to solve

$$\mathcal{A}^* = \operatorname{argmax}_{|\mathcal{A}| \leq C} \sigma(\mathcal{A}, T)$$

which is NP-hard in general.

- $\sigma(\mathcal{A}, T)$  is a non-negative, monotonic, submodular function.

I

# INFLUENCE MAXIMIZATION

- We seek to solve

$$\mathcal{A}^* = \operatorname{argmax}_{|\mathcal{A}| \leq C} \sigma(\mathcal{A}, T)$$

which is NP-hard in general.

- $\sigma(\mathcal{A}, T)$  is a non-negative, monotonic, submodular function.
- Greedy algorithm achieves at least a fraction  $(1 - 1/e)$  of the optimal value (OPT)

I

# INFLUENCE MAXIMIZATION

- We seek to solve

$$\mathcal{A}^* = \operatorname{argmax}_{|\mathcal{A}| \leq C} \sigma(\mathcal{A}, T)$$

which is NP-hard in general.

- $\sigma(\mathcal{A}, T)$  is a non-negative, monotonic, submodular function.
- Greedy algorithm achieves at least a fraction  $(1 - 1/e)$  of the optimal value (OPT)

## THEOREM

*Suppose the influence  $\sigma(\mathcal{A}, T)$  for all  $\mathcal{A}$  with  $|\mathcal{A}| \leq C$  are estimated uniformly with error  $\epsilon$  and confidence  $1 - \delta$ , the greedy algorithm returns a set of sources  $\hat{\mathcal{A}}$  such that*

$$\sigma(\hat{\mathcal{A}}, T) \geq (1 - 1/e)OPT - 2C\epsilon$$

*with probability at least  $1 - \delta$ .*

Page 39 of 49

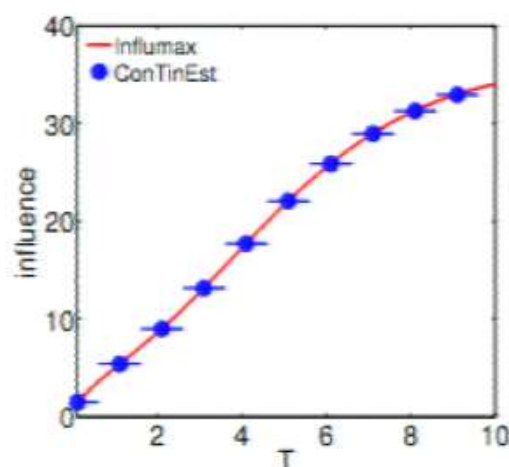
# EXPERIMENTAL EVALUATION

- Synthetic dataset
  - Generate network structure.
  - Weibull pairwise transmission function.
- Real dataset
  - MemeTracker data (172m news articles 08/2009-09/2009).
- Evaluation
  - Accuracy of estimated influence.
  - Quality of selected sources.
  - Scalability.

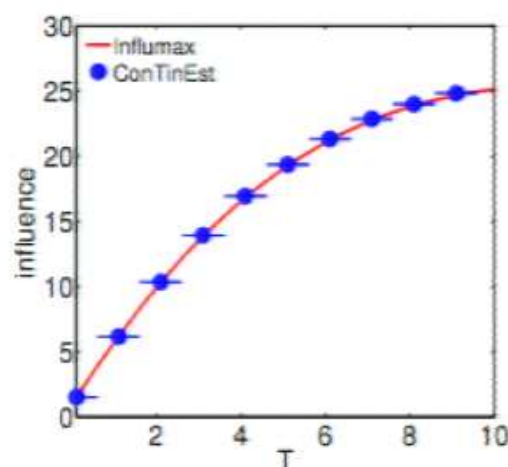


# SYNTHETIC DATASET

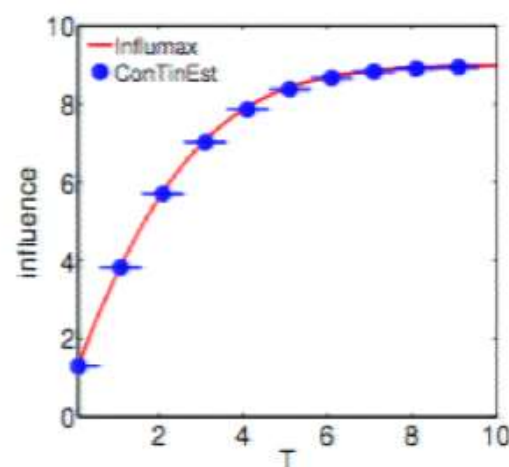
Accuracy of the estimated influence (highest out-degree node)



