



Understanding and Managing Cascades on Large Graphs

B. Aditya Prakash

Virginia Tech.

Christos Faloutsos

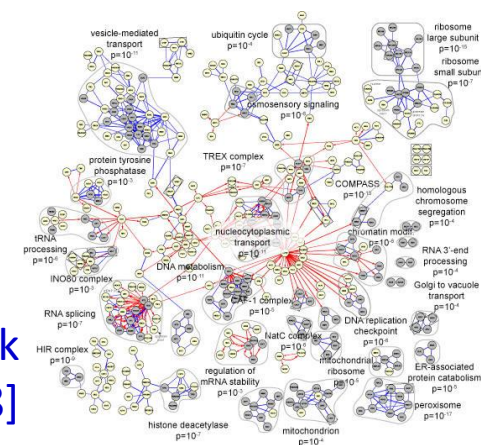
Carnegie Mellon University



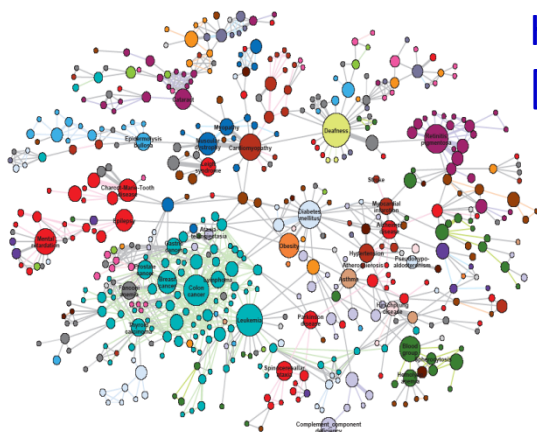
Networks are everywhere!



Facebook Network [2010]

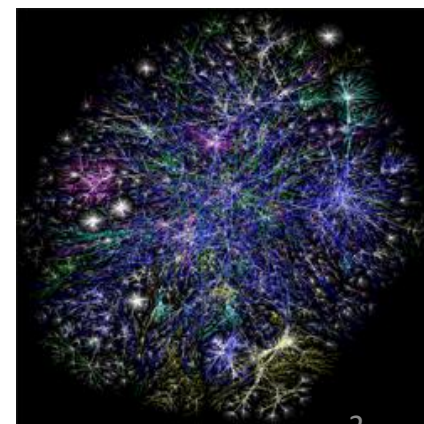


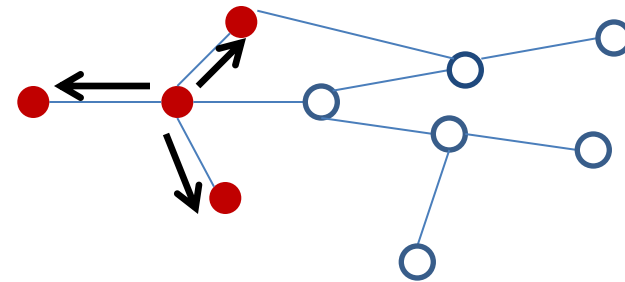
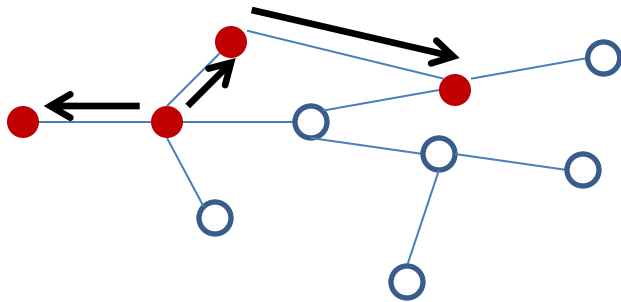
Gene Regulatory Network [Decourty 2008]



Human Disease Network [Barabasi 2007]

The Internet [2005]





Dynamical Processes *over* networks
are also everywhere!

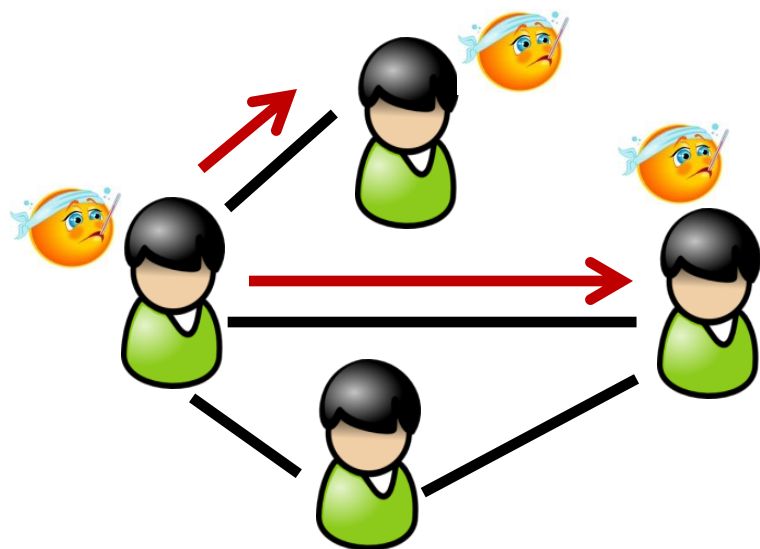
Why do we care?

- Social collaboration
- Information Diffusion
- Viral Marketing
- Epidemiology and Public Health
- Cyber Security
- Human mobility
- Games and Virtual Worlds
- Ecology



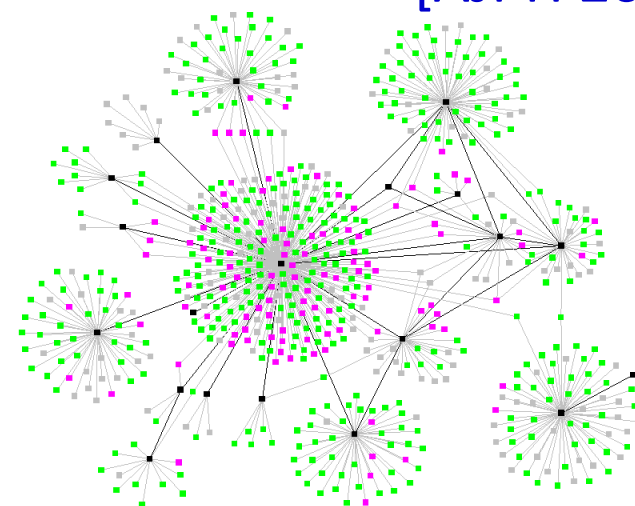
Why do we care? (1: Epidemiology)

- Dynamical Processes over networks



Diseases over contact networks

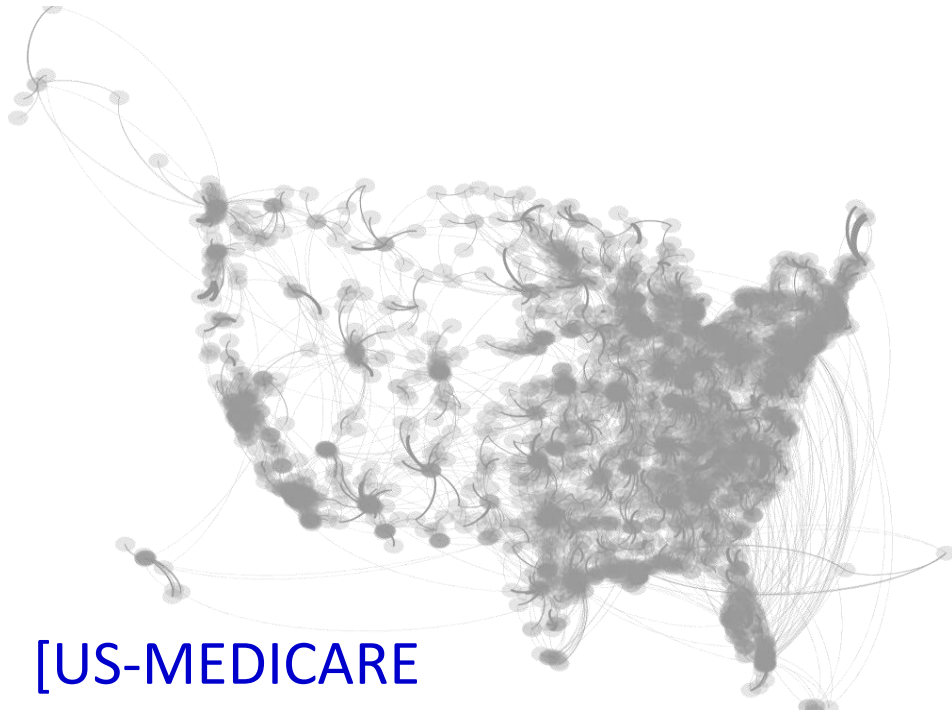
[AJPH 2007]



CDC data: Visualization of the first 35 tuberculosis (TB) patients and their 1039 contacts

Why do we care? (1: Epidemiology)

- Dynamical Processes over networks



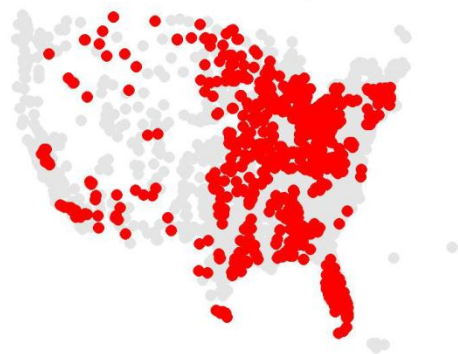
[US-MEDICARE
NETWORK 2005]

- Each circle is a hospital
- ~3000 hospitals
- More than 30,000 patients transferred

Problem: Given k units of disinfectant, whom to immunize?

Why do we care? (1: Epidemiology)

**~6x
fewer!**



CURRENT PRACTICE

[US-MEDICARE NETWORK 2005]



OUR METHOD

Why do we care? (2: Online Diffusion)



> 800m users, ~\$1B revenue [WSJ 2010]



~100m active users

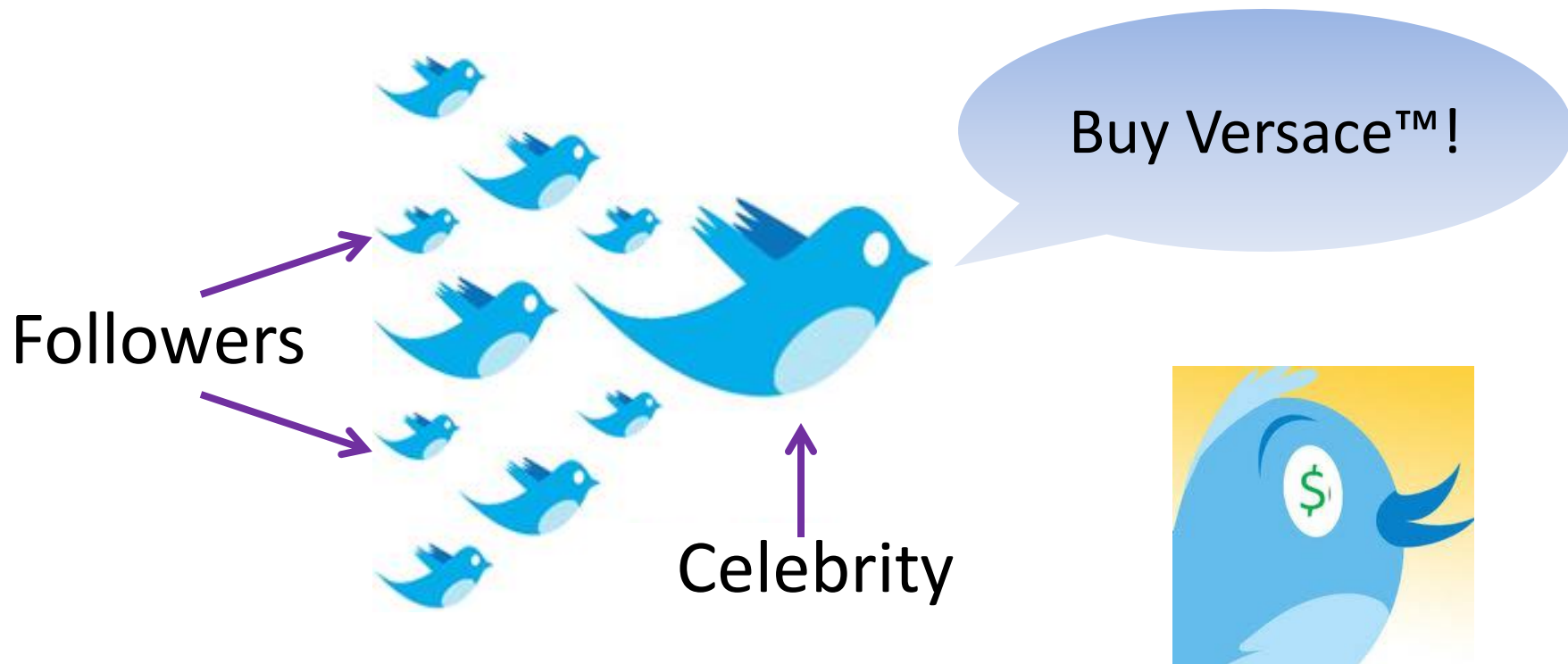


> 50m users



Why do we care? (2: Online Diffusion)

- Dynamical Processes over networks



Social Media Marketing

Prakash and Faloutsos 2012

Why do we care? (3: To change the world?)

- Dynamical Processes over networks



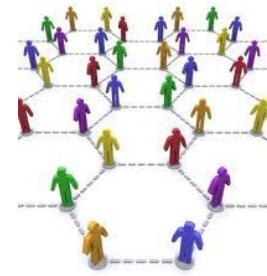
Social networks and Collaborative Action

High Impact – Multiple Settings

Q. How to squash rumors faster?



Q. How do opinions spread?

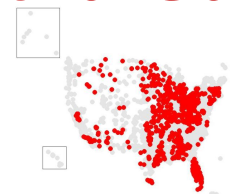


Q. How to market better?

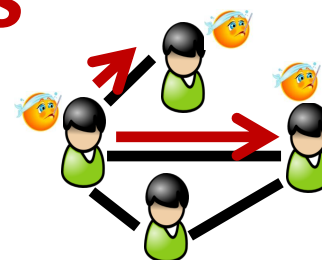


High Impact – Multiple Settings

Q. How to squash ~~rumors~~ **epidemic out-breaks** faster?



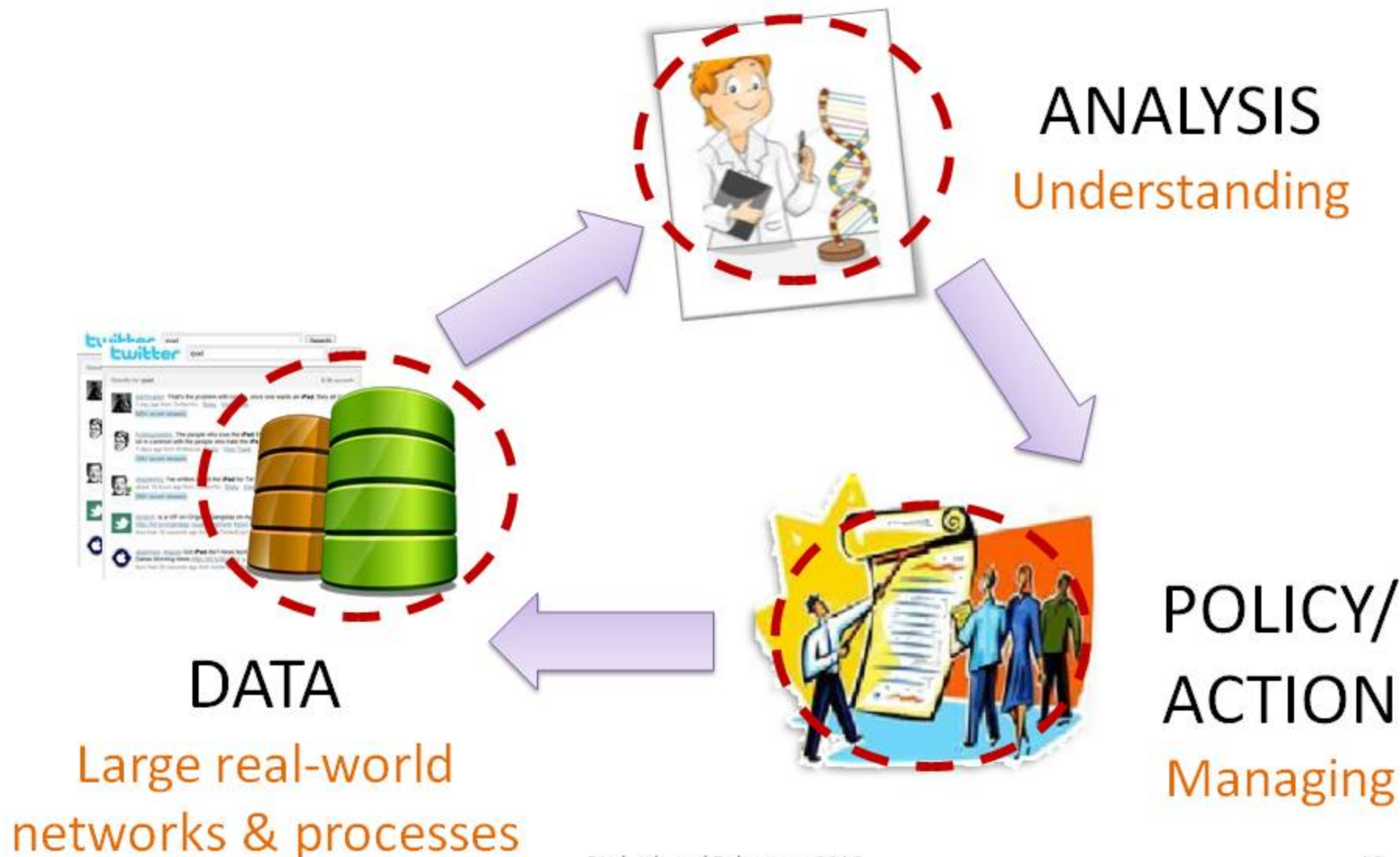
Q. How do ~~opinions~~ **products/viruses** spread?



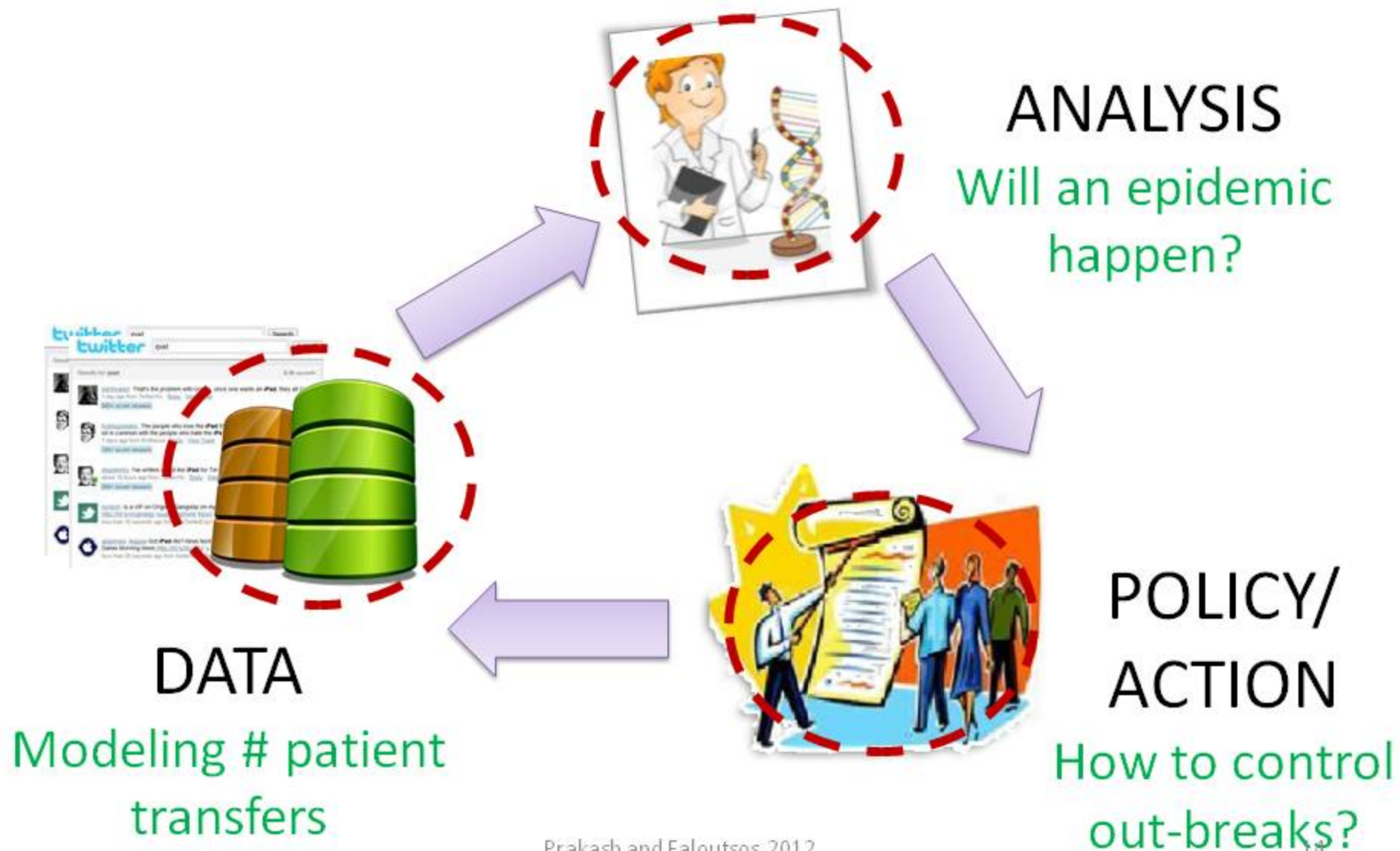
Q. How to ~~market~~ **transmit s/w patches** better?



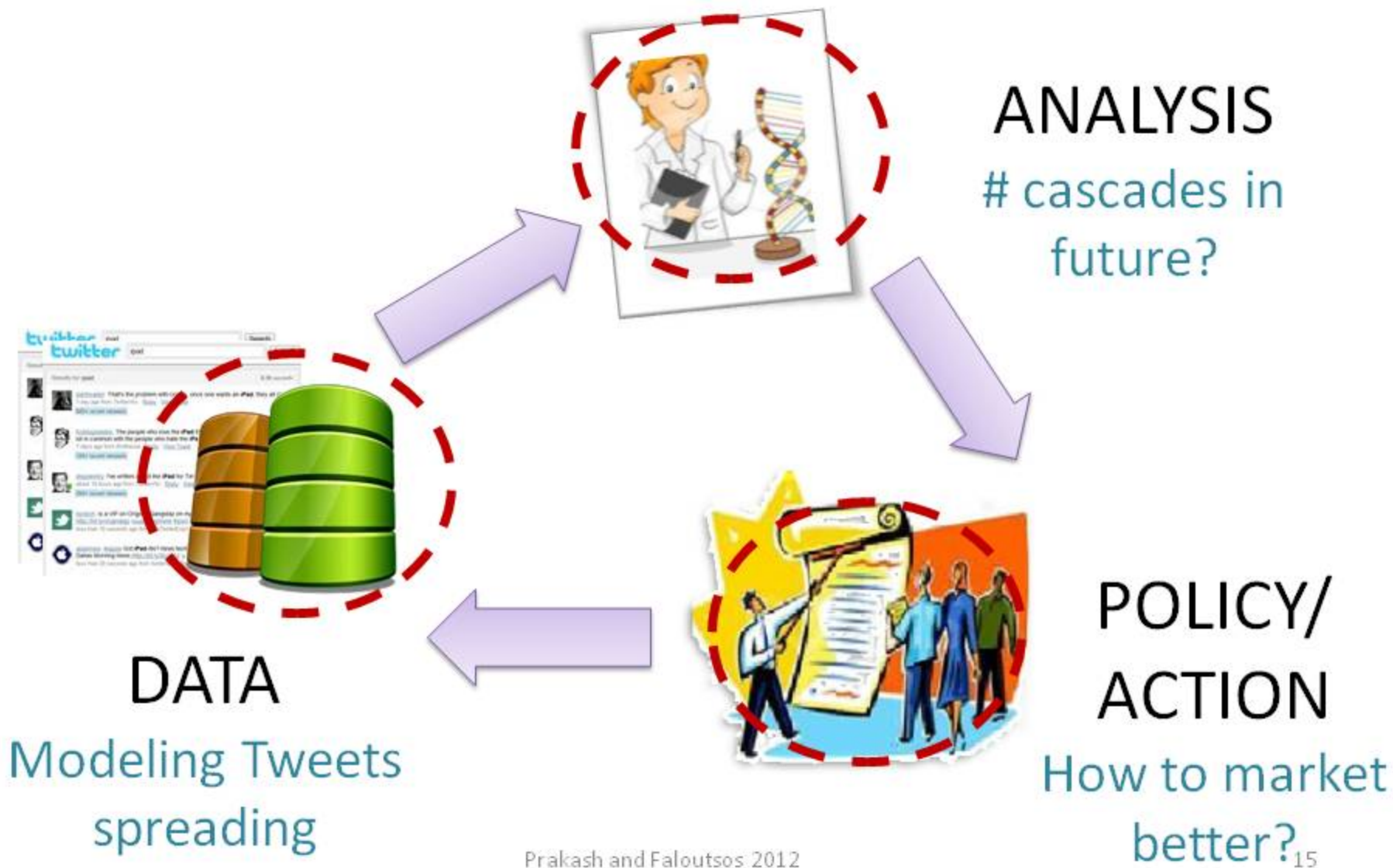
Research Theme



Research Theme – Public Health



Research Theme – Social Media



In this tutorial



ANALYSIS

Understanding

Given propagation models:

Q1: What is the epidemic threshold?

Q2: How do viruses compete?

In this tutorial



**POLICY/
ACTION**
Managing

Q3: How to immunize and control out-breaks better?

Q4: How to detect outbreaks?

Q5: Who are the culprits?

In this tutorial



DATA

Large real-world
networks & processes

Q6: How do cascades look like?

Q7: How does activity evolve over time?

Q8: How does external influence act?

Outline

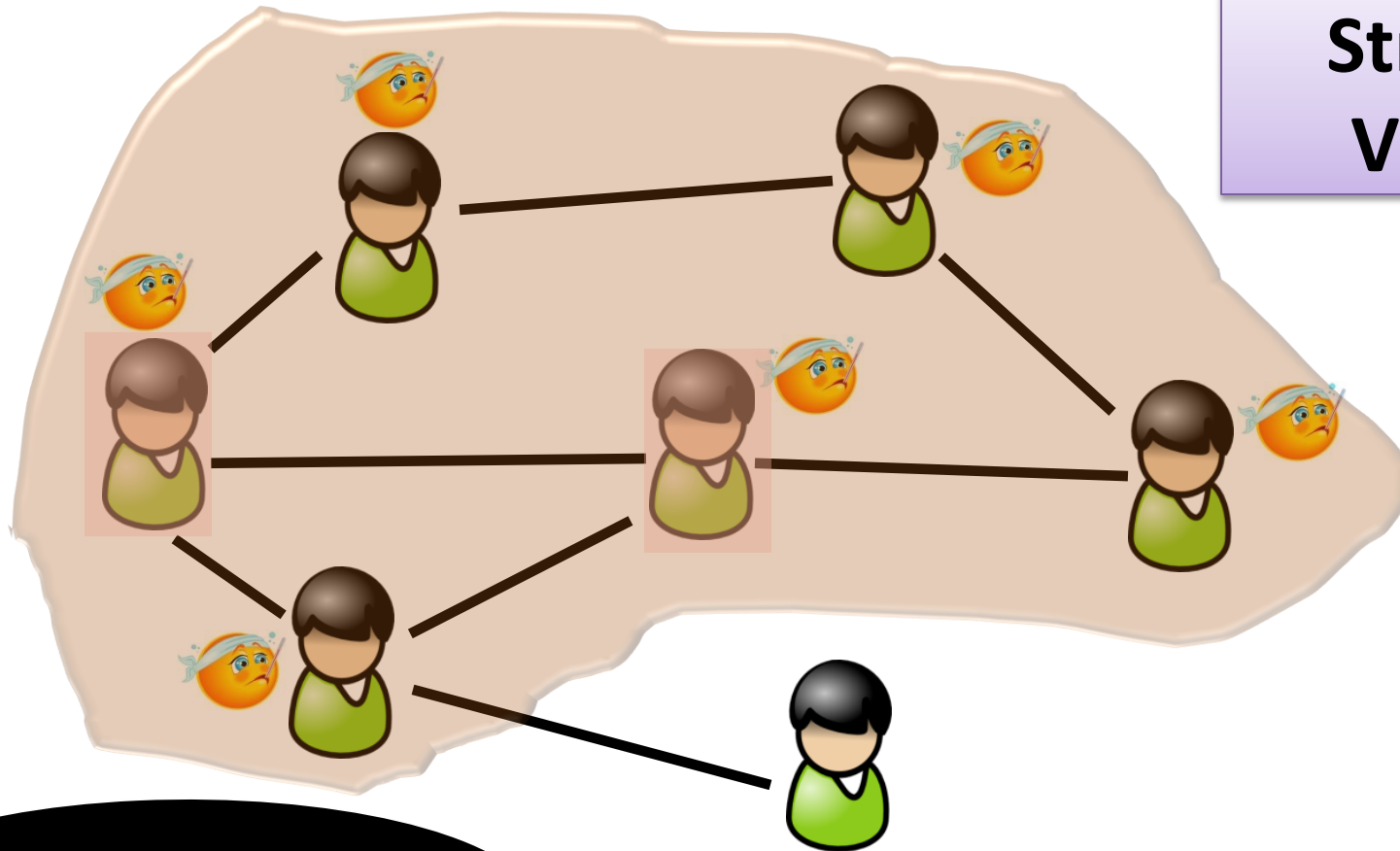
- Motivation
- **Part 1: Understanding Epidemics (Theory)**
- Part 2: Policy and Action (Algorithms)
- Part 3: Learning Models (Empirical Studies)
- Conclusion

Part 1: Theory

- **Q1: What is the epidemic threshold?**
- Q2: How do viruses compete?

A fundamental question

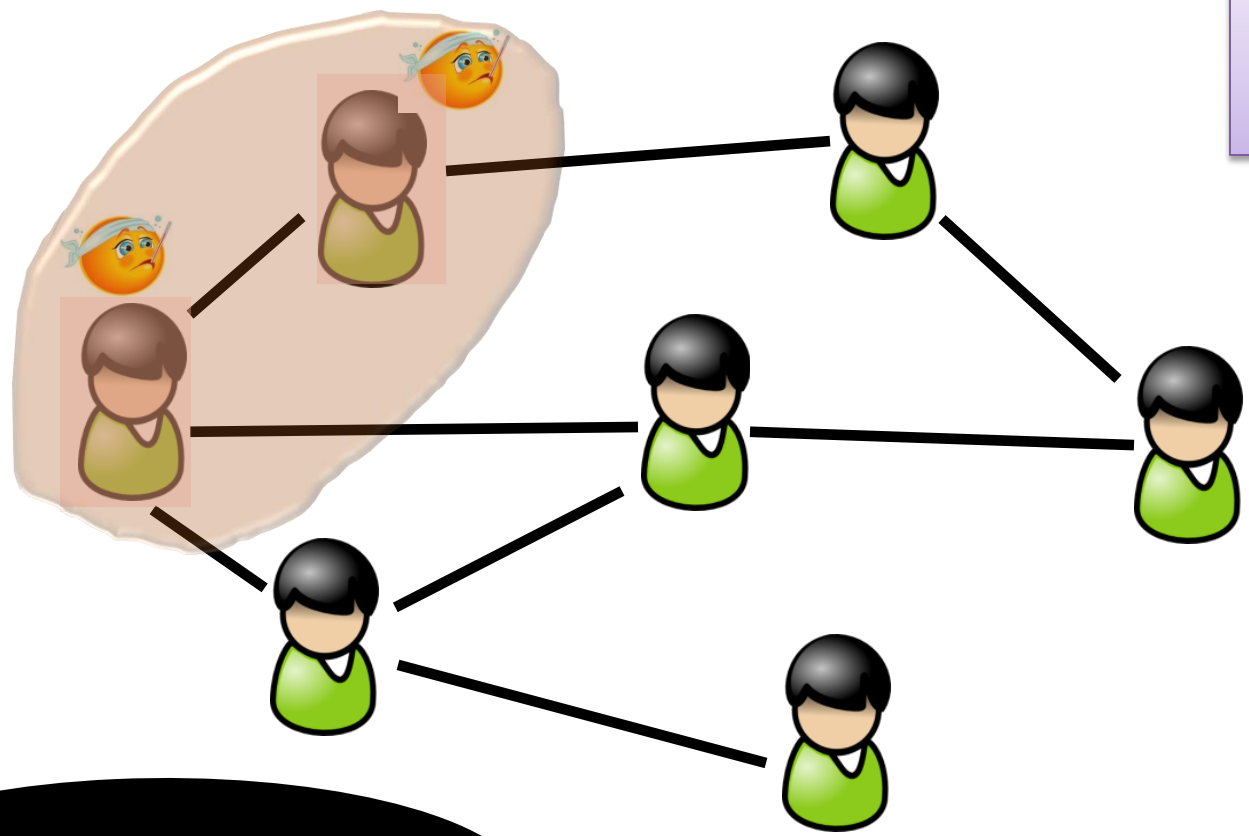
**Strong
Virus**



Epidemic?

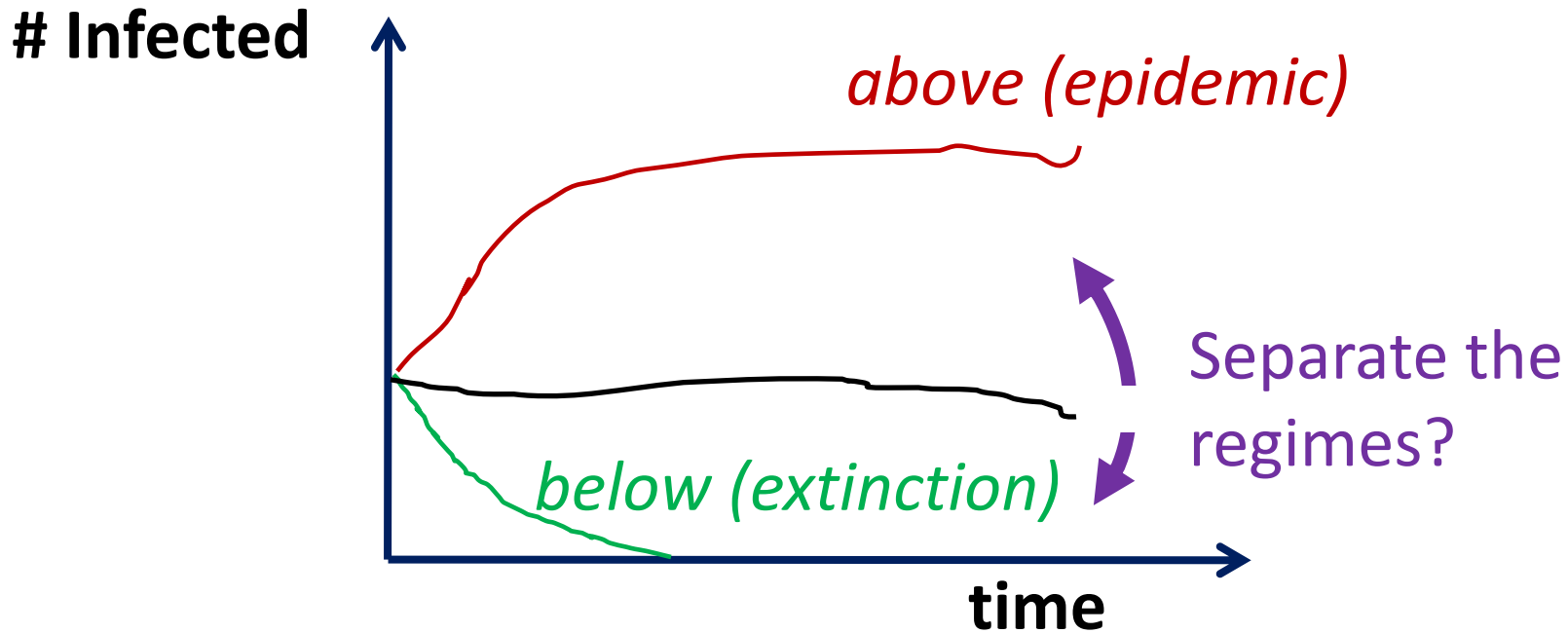
example (static graph)

Weak Virus



Epidemic?

Problem Statement



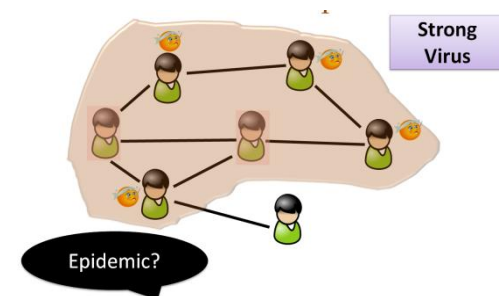
Find, a condition under which

- virus will die out exponentially quickly*
- regardless of initial infection condition*

Threshold (static version)

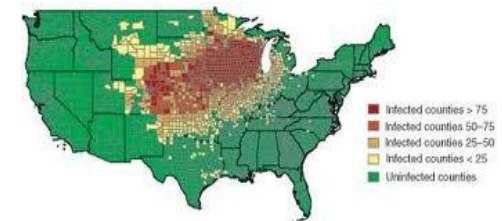
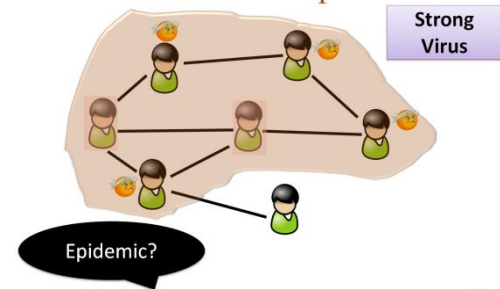
Problem Statement

- Given:
 - Graph G , *and*
 - Virus specs (attack prob. etc.)
- Find:
 - A condition for virus extinction/invasion



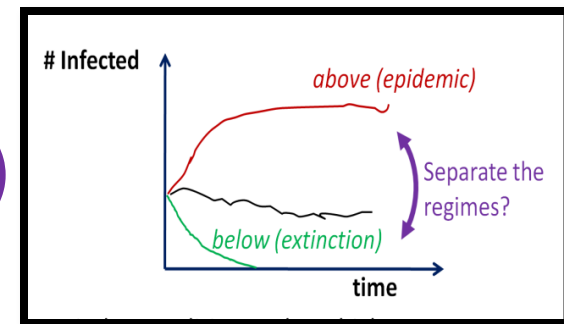
Threshold: Why important?

- Accelerating simulations
- Forecasting ('What-if' scenarios)
- Design of contagion and/or topology
- A great handle to manipulate the spreading
 - Immunization
 - Maximize collaboration
-



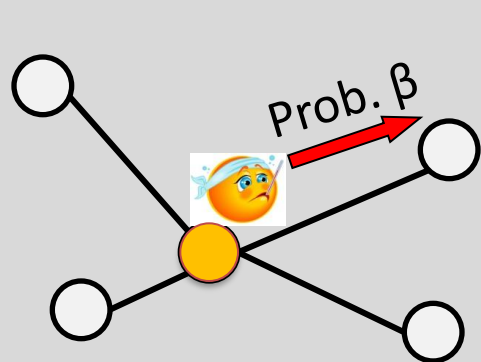
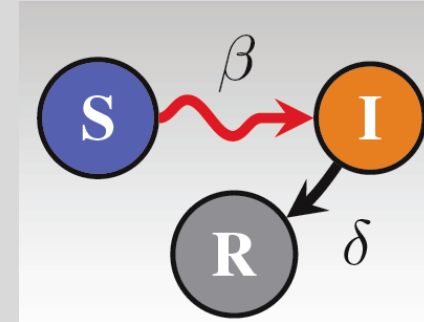
Part 1: Theory

- Q1: What is the epidemic threshold?
 - **Background**
 - Result and Intuition (Static Graphs)
 - Proof Ideas (Static Graphs)
 - Bonus: Dynamic Graphs
- Q2: How do viruses compete?

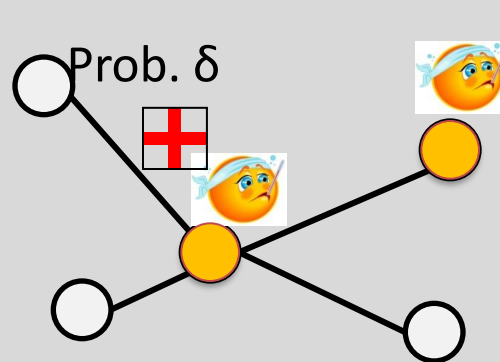


"SIR" model: life immunity (mumps)

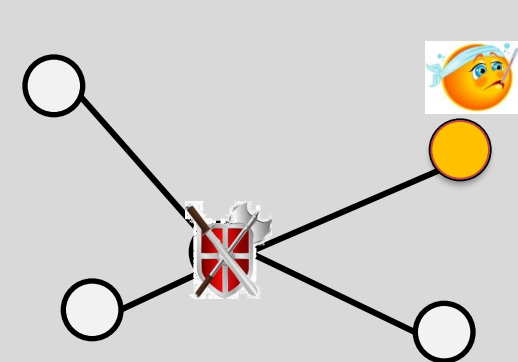
- Each node in the graph is in one of three states
 - **S**usceptible (i.e. healthy) ○
 - **I**nfected ●
 - **R**emoved (i.e. can't get infected again) 🏳️



$t = 1$



$t = 2$



$t = 3$

Terminology: continued

- Other virus propagation models (“VPM”)
 - SIS : susceptible-infected-susceptible, flu-like
 - SIRS : **temporary** immunity, like pertussis
 - SEIR : mumps-like, with virus **incubation**
(E = Exposed)

.....
- Underlying contact-network – ‘who-can-infect-whom’

Related Work

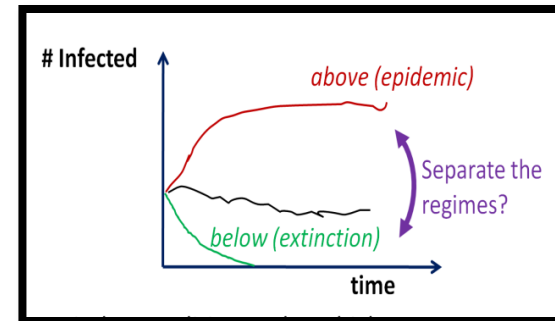
- ❑ R. M. Anderson and R. M. May. Infectious Diseases of Humans. Oxford University Press, 1991.
- ❑ A. Barrat, M. Barthélemy, and A. Vespignani. Dynamical Processes on Complex Networks. Cambridge University Press, 2010.
- ❑ F. M. Bass. A new product growth for model consumer durables. Management Science, 15(5):215–227, 1969.
- ❑ D. Chakrabarti, Y. Wang, C. Wang, J. Leskovec, and C. Faloutsos. Epidemic thresholds in real networks. ACM TISSEC, 10(4), 2008.
- ❑ D. Easley and J. Kleinberg. Networks, Crowds, and Markets: Reasoning About a Highly Connected World. Cambridge University Press, 2010.
- ❑ A. Ganesh, L. Massoulié, and D. Towsley. The effect of network topology in spread of epidemics. IEEE INFOCOM, 2005.
- ❑ Y. Hayashi, M. Minoura, and J. Matsukubo. Recoverable prevalence in growing scale-free networks and the effective immunization. arXiv:cond-at/0305549 v2, Aug. 6 2003.
- ❑ H. W. Hethcote. The mathematics of infectious diseases. SIAM Review, 42, 2000.
- ❑ H. W. Hethcote and J. A. Yorke. Gonorrhoea transmission dynamics and control. Springer Lecture Notes in Biomathematics, 46, 1984.
- ❑ J. O. Kephart and S. R. White. Directed-graph epidemiological models of computer viruses. IEEE Computer Society Symposium on Research in Security and Privacy, 1991.
- ❑ J. O. Kephart and S. R. White. Measuring and modeling computer virus prevalence. IEEE Computer Society Symposium on Research in Security and Privacy, 1993.
- ❑ R. Pastor-Santorrás and A. Vespignani. Epidemic spreading in scale-free networks. Physical Review Letters 86, 14, 2001.
- ❑
- ❑
- ❑

All are about *either*:

- **Structured topologies** (cliques, block-diagonals, hierarchies, random)
- **Specific virus propagation models**
- **Static graphs**

Part 1: Theory

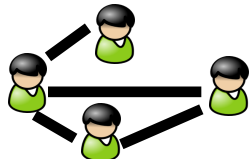
- Q1: What is the epidemic threshold?
 - Background
 - **Result and Intuition (Static Graphs)**
 - Proof Ideas (Static Graphs)
 - Bonus: Dynamic Graphs
- Q2: How do viruses compete?



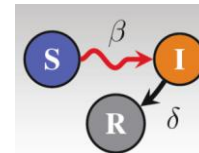
How should the answer look like?

- Answer should depend on:

- Graph



- Virus Propagation Model (VPM)



- But how??

- Graph – average degree? max. degree? diameter?

- VPM – which parameters?