(a) Core-periphery



(b) Random

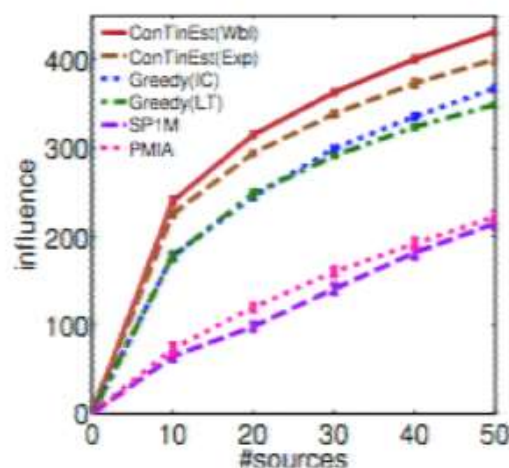


(c) Hierarchical

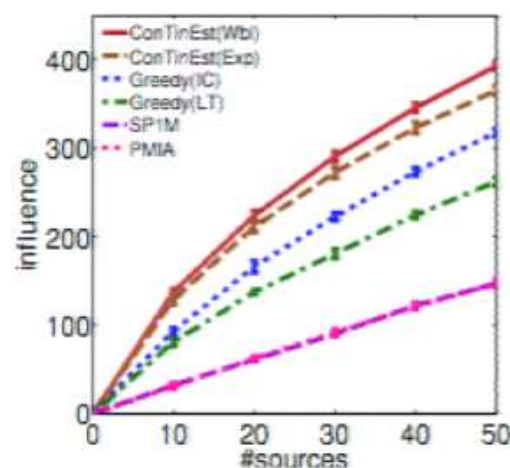
- CONTINEST is close to INFLUMAX (sparse small networks, exponential transmission functions).
- accuracy does not depend on network structure (128 nodes, 141 edges).

# SYNTHETIC DATASET

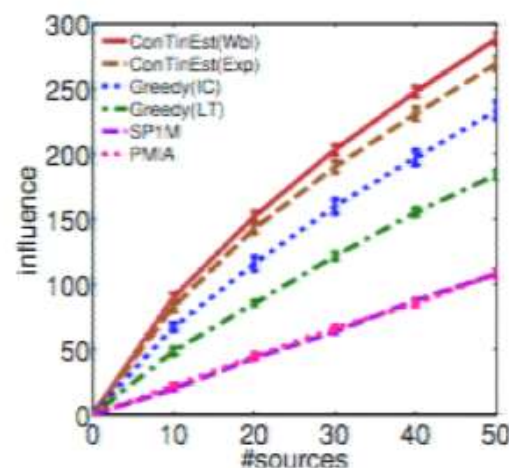
## Quality of the selected nodes for influence maximization



(a) Core-periphery



(b) Random

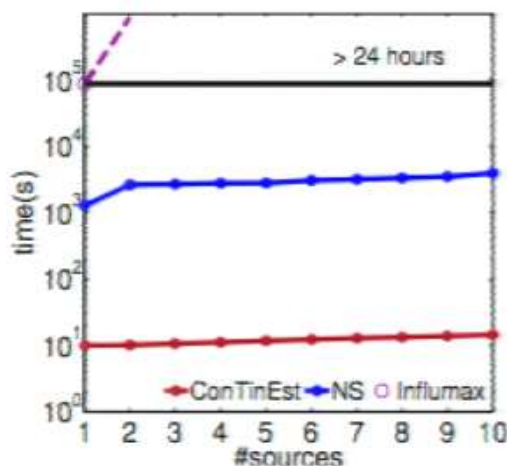


(c) Hierarchical

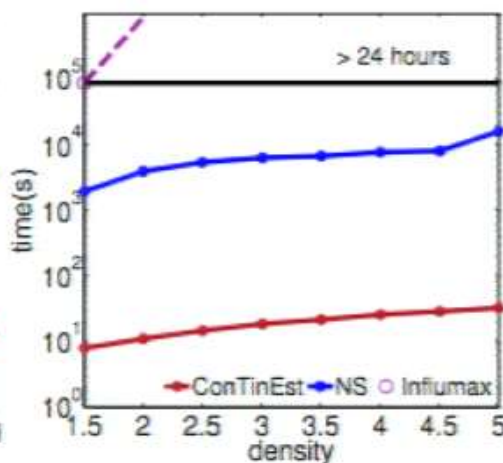
- CONTINEST typically outperforms competitive methods by 20%.
- Performance does not depend on network structure (1024 nodes, 2048 edges).