- How to combine – linear? quadratic? exponential?

$$\beta d_{avg} + \delta \sqrt{diameter} ? (\beta^2 d_{avg}^2 - \delta d_{avg}) / d_{max} ? \dots$$

Static Graphs: Our Main Result

- Informally,

For,

➤ any arbitrary topology (adjacency matrix A)

➤ any virus propagation model (VPM) in standard literature

the epidemic threshold depends only

1. on the λ , first eigenvalue of A , and
2. some constant C_{VPM} , determined by the virus propagation model

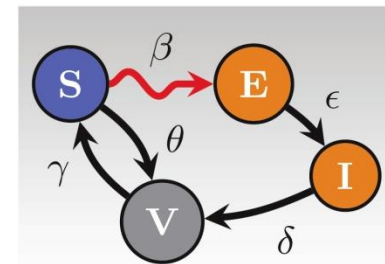
λ

C_{VPM}

No epidemic if
 $\lambda * C_{VPM} < 1$

Our thresholds for some models

- $s = \text{effective strength}$
- $s < 1$: *below threshold*



Models	Effective Strength (s)	Threshold (tipping point)
SIS, SIR, SIRS, SEIR	$s = \lambda \cdot \left(\frac{\beta}{\delta} \right)$	$s = 1$
SIV, SEIV	$s = \lambda \cdot \left(\frac{\beta\gamma}{\delta(\gamma + \theta)} \right)$	
$SI_1I_2V_1V_2$ (<u>H.I.V.</u>)	$s = \lambda \cdot \left(\frac{\beta_1v_2 + \beta_2\varepsilon}{v_2(\varepsilon + v_1)} \right)$	

Our result: Intuition for λ

“Official” definition:

- Let A be the adjacency matrix. Then λ is the root with the largest magnitude of the characteristic polynomial of A [$\det(A - xI)$].
- Doesn't give much intuition!

“Un-official” Intuition 😊

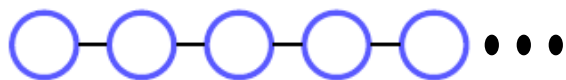
- $\lambda \sim \#$ paths in the graph

$$A^k \approx \lambda^k \cdot u \cdot u$$

$A^k(i, j) = \#$ of paths $i \rightarrow j$
of length k

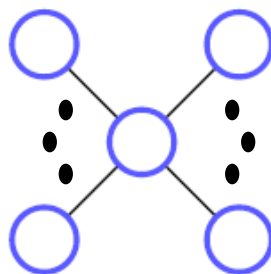
Largest Eigenvalue (λ)

better connectivity \longrightarrow higher λ



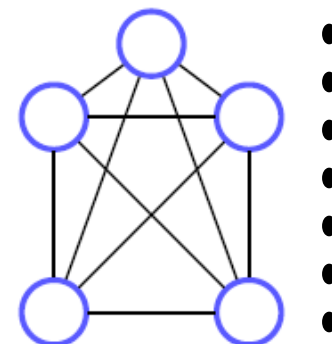
$$\lambda \approx 2$$

(a) Chain



$$\lambda = \sqrt{N}$$

(b) Star



$$\lambda = N-1$$

(c) Clique

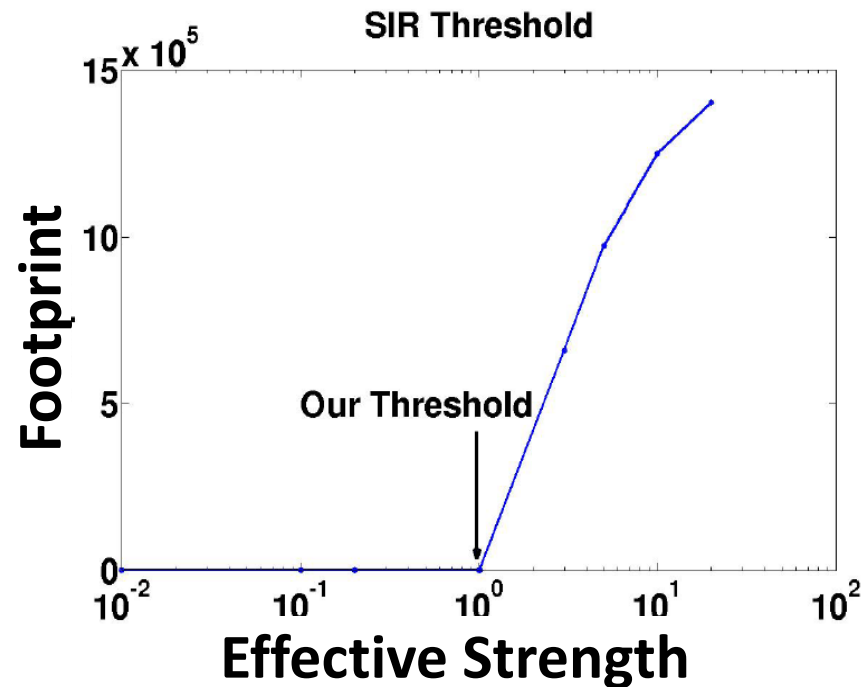
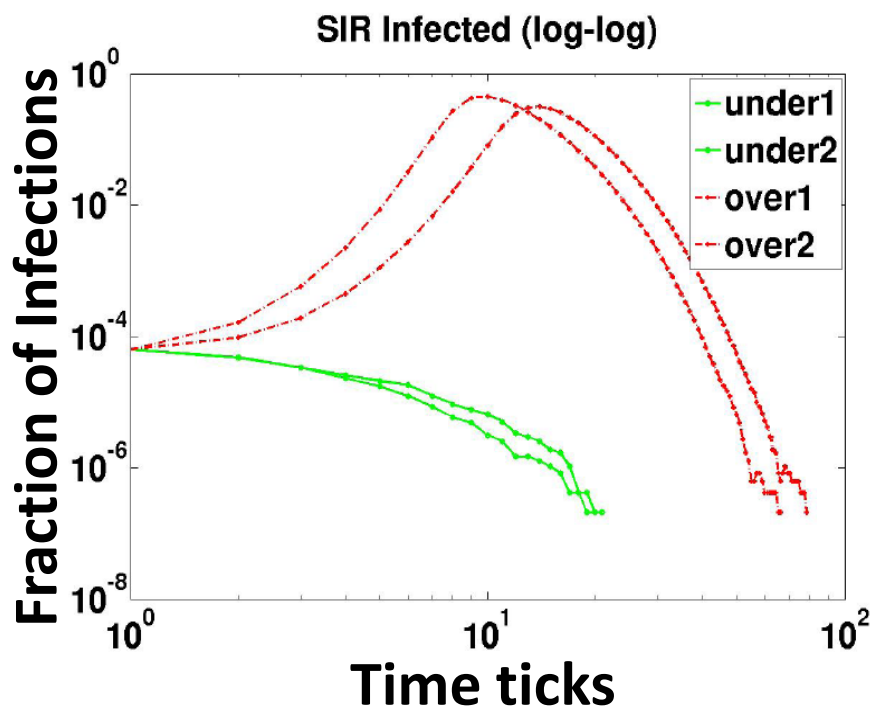
$$\lambda \approx 2$$

$$\lambda = 31.67$$

$$\lambda = 999$$

$N = 1000$

Examples: Simulations – SIR (mumps)



(a) Infection profile

(b) "Take-off" plot

PORTLAND graph

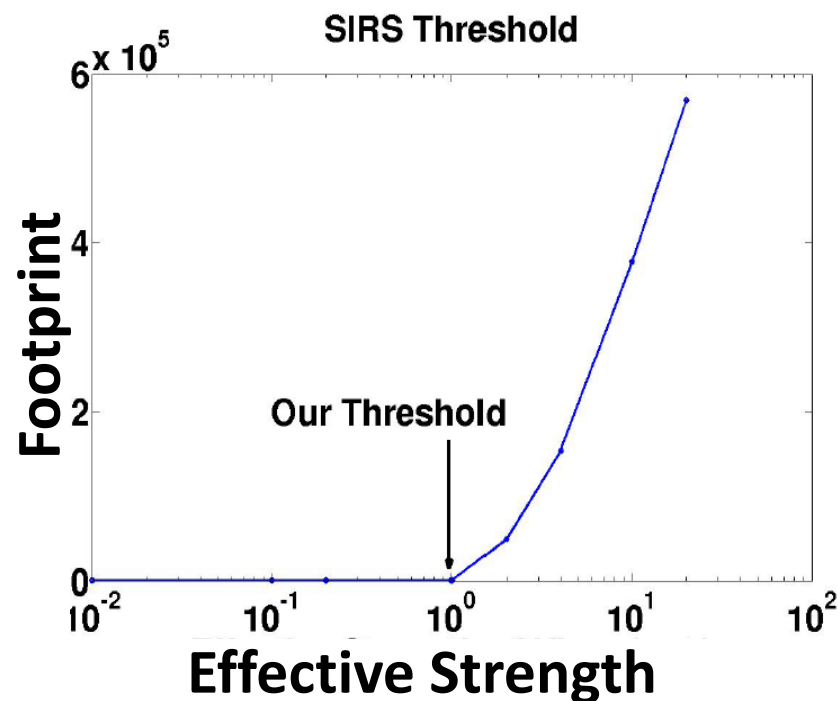
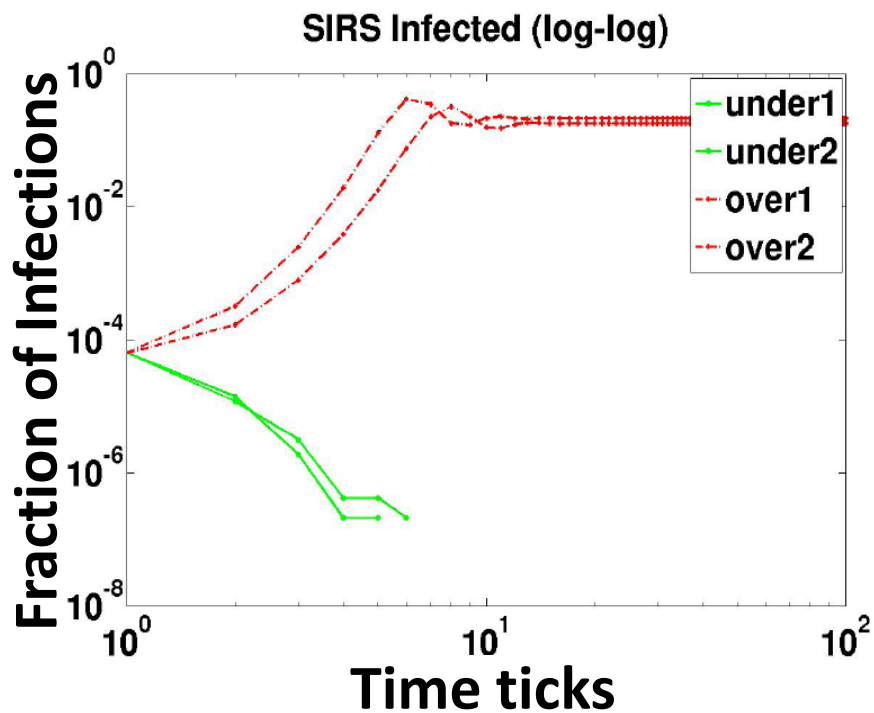


31 million links, 6 million nodes

Prakash and Faloutsos 2012



Examples: Simulations – SIRS (pertusis)



(a) Infection profile

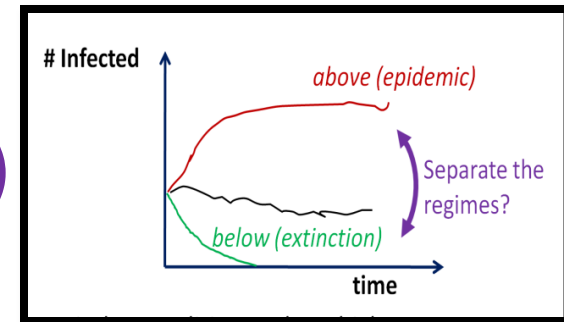
(b) "Take-off" plot

PORTLAND graph

31 million links, 6 million nodes

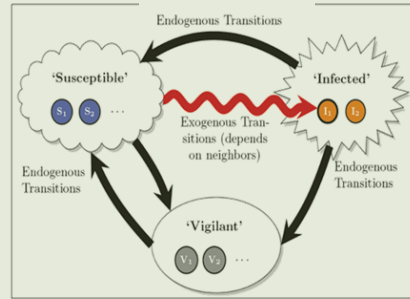
Part 1: Theory

- Q1: What is the epidemic threshold?
 - Background
 - Result and Intuition (Static Graphs)
 - **Proof Ideas (Static Graphs)**
 - Bonus: Dynamic Graphs
- Q2: How do viruses compete?



Proof Sketch

Model	Used for
SIR	Mumps
SIS	Flu
SIRS	Pertussis
SEIR	Varicella
.....	
SICR	Tuberculosis
MSIR	Measles
SIV	Sensor Stability
$SI_1I_2V_1V_2$	H.I.V.
.....	



General VPM structure

47

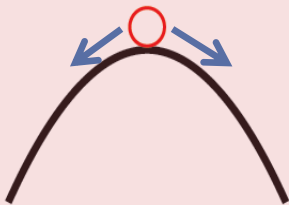
Model-based



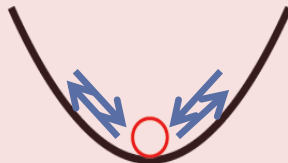
$$\lambda^* C_{VPM} < 1$$

Graph-based

Topology and stability



(A) Unstable



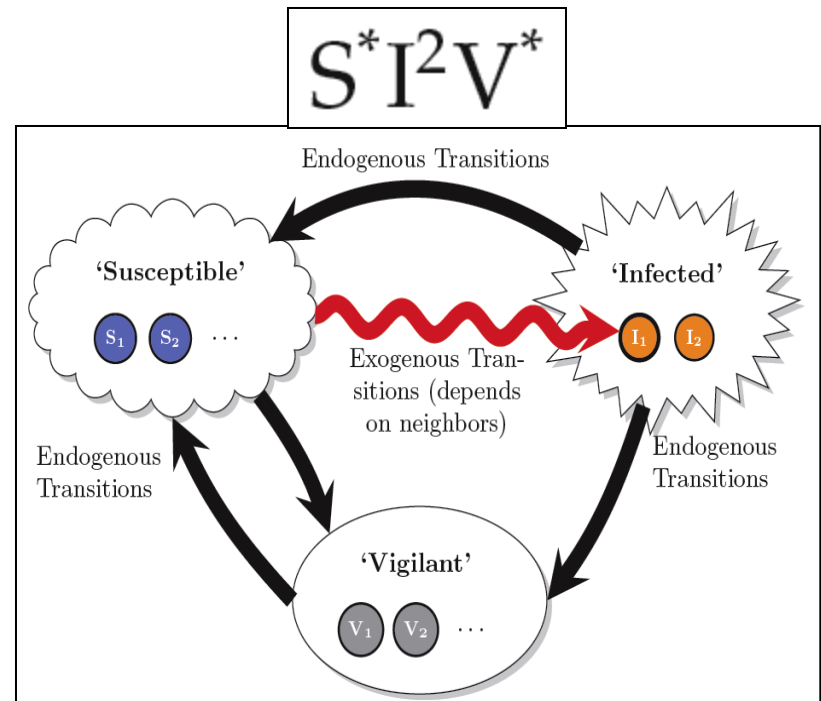
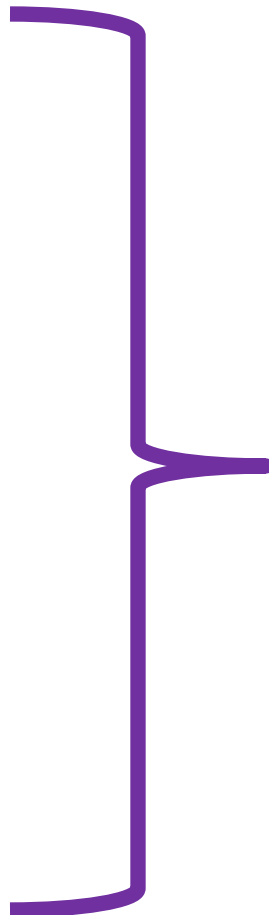
(B) Stable



(C) Neutral (at threshold)

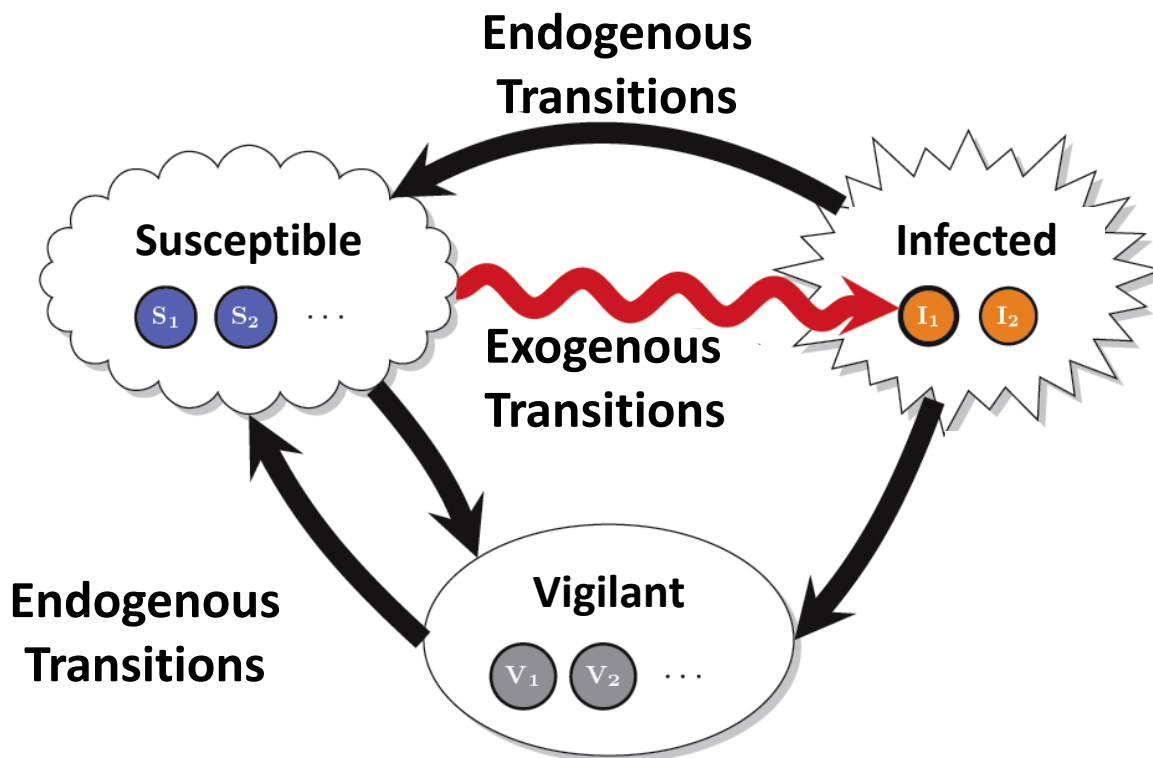
Models and more models

Model	Used for
SIR	Mumps
SIS	Flu
SIRS	Pertussis
SEIR	Chicken-pox
.....	
SICR	Tuberculosis
MSIR	Measles
SIV	Sensor Stability
$SI_1 I_2 V_1 V_2$	H.I.V.
.....	

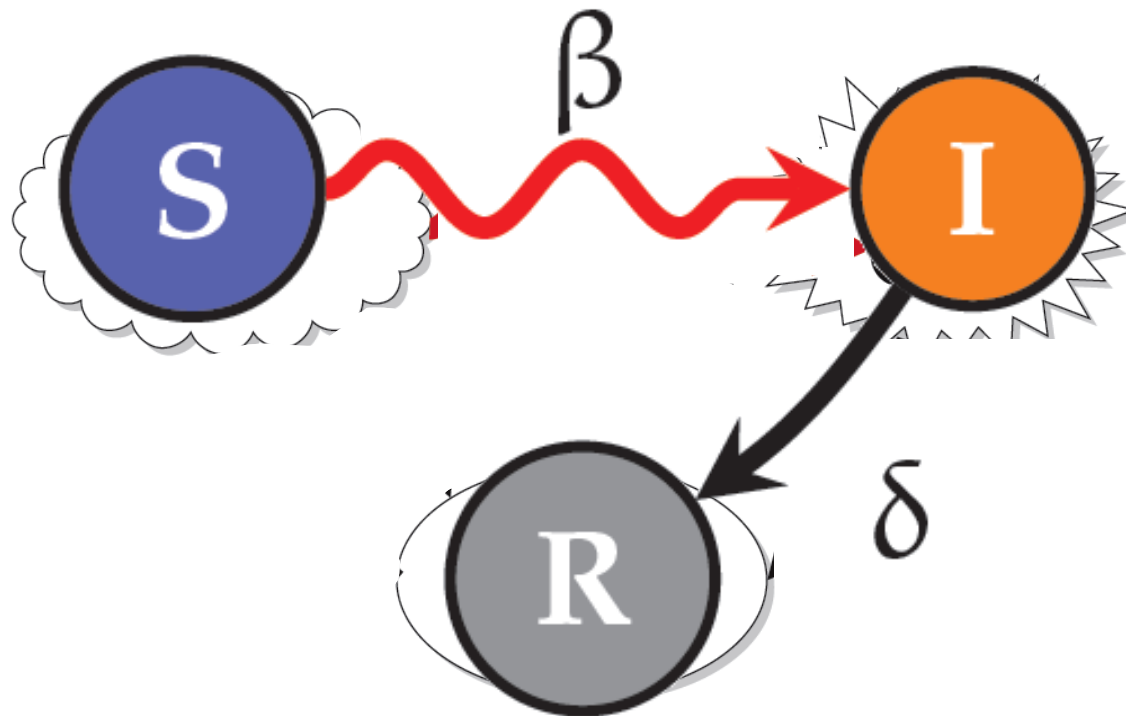


Ingredient 1: Our generalized model

$$S^* I^2 V^* \quad (S^* I^* V^* ?)$$

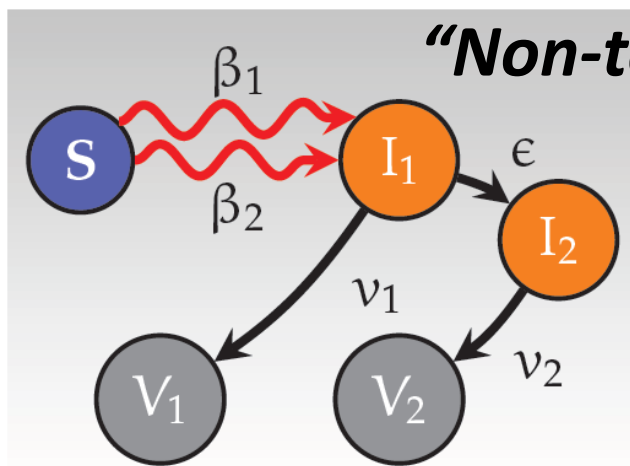


Special case: SIR



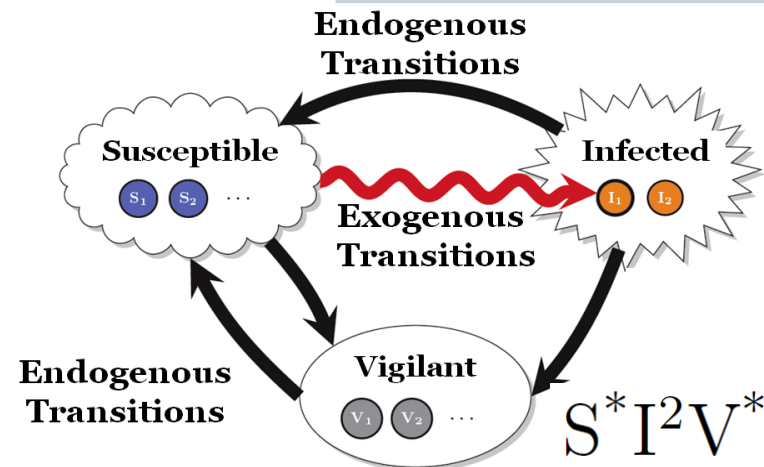
Special case: H.I.V.

$$SI_1I_2V_1V_2$$



“Non-terminal”

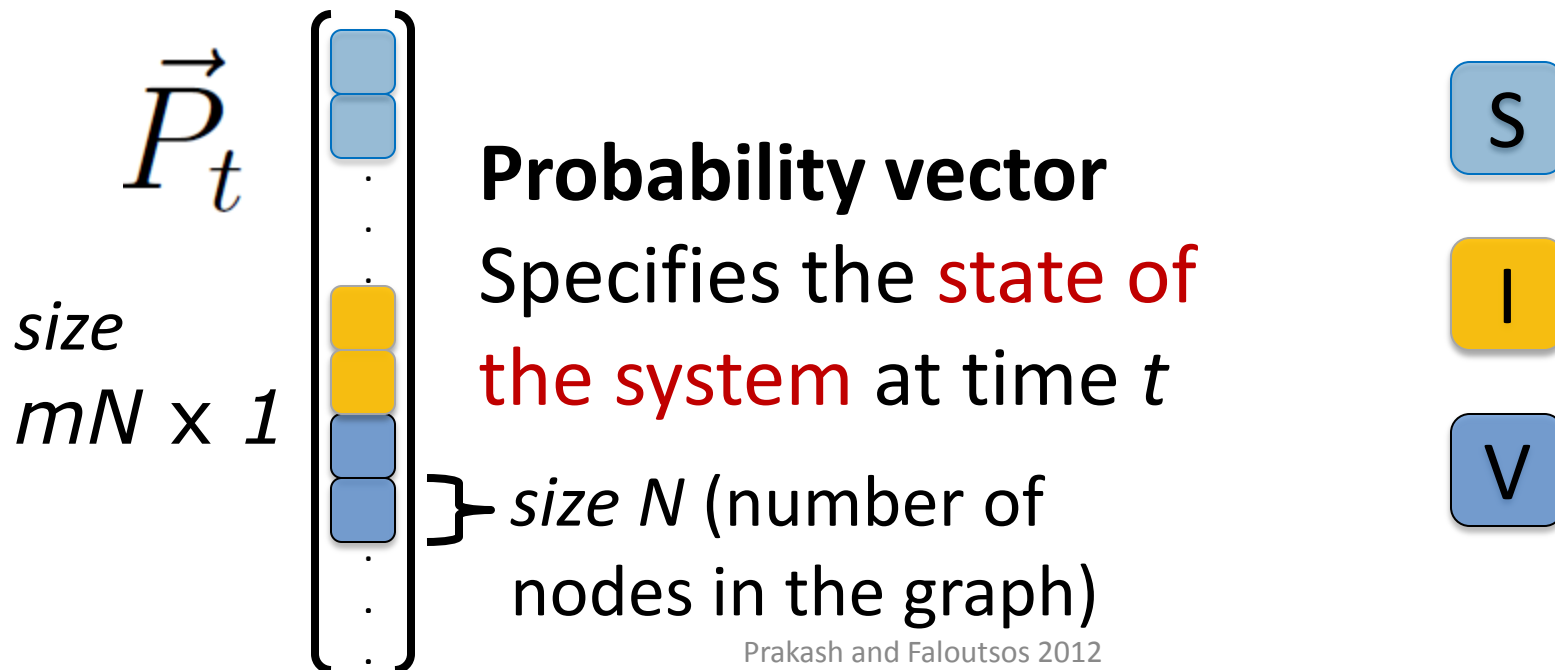
“Terminal”



Multiple Infectious,
Vigilant states

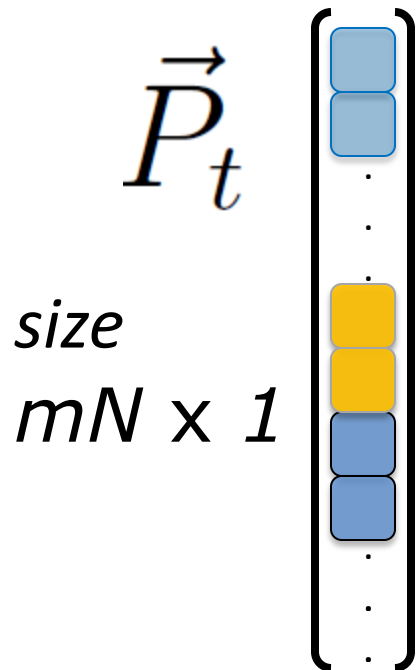
Ingredient 2: NLDS+Stability

- View as a NLDS $\vec{P}_{t+1} = \mathcal{G}(\vec{P}_t)$
 - discrete time
 - non-linear dynamical system (NLDS)



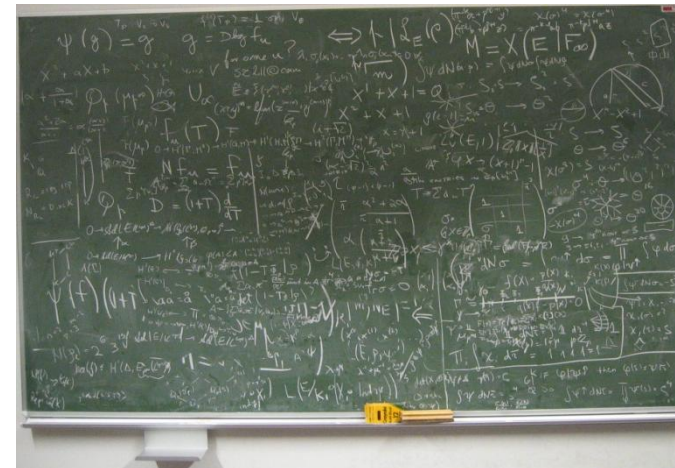
Ingredient 2: NLDS + Stability

- View as a NLDS $\vec{P}_{t+1} = \mathcal{G}(\vec{P}_t)$
 - discrete time
 - non-linear dynamical system (NLDS)



Non-linear function
 Explicitly **gives the evolution** of system

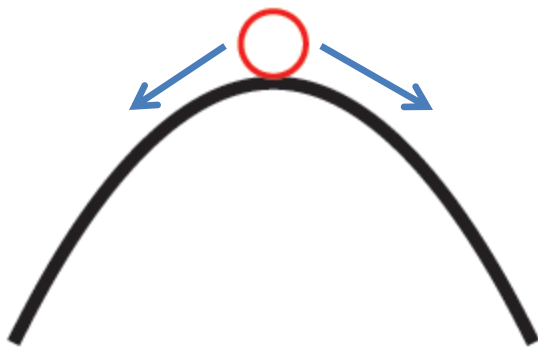
$$\mathcal{G} : \mathbb{R}^{mN} \rightarrow \mathbb{R}^{mN}$$



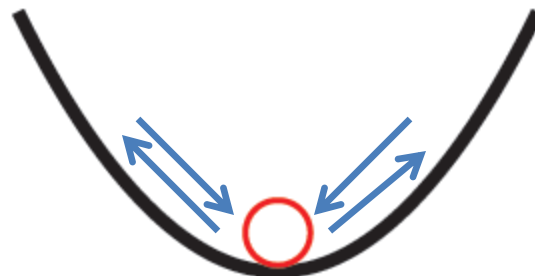
Ingredient 2: NLDS + Stability

- View as a NLDS $\vec{P}_{t+1} = \mathcal{G}(\vec{P}_t)$
 - discrete time
 - non-linear dynamical system (NLDS)

- Threshold \rightarrow Stability of NLDS



(A) Unstable

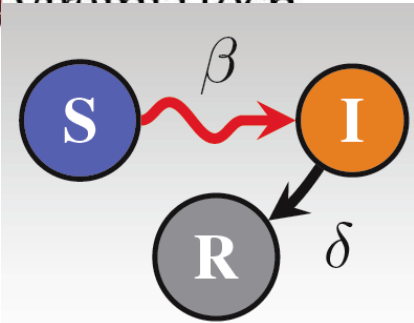


(B) Stable



(C) Neutral (at threshold)

Special case: SIR



$$\vec{P}_{t+1}$$

size
 $3N \times 1$

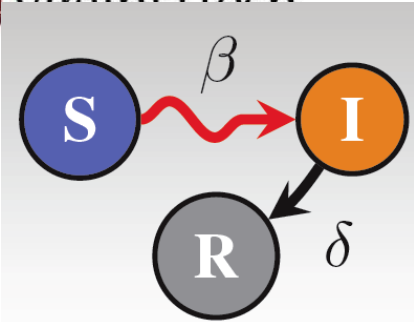
$$\mathcal{G} : \mathbb{R}^{3N} \rightarrow \mathbb{R}^{3N}$$

$$\begin{aligned} P_{S,i,t+1} &= P_{S,i,t} \zeta_{i,t}(I) \\ P_{I,i,t+1} &= P_{S,i,t} (1 - \zeta_{i,t}(I)) + (1 - \delta) P_{I,i,t} \\ P_{R,i,t+1} &= \delta P_{I,i,t} + P_{R,i,t} \end{aligned}$$

$$\vec{P}_t$$

$\zeta_{i,t}(I)$ = probability that node i is not attacked by any of its infectious neighbors

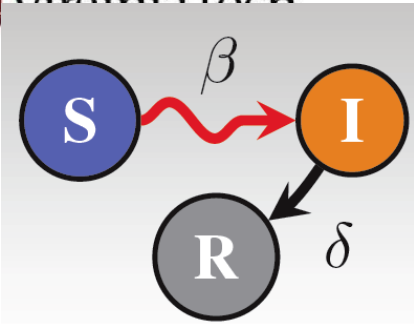
Special case: SIR



$$\begin{array}{l}
 \vec{P}_{t+1} \\
 \text{size} \\
 3N \times 1
 \end{array}
 \begin{array}{|c|}
 \hline
 S \\
 \hline
 I \\
 \hline
 R \\
 \hline
 \end{array}
 =
 \mathcal{G} : \mathbb{R}^{3N} \rightarrow \mathbb{R}^{3N}
 \begin{array}{|c|}
 \hline
 S \\
 \hline
 I \\
 \hline
 R \\
 \hline
 \end{array}
 \vec{P}_t$$

$$\begin{aligned}
 P_{S,i,t+1} &= P_{S,i,t} \zeta_{i,t}(I) \\
 P_{I,i,t+1} &= P_{S,i,t} (1 - \zeta_{i,t}(I)) + (1 - \delta) P_{I,i,t} \\
 P_{R,i,t+1} &= \delta P_{I,i,t} + P_{R,i,t}
 \end{aligned}$$

NLDS

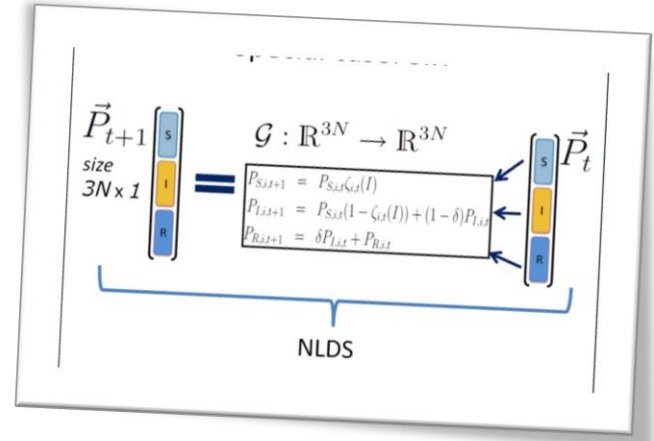


Fixed Point

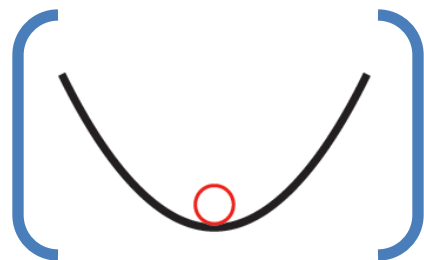
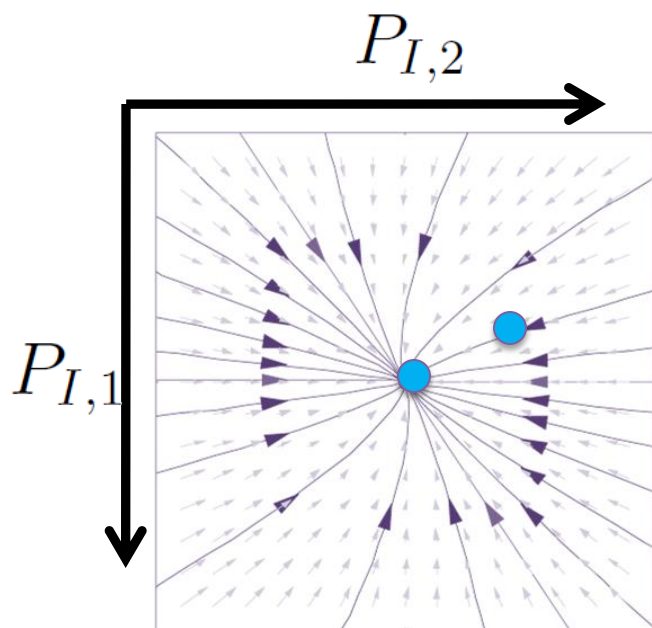
$$\vec{P}^* = \begin{pmatrix} 1 \\ 1 \\ \cdot \\ 0 \\ 0 \\ \cdot \\ 0 \\ 0 \\ \cdot \end{pmatrix}$$

State when **no node is infected**

Q: Is it stable?

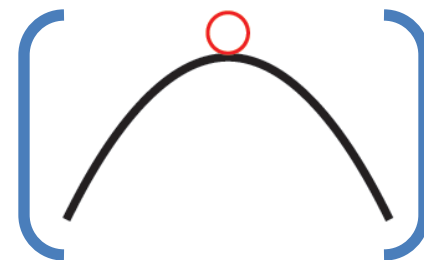
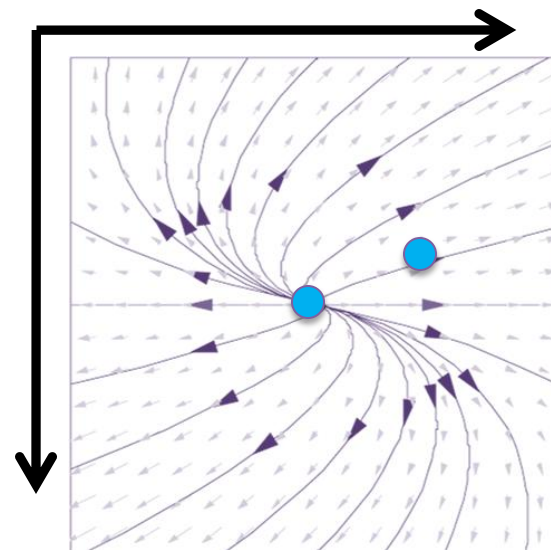


Stability for SIR



Stable

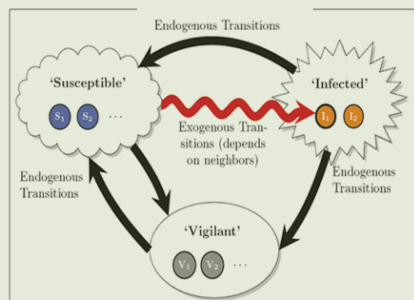
under threshold



Unstable

above threshold

Model	Used for
SIR	Mumps
SIS	Flu
SIRS	Pertussis
SEIR	Varicella
.....	
SICR	Tuberculosis
MSIR	Measles
SIV	Sensor Stability
$SI_1I_2V_1V_2$	H.I.V.
.....	

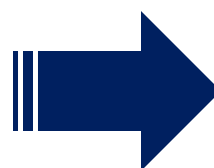


General VPM structure

47

See paper for full proof

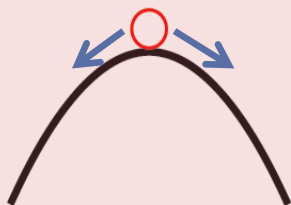
Model-based



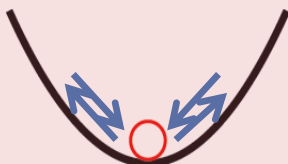
$$\lambda^* C_{VPM} < 1$$

Graph-based

Topology and stability



(A) Unstable



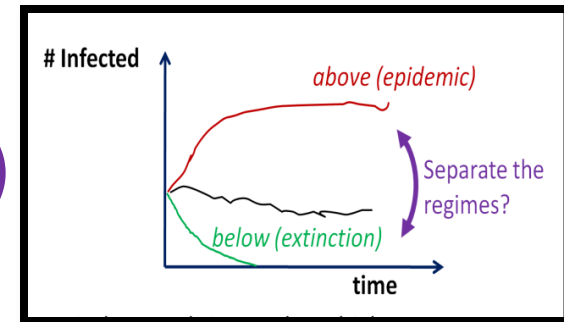
(B) Stable



(C) Neutral (at threshold)

Part 1: Theory

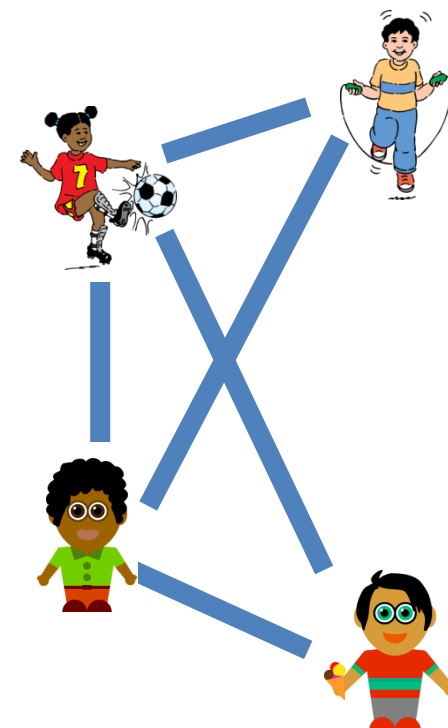
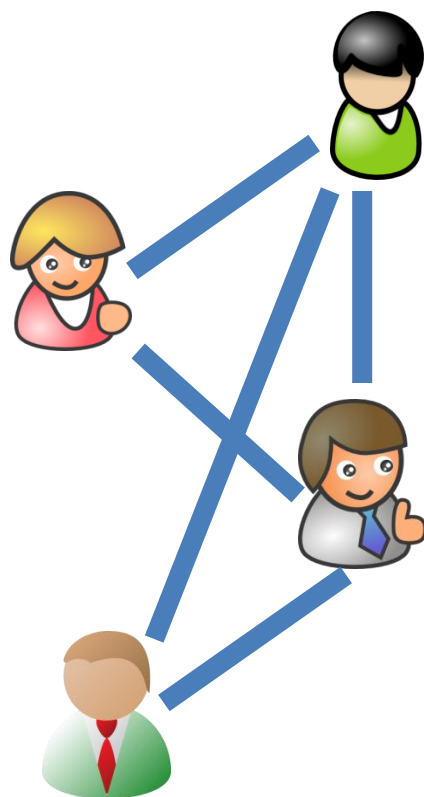
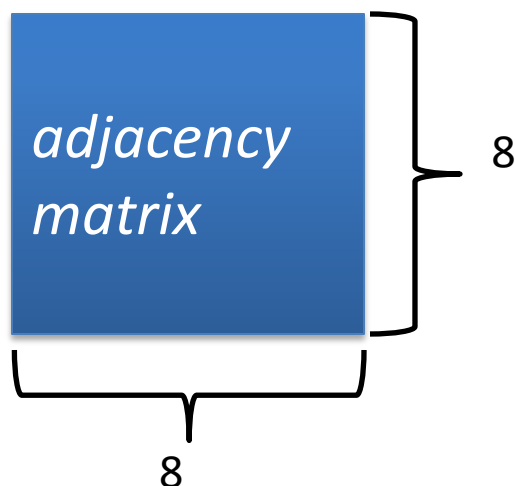
- Q1: What is the epidemic threshold?
 - Background
 - Result and Intuition (Static Graphs)
 - Proof Ideas (Static Graphs)
 - **Bonus: Dynamic Graphs**
- Q2: How do viruses compete?



Dynamic Graphs: Epidemic?

Alternating behaviors

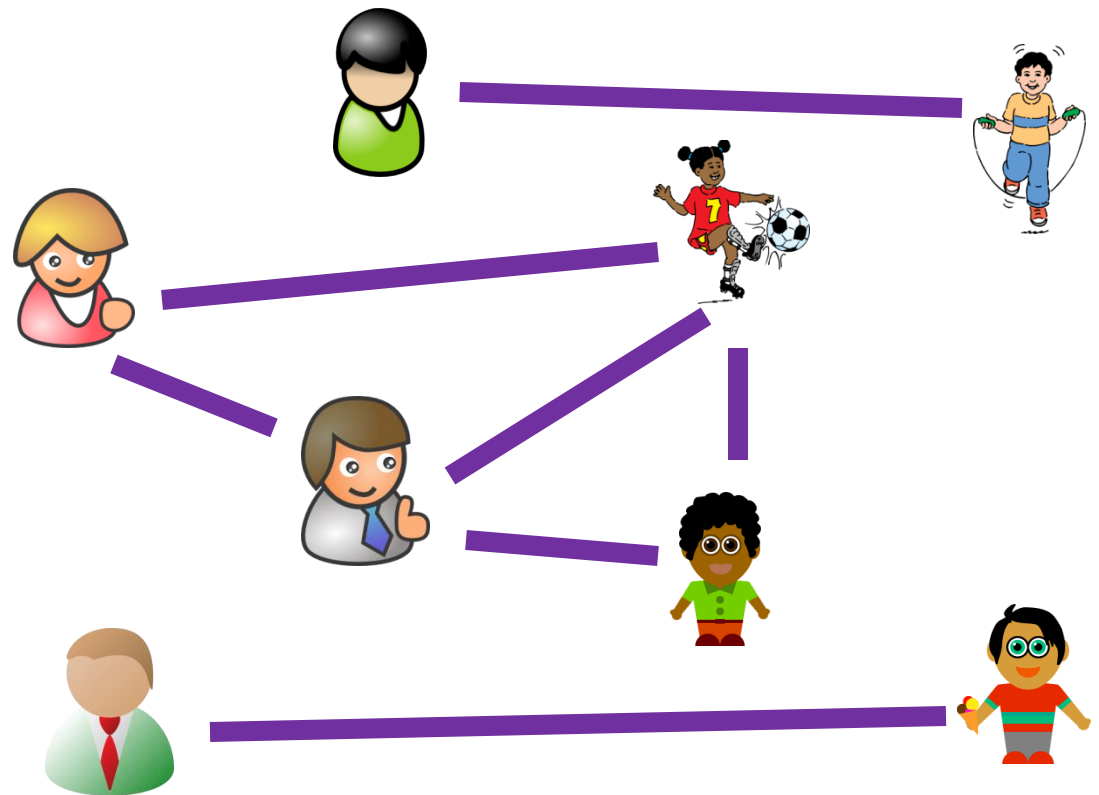
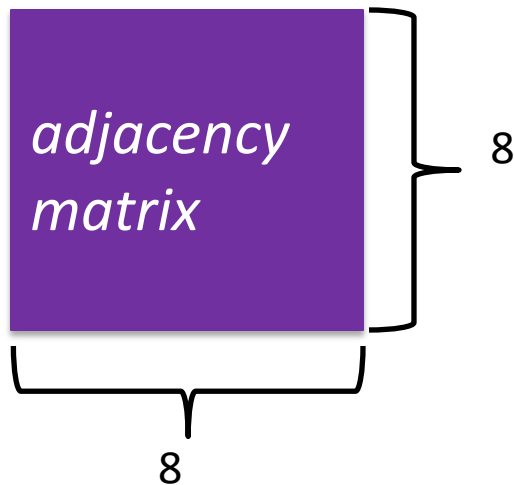
DAY
(e.g., work)



Dynamic Graphs: Epidemic?