# SYNTHETIC DATASET

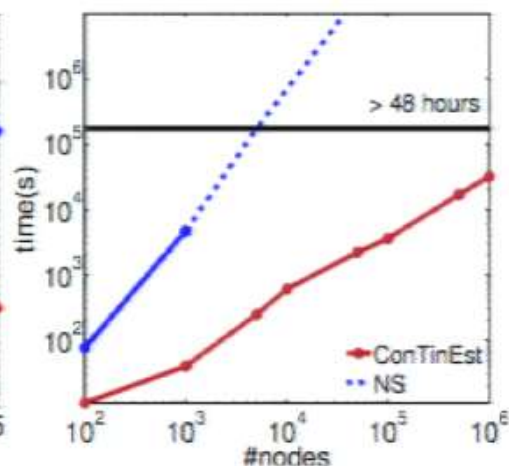
## Scalability of influence maximization



(a) # sources  
Small network



(b) network density  
Small network



(c) network size  
Up to one million nodes

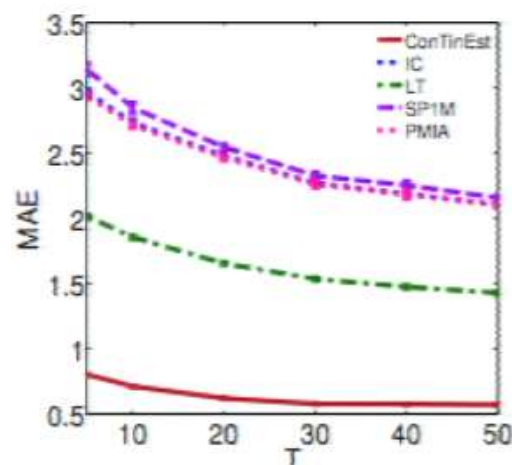
- Small network : 128 nodes.
- Large network : up to 1 million nodes, with density 1.5.
- Our algorithm : sample 10K networks, 5 random labels.

## REAL DATASET

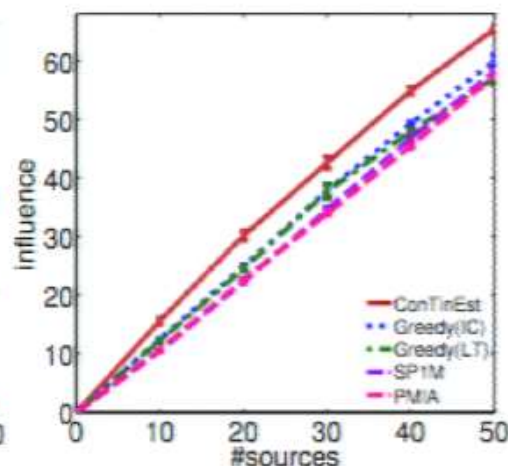
- 10,967 cascades.
- Use 80% cascades for learning continuous-time diffusion model.
- Select sources based on the learnt model.
- Evaluate influence of the sources using 20% test cascades.
- Compared to discrete-time diffusion models and scalable heuristics.



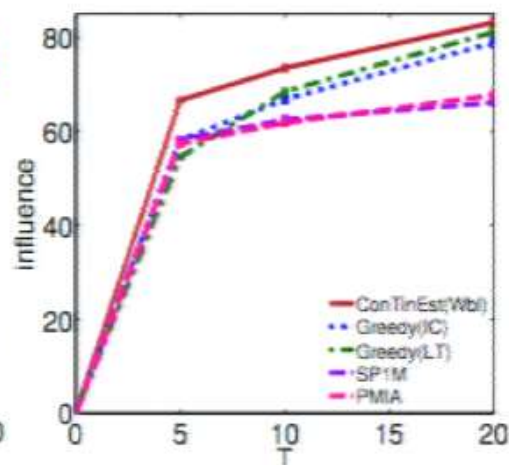
# REAL DATASET



(a) Estimation error



(b) #sources



(c) observation window

- CONTINEST achieves the lowest MAE error.



# CONCLUSION

- A randomized algorithm achieving :
  - the lowest estimation error in real data.
  - the largest influence within short time period.
  - the scaling up to millions of nodes in practice.

# CONCLUSION

- A randomized algorithm achieving :
  - the lowest estimation error in real data.
  - the largest influence within short time period.
  - the scaling up to millions of nodes in practice.
- Future work :
  - User engagement maximization of online systems.
  - Influence minimization and manipulation.
  - More general continuous-time diffusion model.

# CONCLUSION

- A randomized algorithm achieving :
  - the lowest estimation error in real data.
  - the largest influence within short time period.
  - the scaling up to millions of nodes in practice.
- Future work :
  - User engagement maximization of online systems.
  - Influence minimization and manipulation.
  - More general continuous-time diffusion model.
- Welcome to our poster F41 for detailed discussion.