NIGHT
(e.g., home)

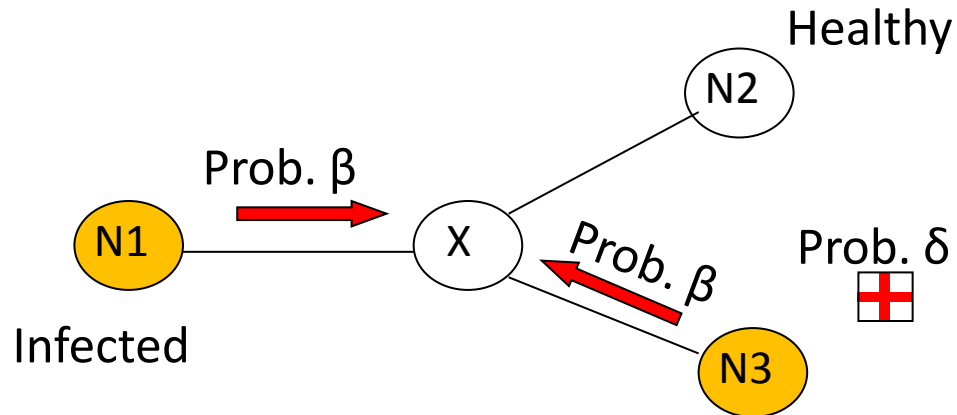
Alternating behaviors



Model Description

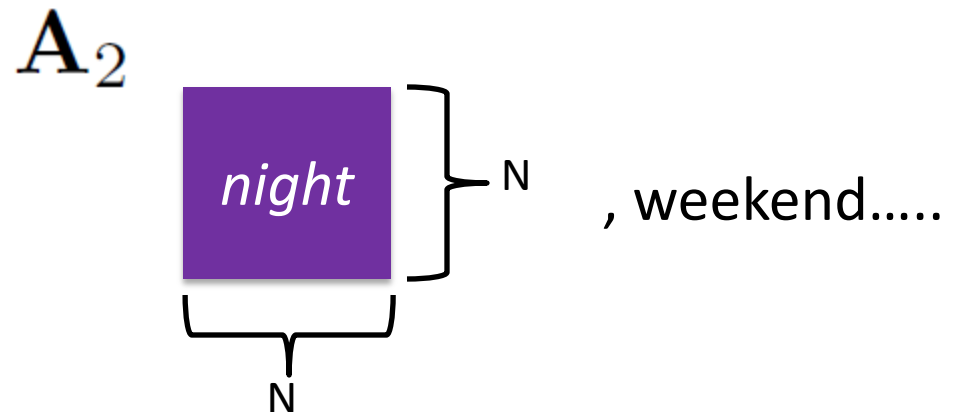
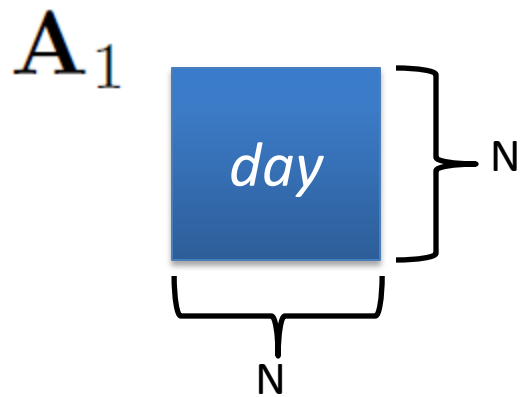
- SIS model

- recovery rate δ
- infection rate β



- Set of T arbitrary graphs

$$\{A_1, A_2 \dots, A_T\}$$



Our result: Dynamic Graphs Threshold

- Informally, *NO* epidemic if

$$\text{eig}(\mathbf{S}) = \lambda_{\mathbf{S}} < 1$$

Single number!
Largest eigenvalue of
The *system matrix* \mathbf{S}

$$\mathbf{S} = \prod_i \mathbf{S}_i$$

Details

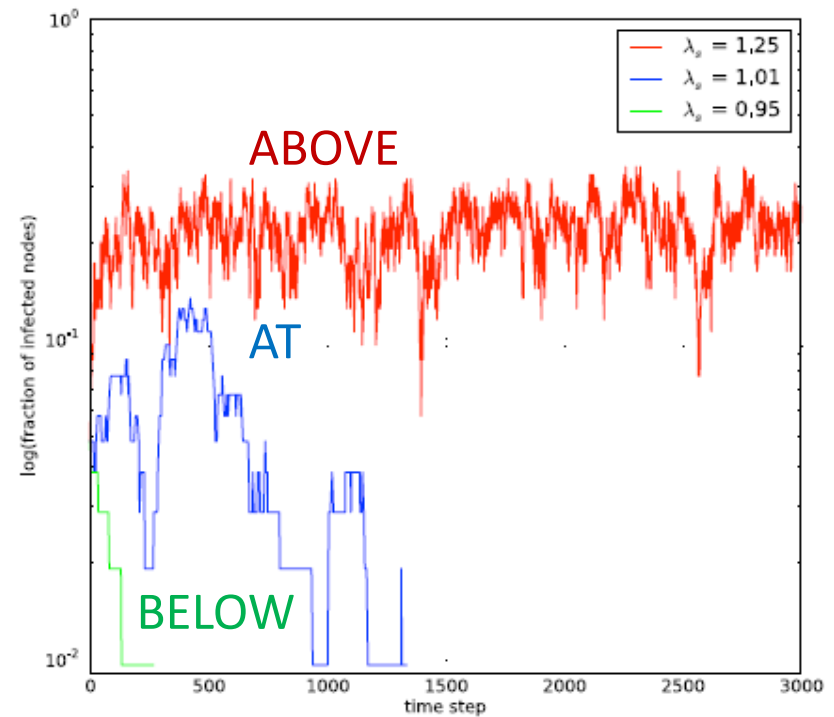
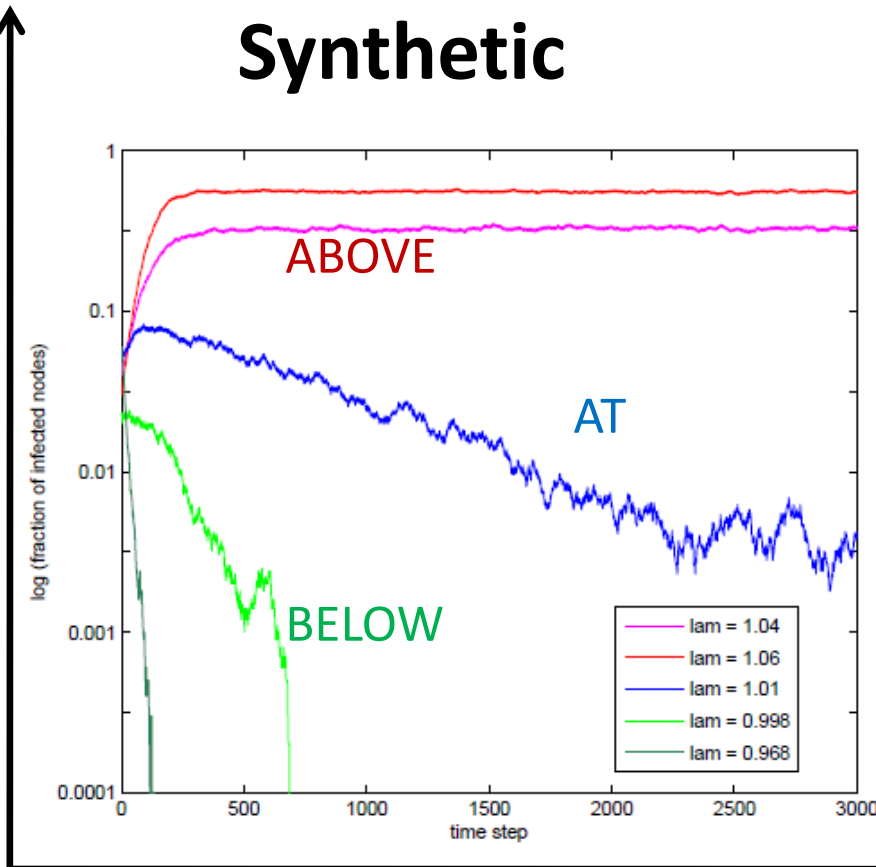
$$\mathbf{S}_i = (1 - \delta)\mathbf{I} + \beta\mathbf{A}_i$$

Infection-profile

$\log(\text{fraction infected})$

Synthetic

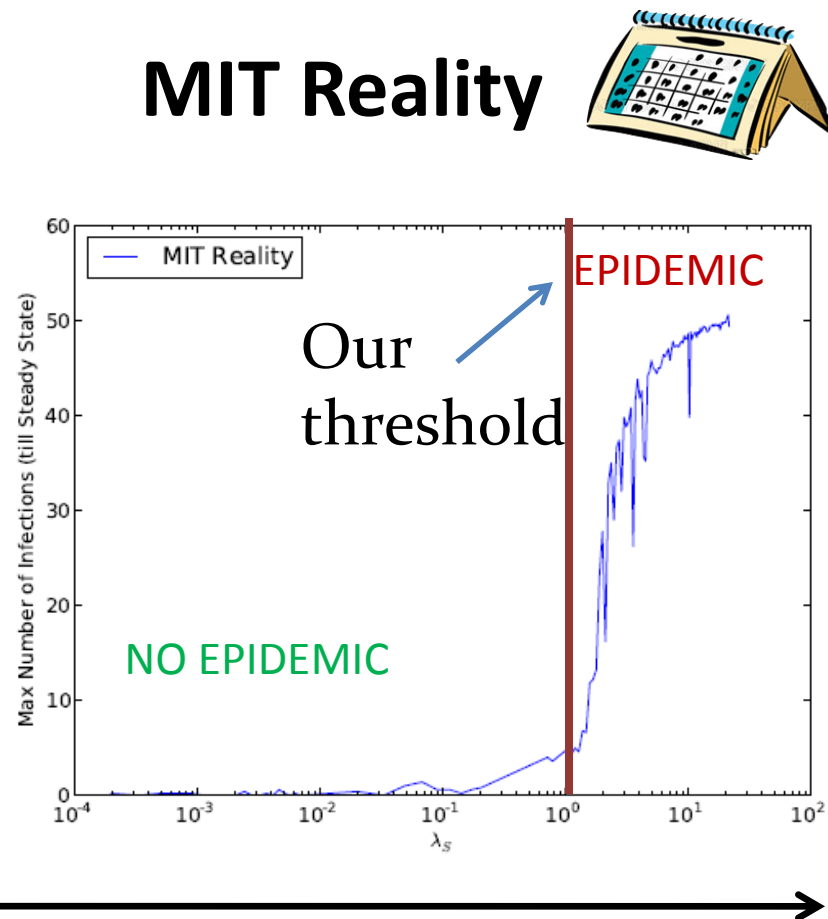
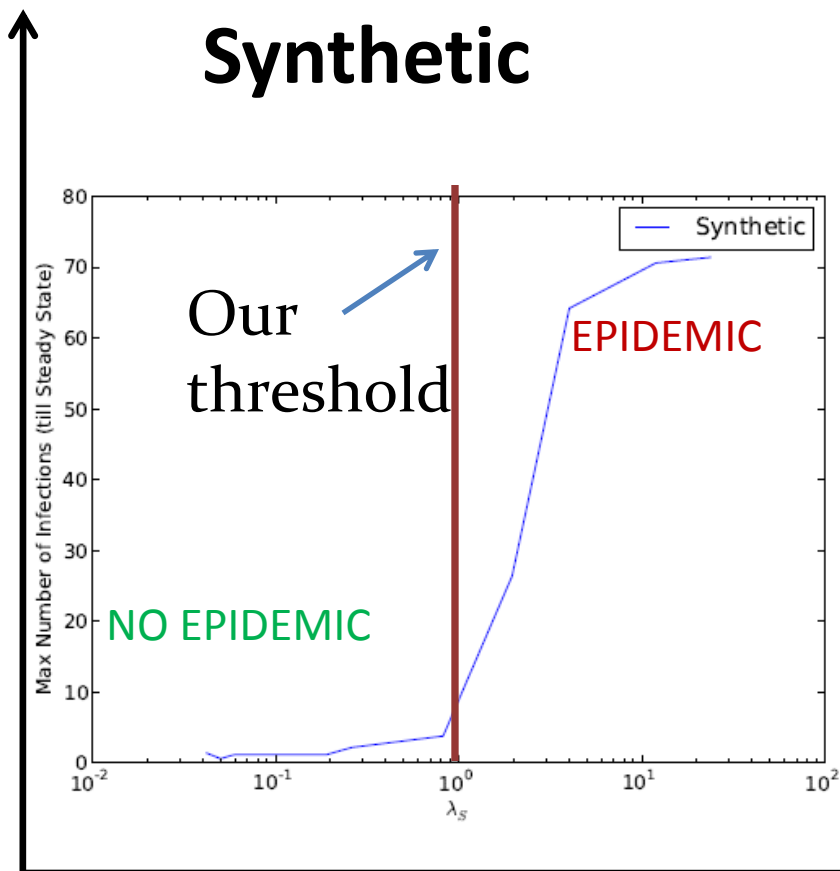
MIT Reality



Time

Footprint (#
infected @
“steady state”)

“Take-off” plots



$$\lambda_{\prod_i S_i} \text{ (log scale)}$$

Part 1: Theory

- Q1: What is the epidemic threshold?
- **Q2: What happens when viruses compete?**
 - Mutually-exclusive viruses
 - Interacting viruses

Competing Contagions



iPhone v Android



Blu-ray v HD-DVD



Attack

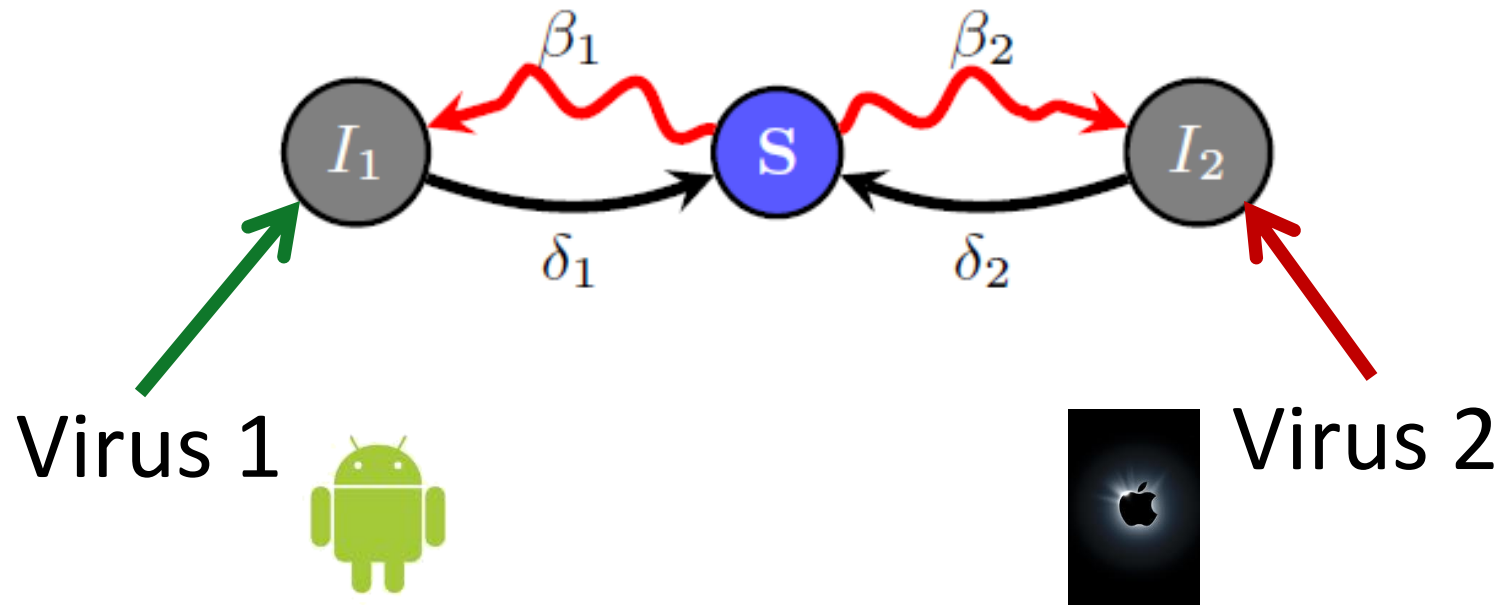
v



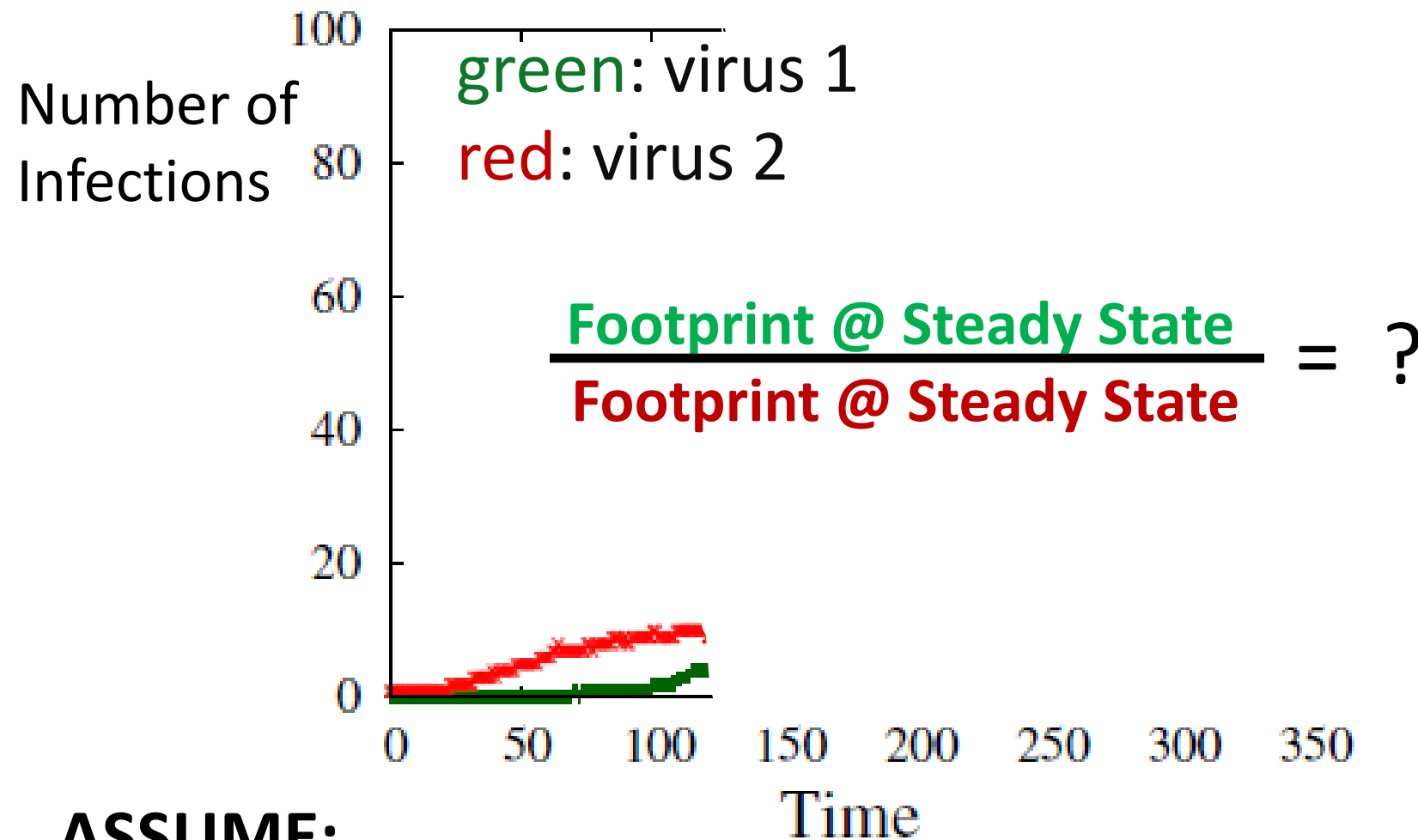
Retreat

A simple model

- Modified flu-like
- Mutual Immunity (“pick one of the two”)
- Susceptible-Infected1-Infected2-Susceptible



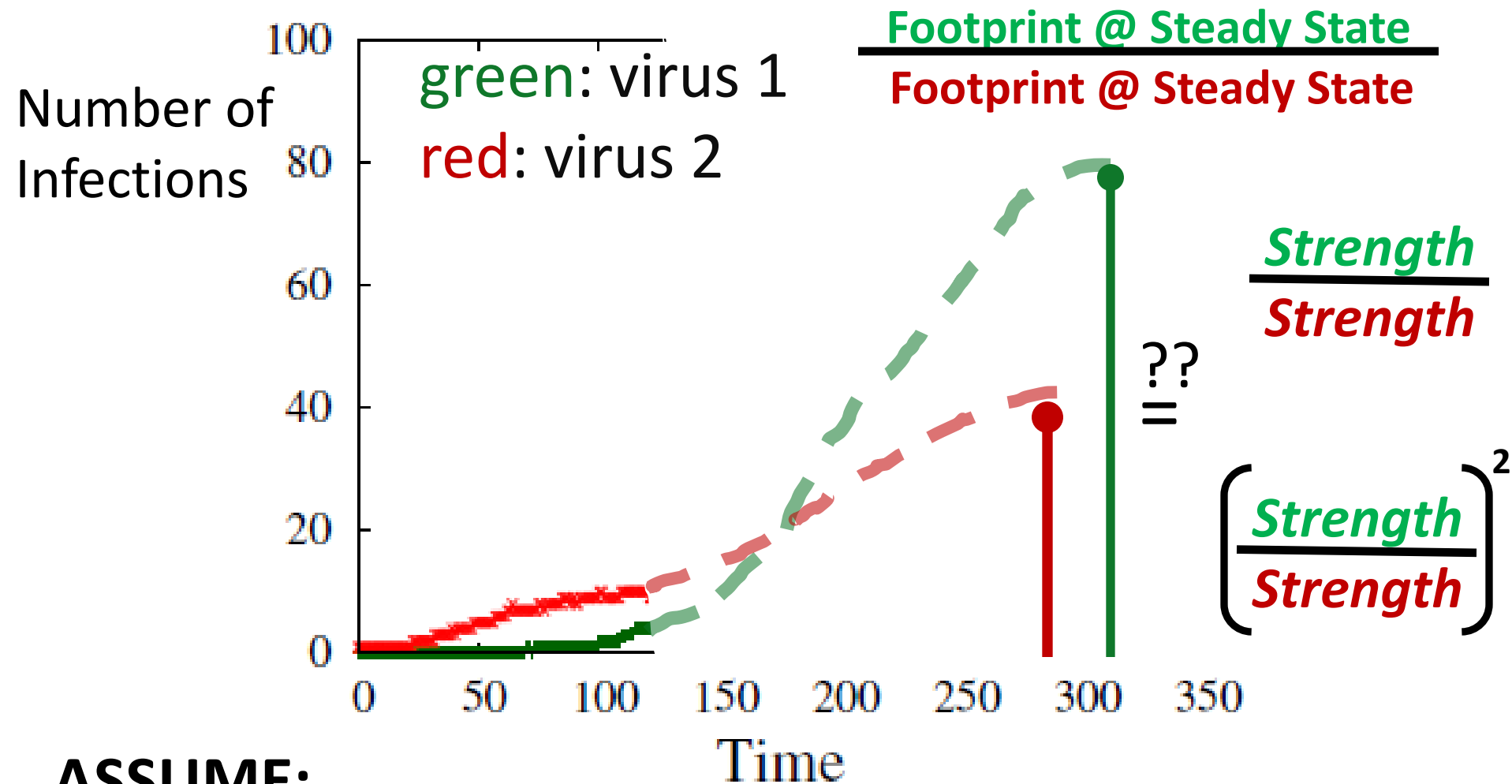
Question: What happens in the end?



ASSUME:

Virus 1 is stronger than Virus 2

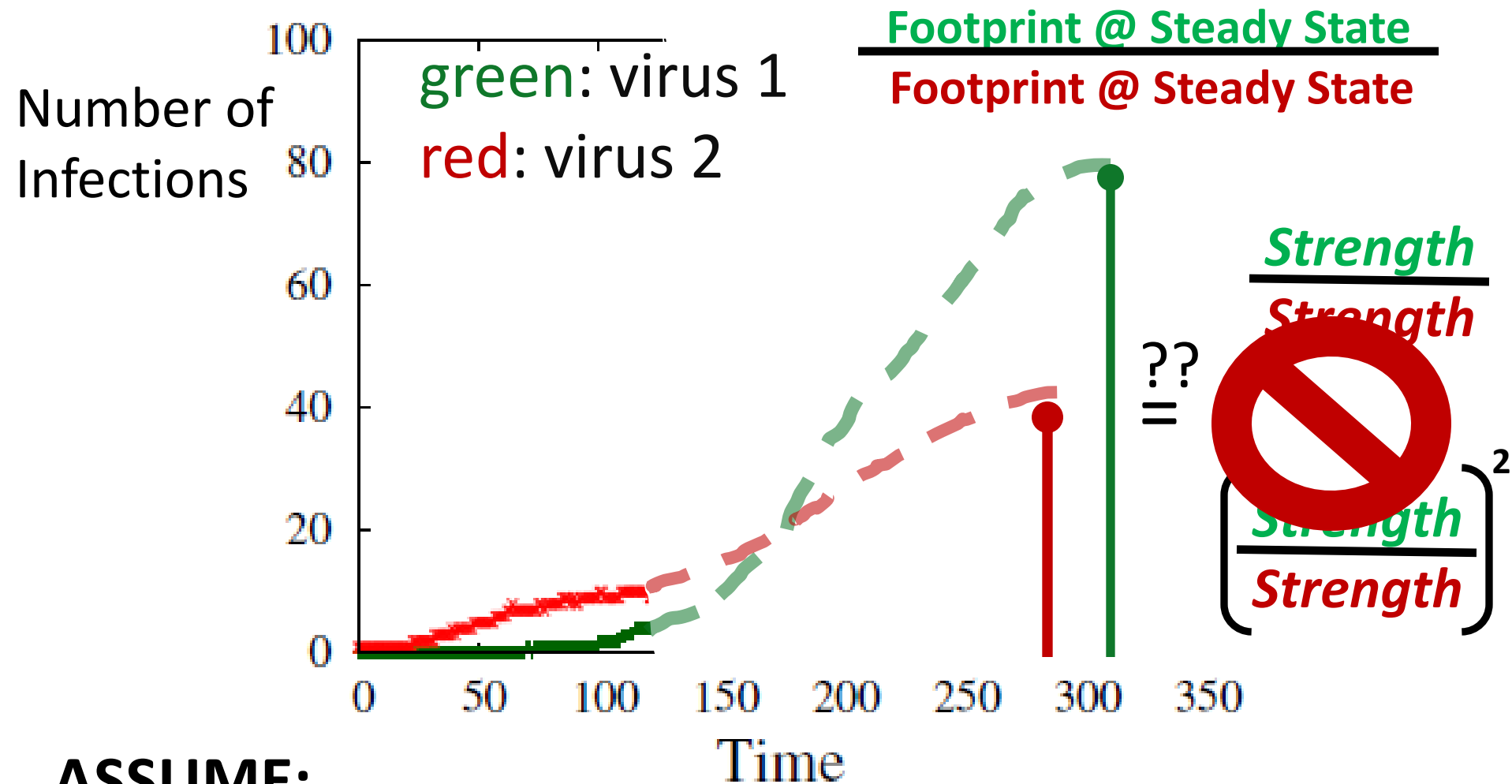
Question: What happens in the end?



ASSUME:

Virus 1 is stronger than Virus 2

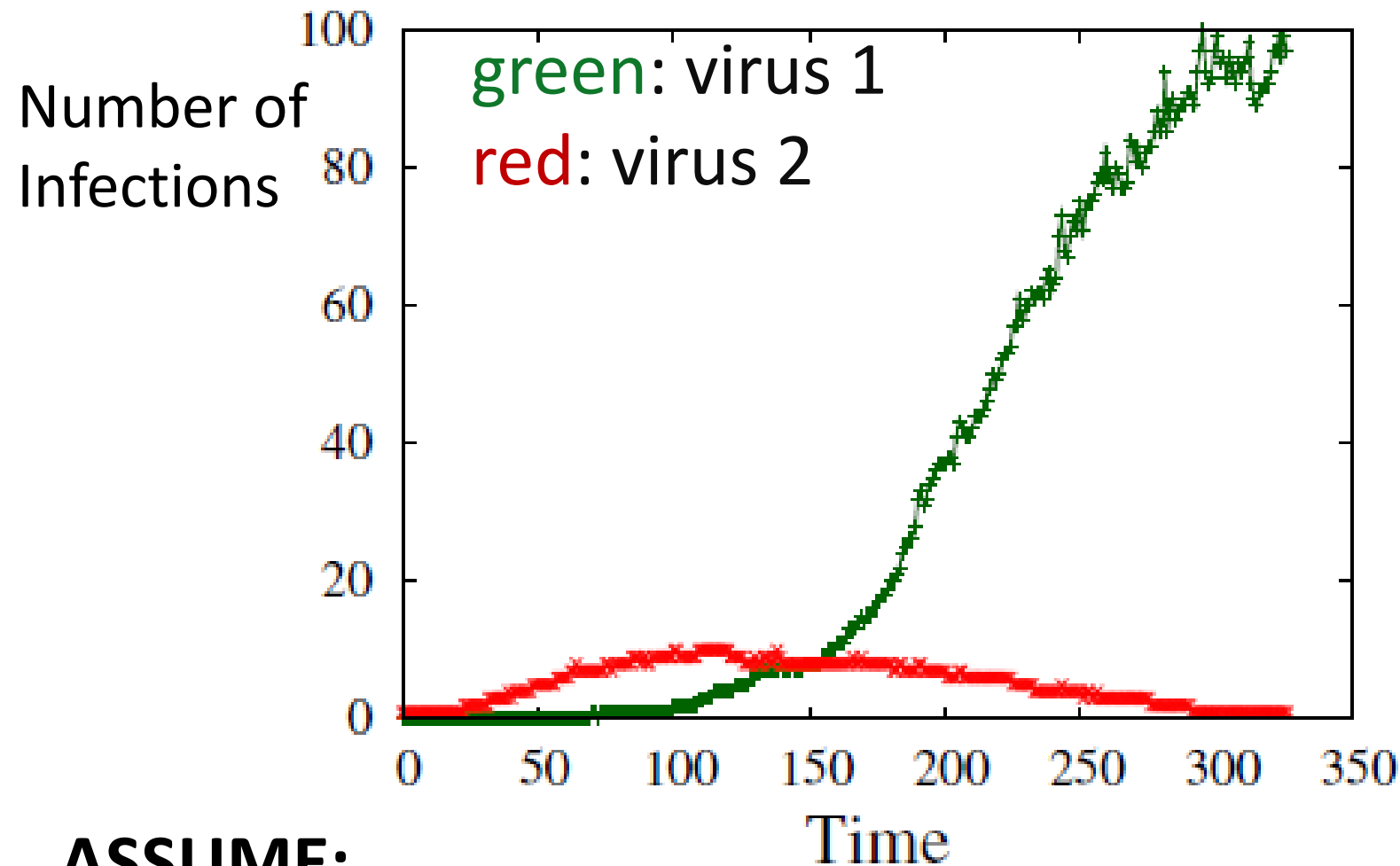
Question: What happens in the end?



ASSUME:

Virus 1 is stronger than Virus 2

Answer: Winner-Takes-All



ASSUME:

Virus 1 is stronger than Virus 2

Our Result: Winner-Takes-All

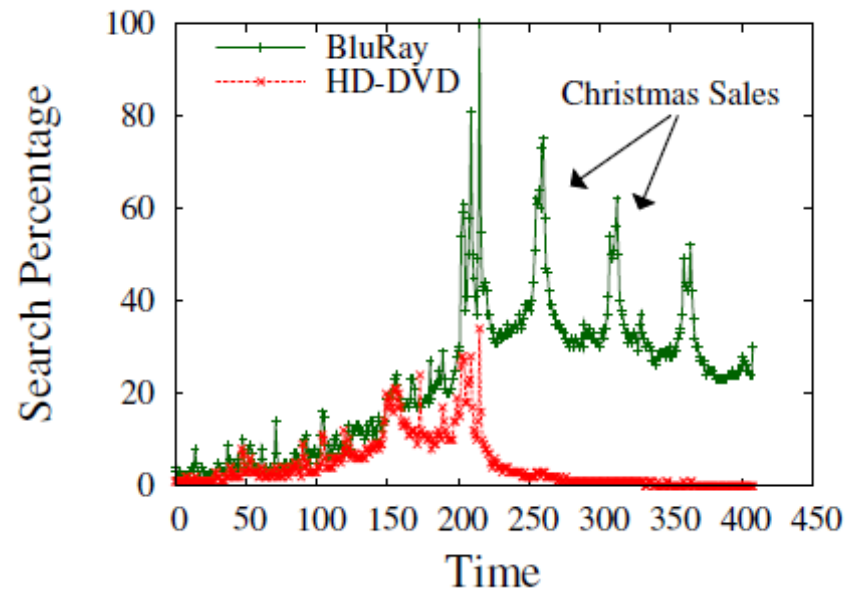
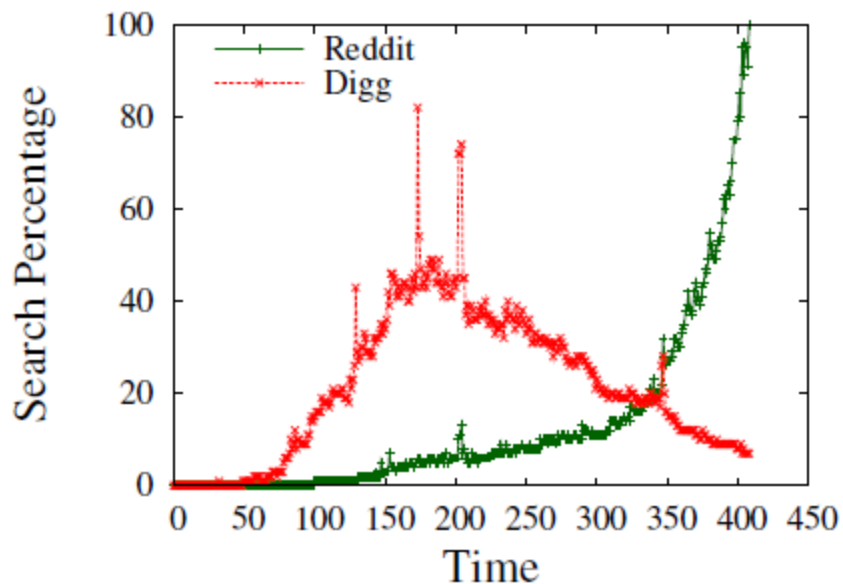
Given our model, and *any graph*, the weaker virus always **dies-out completely**

Details

1. The stronger survives only if it is above threshold
2. Virus 1 is stronger than Virus 2, if:
$$\text{strength}(\text{Virus 1}) > \text{strength}(\text{Virus 2})$$
3. $\text{Strength}(\text{Virus}) = \lambda \beta / \delta \rightarrow$ same as before!

Real Examples

[Google Search Trends data]



Reddit v Digg



Blu-Ray v HD-DVD

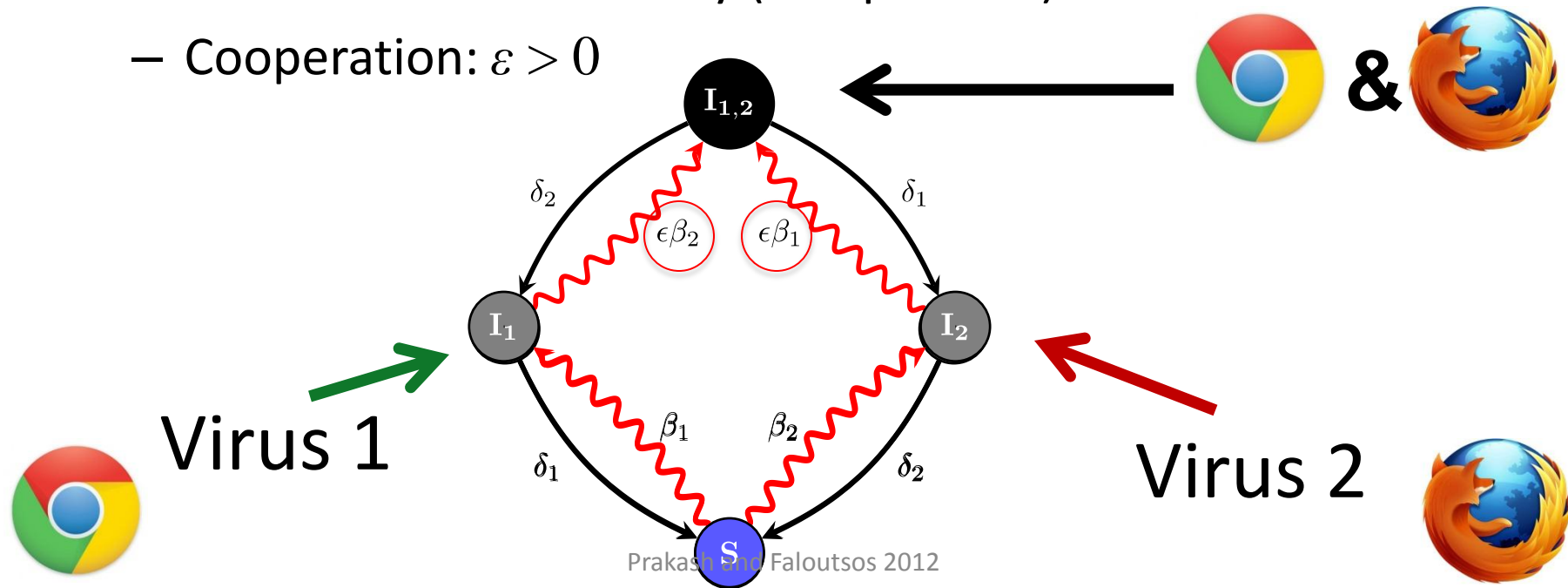


Part 1: Theory

- Q1: What is the epidemic threshold?
- **Q2: What happens when viruses compete?**
 - Mutually-exclusive viruses
 - **Interacting viruses**

A simple model: $SI_{1|2}S$

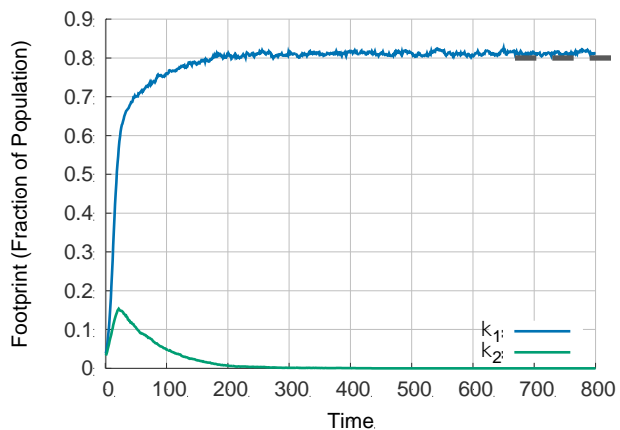
- Modified flu-like (SIS)
- Susceptible-Infected_{1 or 2}-Susceptible
- Interaction Factor ϵ
 - Full Mutual Immunity: $\epsilon = 0$
 - Partial Mutual Immunity (competition): $\epsilon < 0$
 - Cooperation: $\epsilon > 0$



Question: What happens in the end?

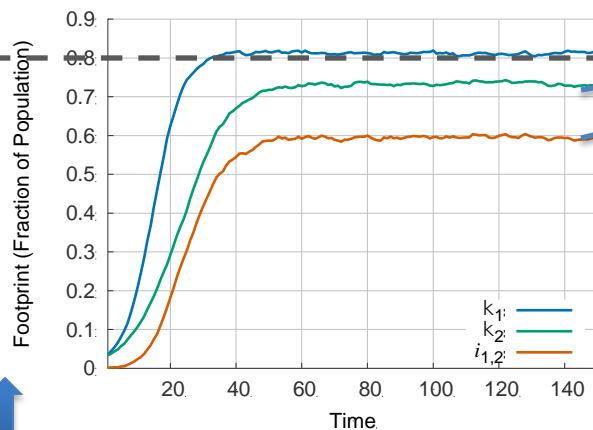
$\epsilon = 0$

Winner takes all



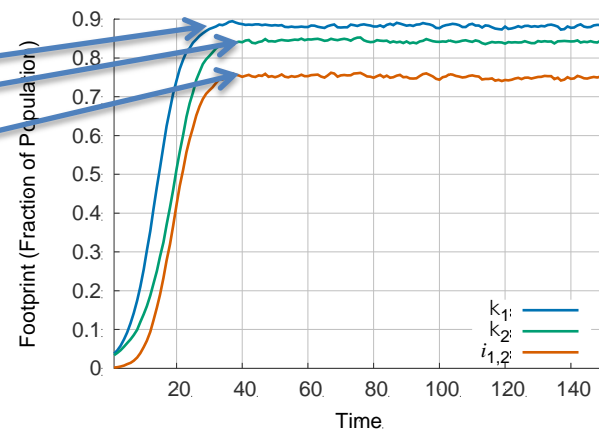
$\epsilon = 1$

Co-exist independently



$\epsilon = 2$

Viruses cooperate



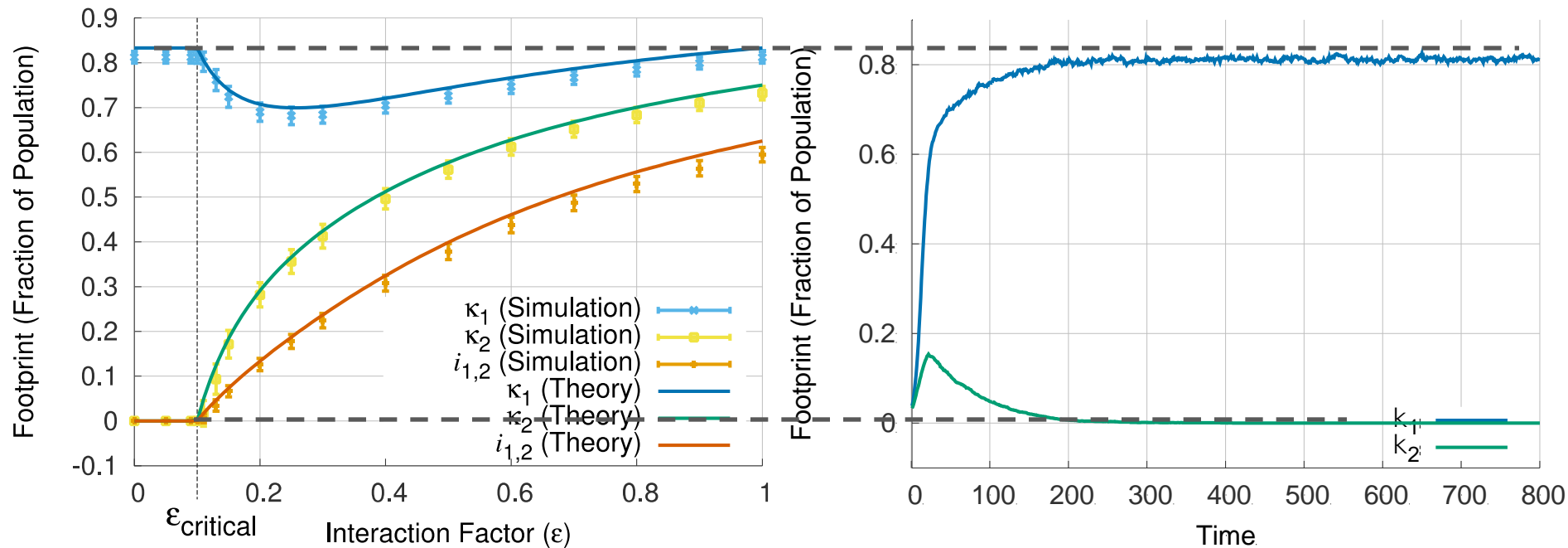
What about for $0 < \epsilon < 1$?
Is there a point at which both viruses can *co-exist*?

ASSUME:

Virus 1 is stronger than Virus 2

Answer: Yes!

There is a phase transition

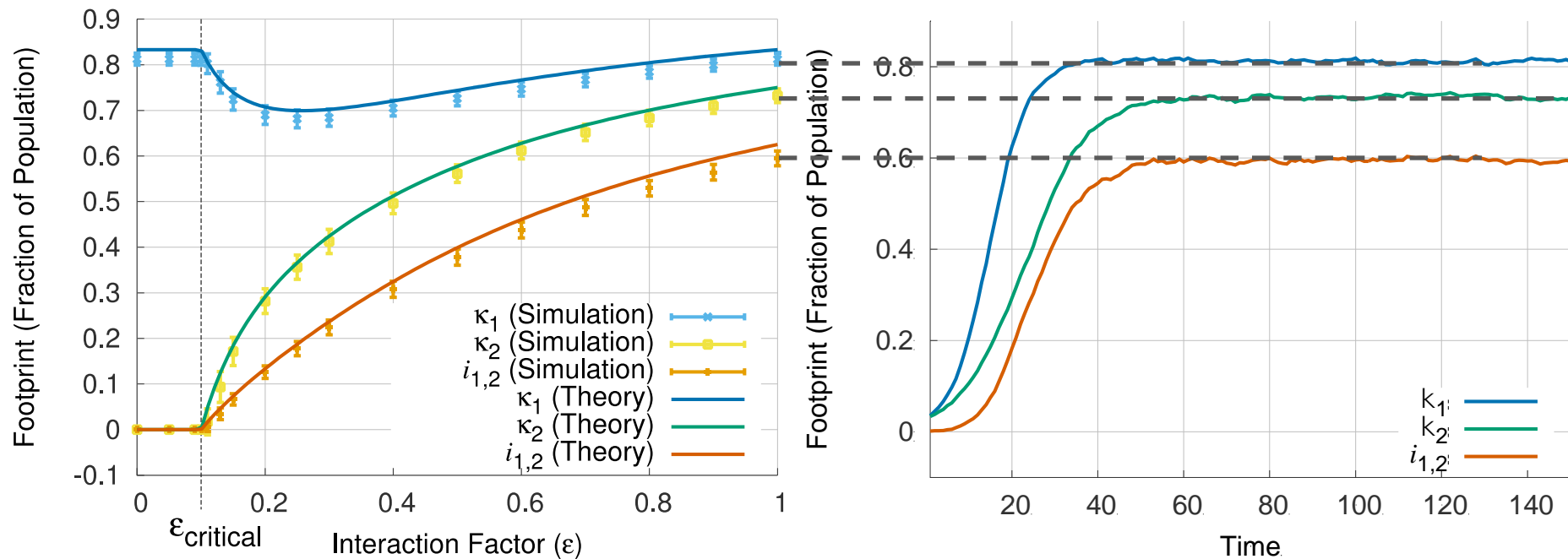


ASSUME:

Virus 1 is stronger than Virus 2

Answer: Yes!

There is a phase transition

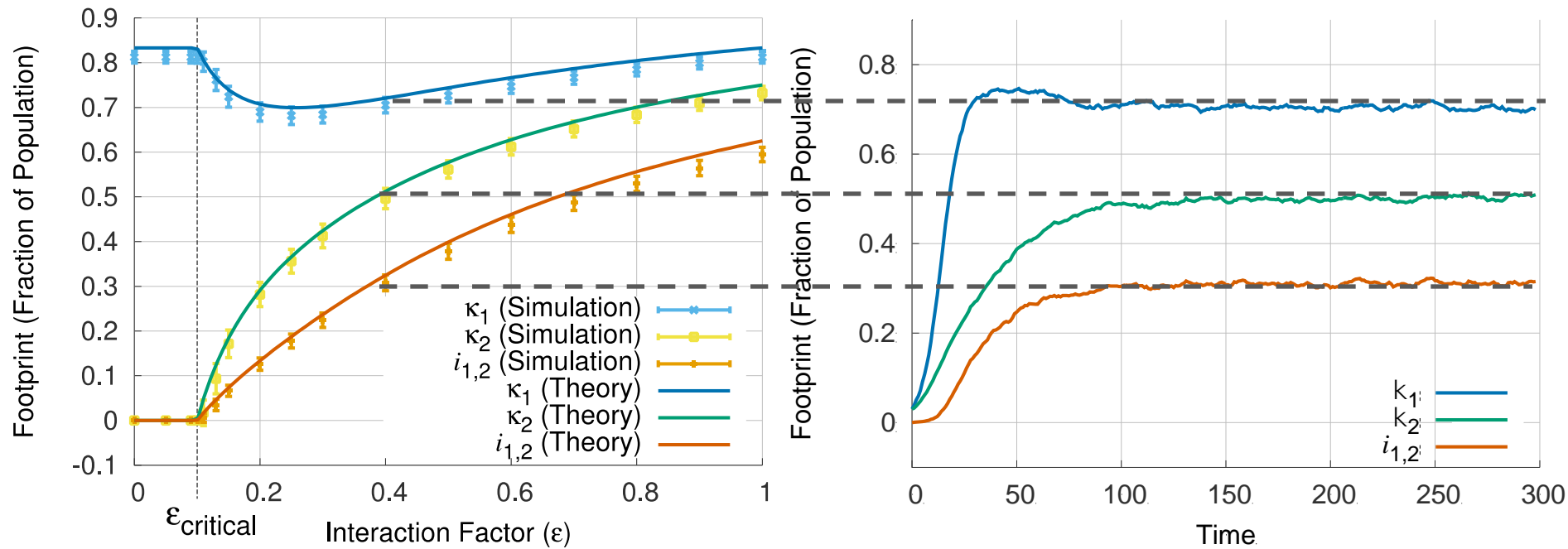


ASSUME:

Virus 1 is stronger than Virus 2

Answer: Yes!

There is a phase transition



ASSUME:

Virus 1 is stronger than Virus 2

Our Result: Viruses can Co-exist

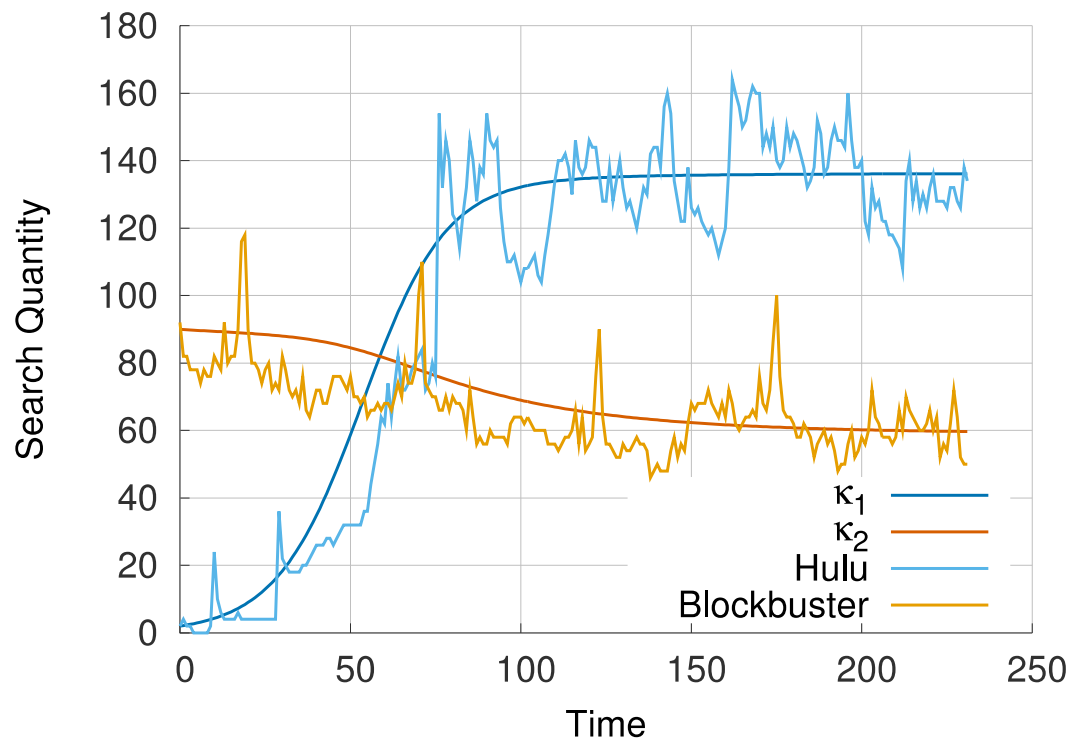
Given our model and a fully connected graph, there exists an $\varepsilon_{\text{critical}}$ such that for $\varepsilon \geq \varepsilon_{\text{critical}}$, there is a fixed point where both viruses survive.

Details

1. The stronger survives only if it is above threshold
2. Virus 1 is stronger than Virus 2, if:
$$\text{strength}(\text{Virus 1}) > \text{strength}(\text{Virus 2})$$
3. $\text{Strength}(\text{Virus}) \sigma = N \beta / \delta$

Real Examples

[Google Search Trends data]



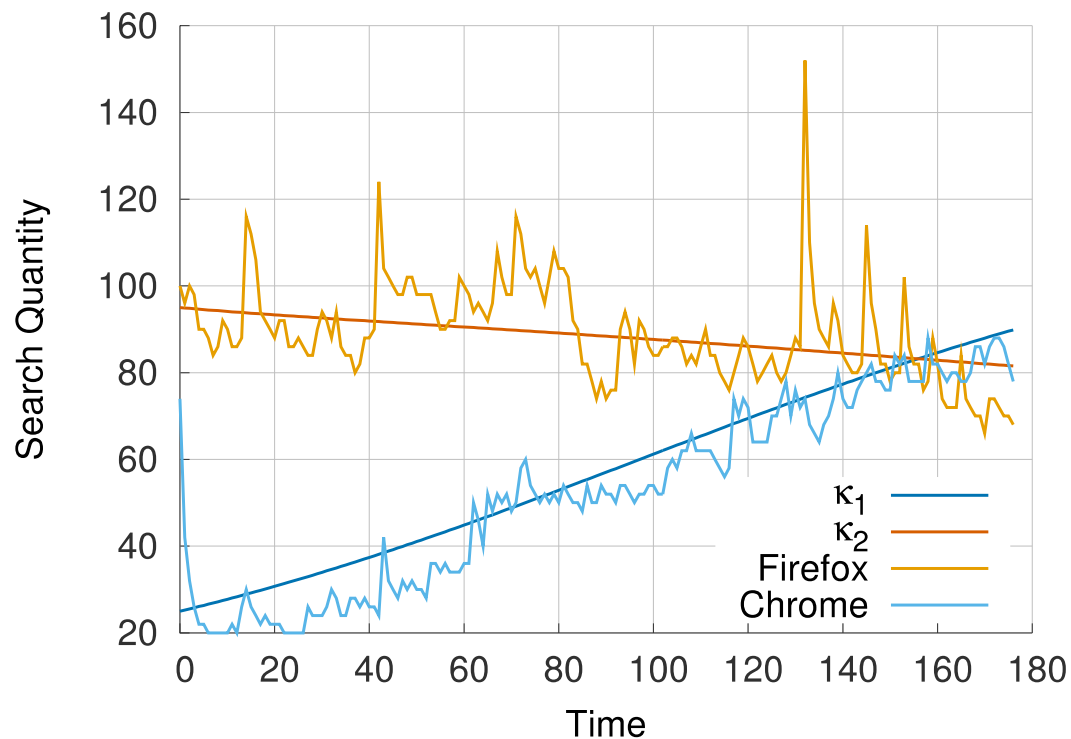
Hulu v Blockbuster

hulu



Real Examples

[Google Search Trends data]



Chrome v Firefox



Outline

- Motivation
- Part 1: Understanding Epidemics (Theory)
- **Part 2: Policy and Action (Algorithms)**
- Part 3: Learning Models (Empirical Studies)
- Conclusion

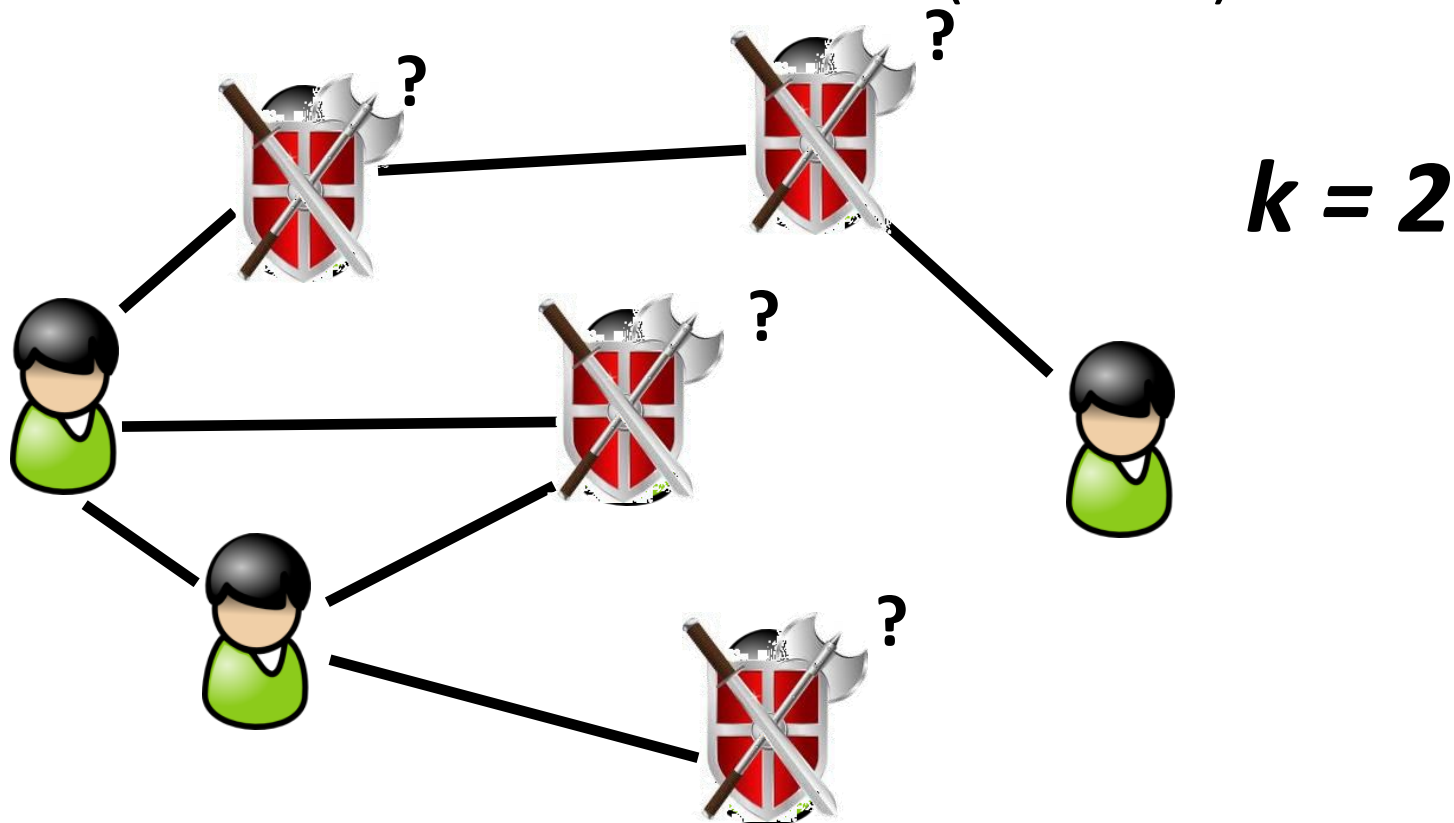
Part 2: Algorithms

- **Q3: Whom to immunize?**
- Q4: How to detect outbreaks?
- Q5: Who are the culprits?

Full Static Immunization

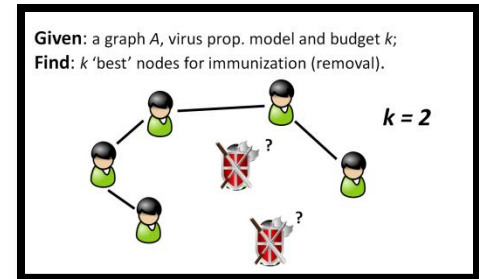
Given: a graph A , virus prop. model and budget k ;

Find: k 'best' nodes for immunization (removal).



Part 2: Algorithms

- **Q3: Whom to immunize?**
 - Full Immunization (Static Graphs)
 - Full Immunization (Dynamic Graphs)
 - Fractional Immunization
- Q4: How to detect outbreaks?
- Q5: Who are the culprits?



Challenges

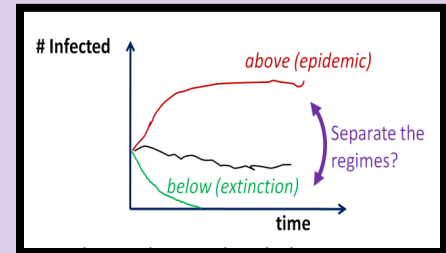
- Given a graph A , budget k ,

Q1 (Metric) How to measure the ‘shield-value’ for a set of nodes (S)?

Q2 (Algorithm) How to find a set of k nodes with highest ‘shield-value’?

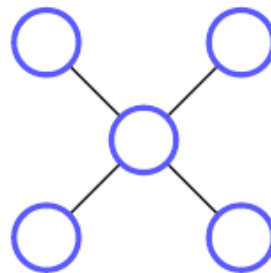
Proposed vulnerability measure λ

λ is the epidemic threshold



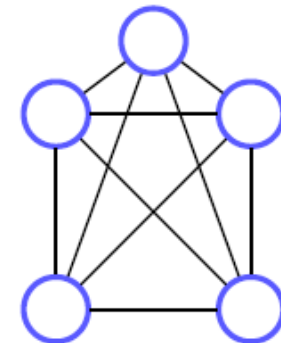
“Safe”

(a) Chain ($\lambda = 1.73$)



“Vulnerable”

(b) Star ($\lambda = 2$)



“Deadly”

(c) Clique ($\lambda = 4$)



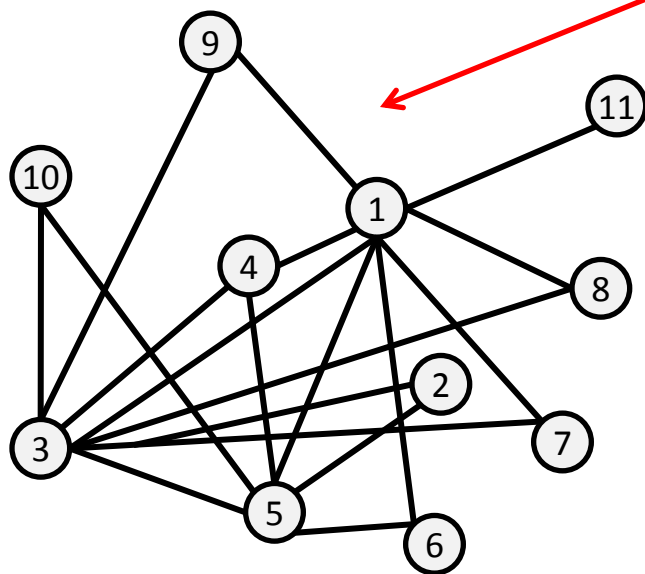
Increasing λ

Increasing vulnerability

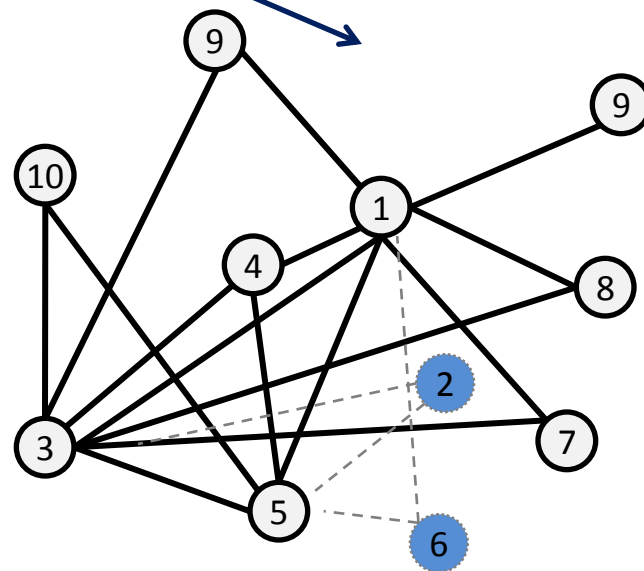
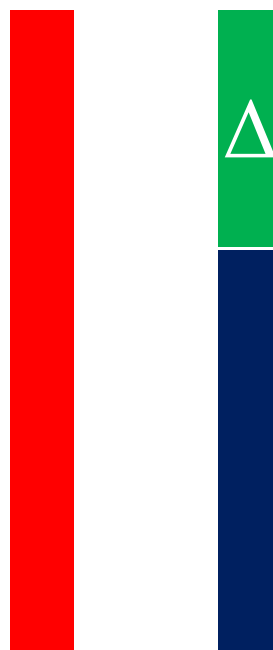
A1: "Eigen-Drop": an ideal shield value

Eigen-Drop(S)

$$\Delta \lambda = \lambda - \lambda_s$$



Original Graph



Without {2, 6}

(Q2) - Direct Algorithm too expensive!

- Immunize k nodes which maximize $\Delta \lambda$

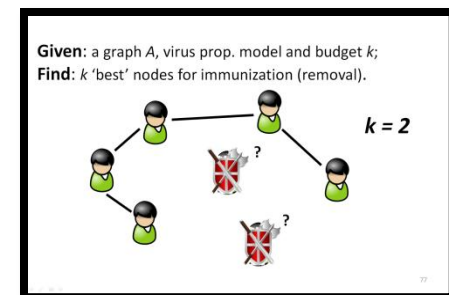
$$S = \operatorname{argmax} \Delta \lambda$$

- Combinatorial!

- Complexity: $O\left(\binom{n}{k} \cdot m\right)$

– Example:

- 1,000 nodes, with 10,000 edges
- It takes 0.01 seconds to compute λ
- It takes **2,615 years** to find 5-best nodes!



A2: Our Solution

- Part 1: Shield Value
 - Carefully approximate Eigen-drop ($\Delta \lambda$)
 - Matrix **perturbation** theory
- Part 2: Algorithm
 - Greedily pick best node at each step
 - Near-optimal due to **submodularity**
- *NetShield* (linear complexity)
 - **$O(nk^2+m)$** $n = \#$ nodes; $m = \#$ edges

Our Solution: Part 1

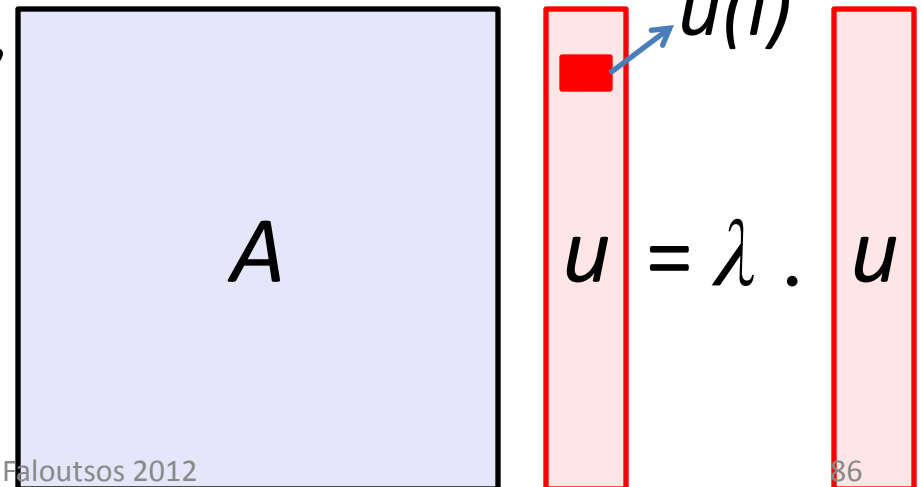
- Approximate Eigen-drop ($\Delta \lambda$)

- $$\Delta \lambda \approx \widehat{S\hat{V}}(S) = \sum_{i \in S} 2\lambda \mathbf{u}(i)^2 - \sum_{i, j \in S} \mathbf{A}(i, j) \mathbf{u}(i) \mathbf{u}(j)$$

– Result using Matrix perturbation theory

– $u(i) ==$ ‘eigenscore’

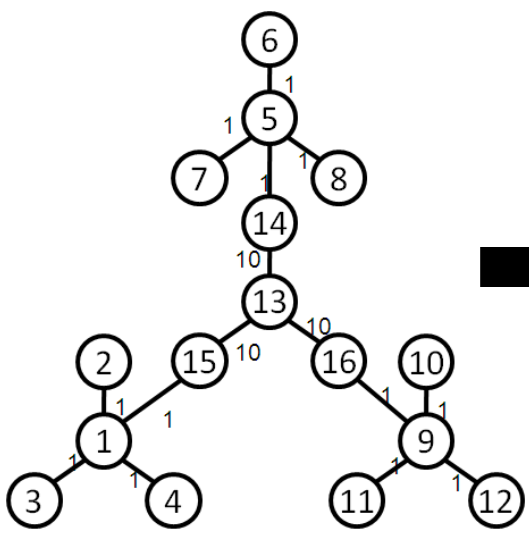
~~ pagerank(i)



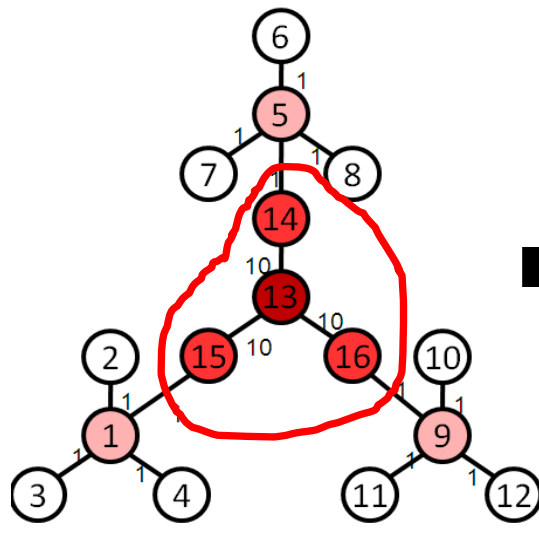
$$\sum_{i \in S} 2\lambda \mathbf{u}(i)^2 - \sum_{i, j \in S} \mathbf{A}(i, j) \mathbf{u}(i) \mathbf{u}(j)$$

P1: node importance

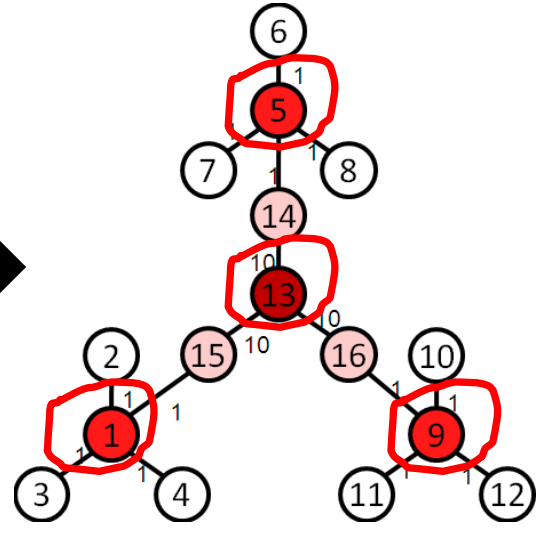
P2: set diversity



Original Graph



Select by P1



Select by P1+P2

Our Solution: Part 2: NetShield

- We prove that:

$\widehat{SV}(S)$ is sub-modular (& monotone non-decreasing)



Corollary: Greedy algorithm works

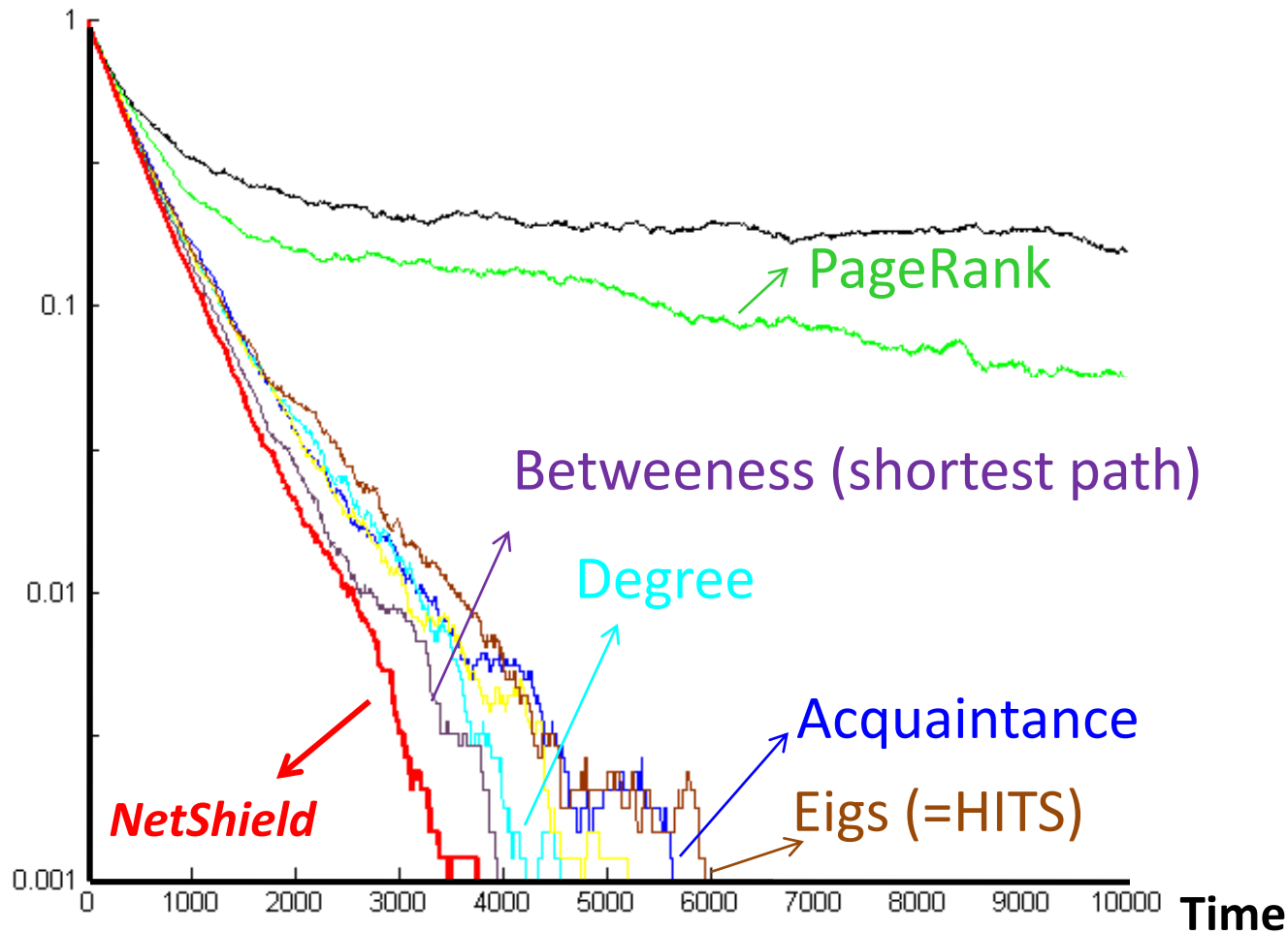
1. NetShield is near-optimal (w.r.t. $\max \widehat{SV}(S)$)
2. NetShield is $O(nk^2+m)$

- NetShield: Greedily add best node at each step

Footnote: near-optimal means $\widehat{SV}(S^{NetShield}) \geq (1-1/e) \widehat{SV}(S^{Opt})$

Experiment: Immunization quality

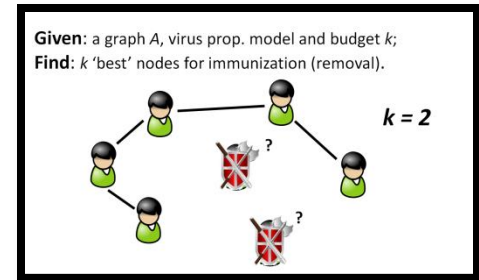
Log(fraction of infected nodes)



Lower is better

Part 2: Algorithms

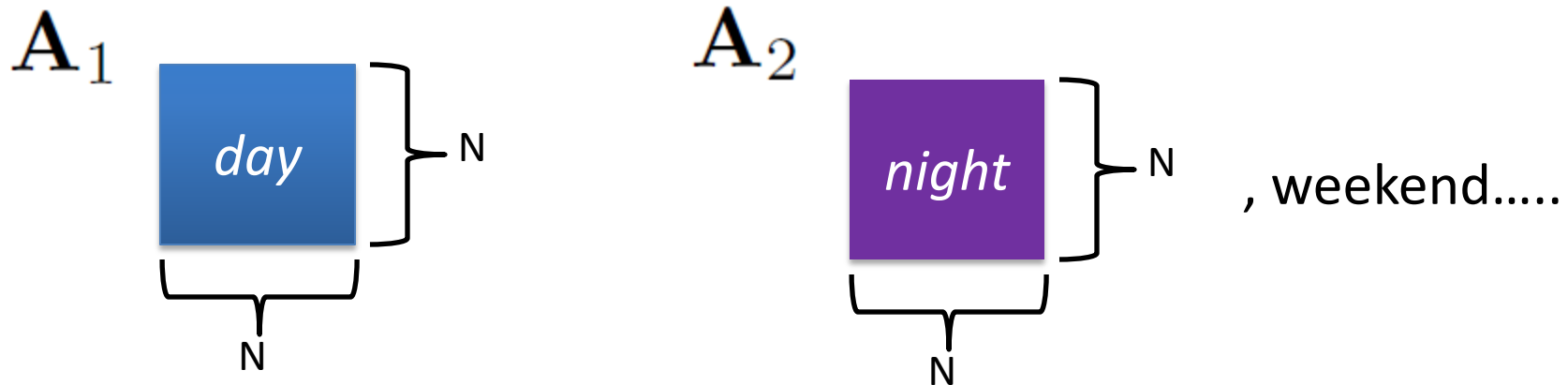
- **Q3: Whom to immunize?**
 - Full Immunization (Static Graphs)
 - Full Immunization (Dynamic Graphs)
 - Fractional Immunization
- Q4: How to detect outbreaks?
- Q5: Who are the culprits?



Full Dynamic Immunization

- *Given:*

Set of T arbitrary graphs $\{\mathbf{A}_1, \mathbf{A}_2 \dots, \mathbf{A}_T\}$



- *Find:*

k 'best' nodes to immunize (remove)

Full Dynamic Immunization

- Our solution

- Recall theorem

- Simple: reduce $\lambda \prod_i s_i$ ($=\lambda$)

*Matrix
Product*

A_1
day



A_2
night



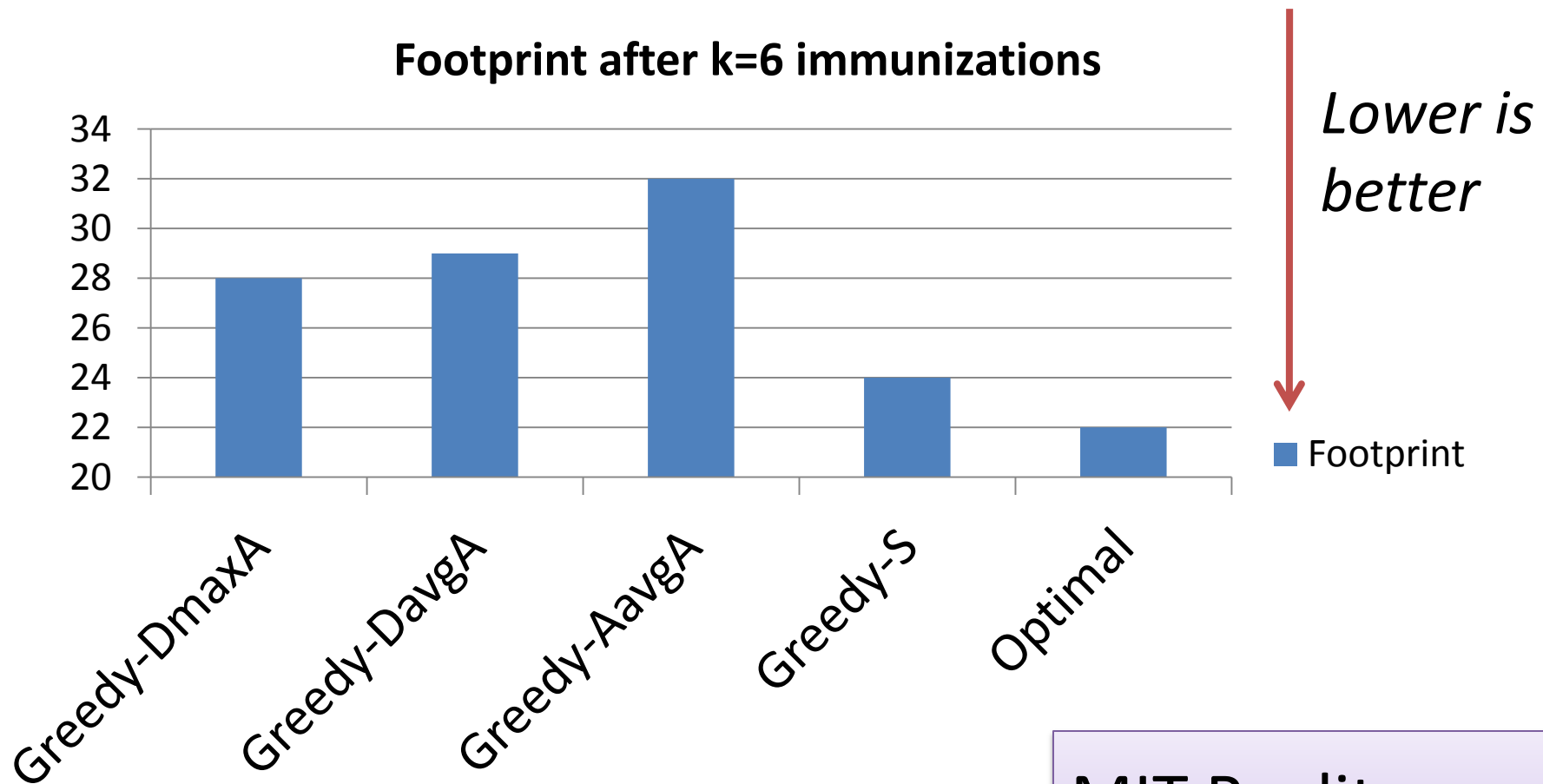
- Goal: max eigendrop $\Delta \lambda$

$$\Delta \lambda = \lambda_{\text{before}} - \lambda_{\text{after}}$$

- No competing policy for comparison

- We propose and evaluate many policies

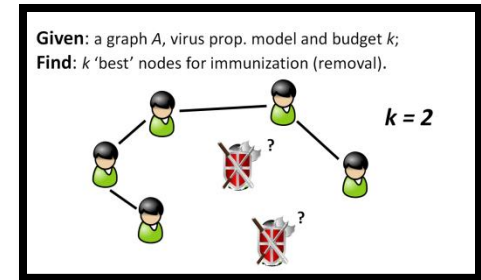
Performance of Policies



MIT Reality
Mining

Part 2: Algorithms

- **Q3: Whom to immunize?**
 - Full Immunization (Static Graphs)
 - Full Immunization (Dynamic Graphs)
 - Fractional Immunization
- Q4: How to detect outbreaks?
- Q5: Who are the culprits?



Fractional Immunization of Networks

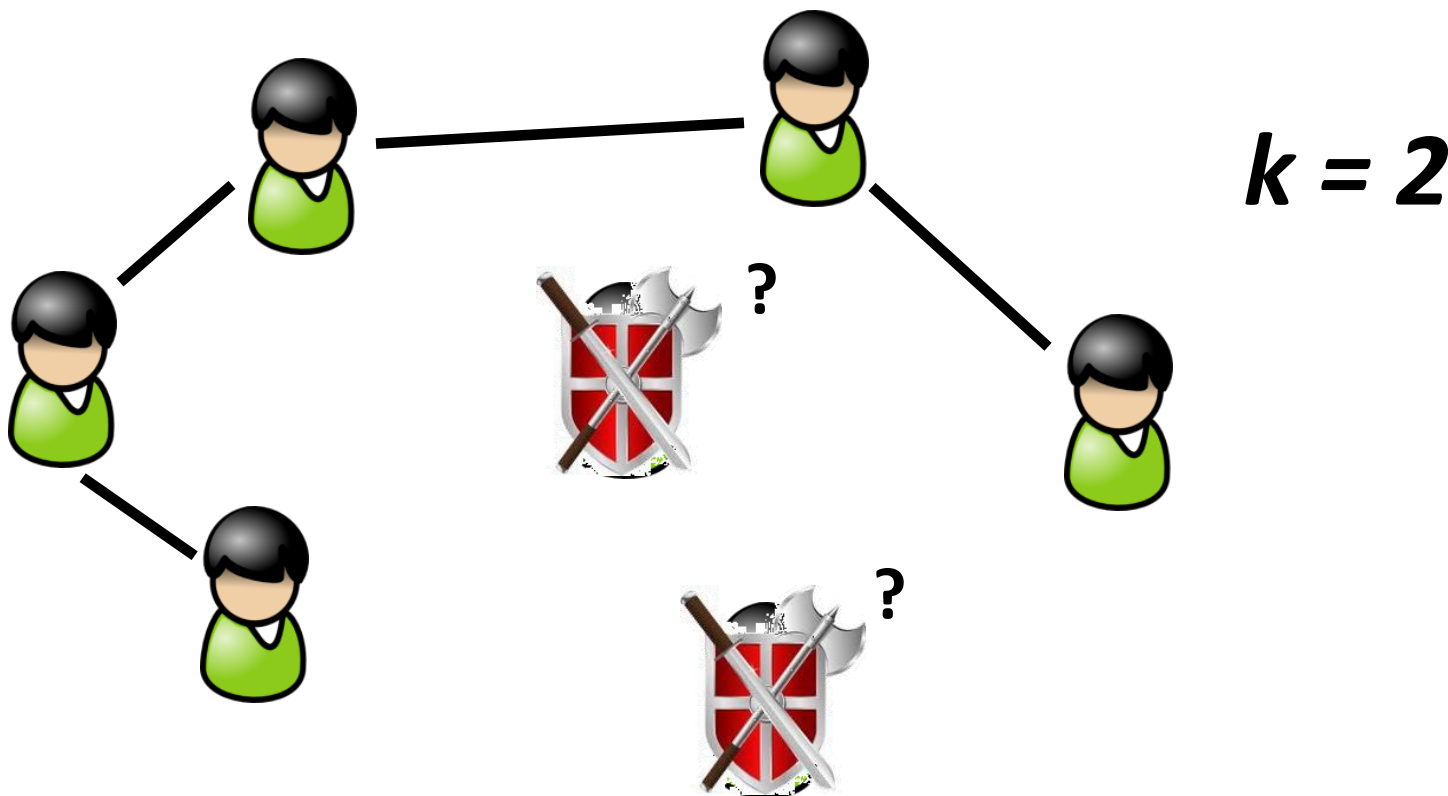
B. Aditya Prakash, Lada Adamic, Theodore Iwashyna (M.D.), Hanghang Tong, Christos Faloutsos

Under Submission

Previously: Full Static Immunization

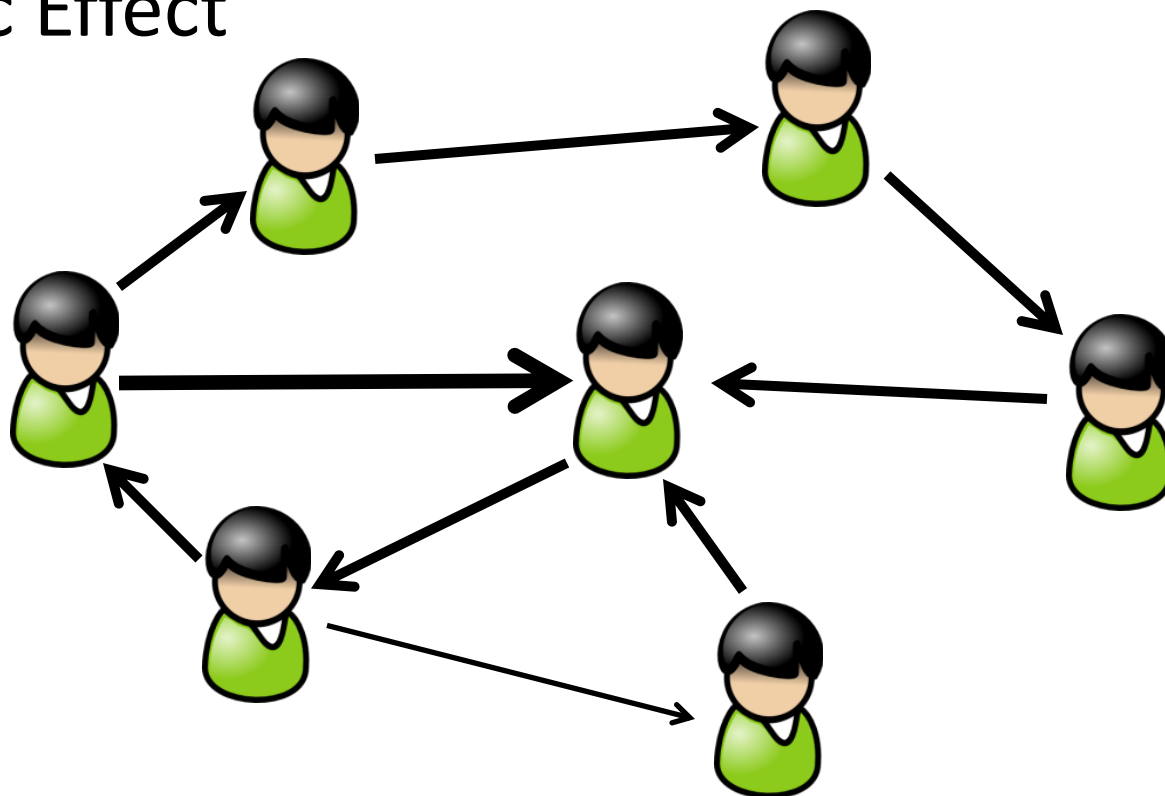
Given: a graph A , virus prop. model and budget k ;

Find: k 'best' nodes for immunization (removal).



Fractional Asymmetric Immunization

- Fractional Effect [$f(x) = 0.5^x$]
- Asymmetric Effect

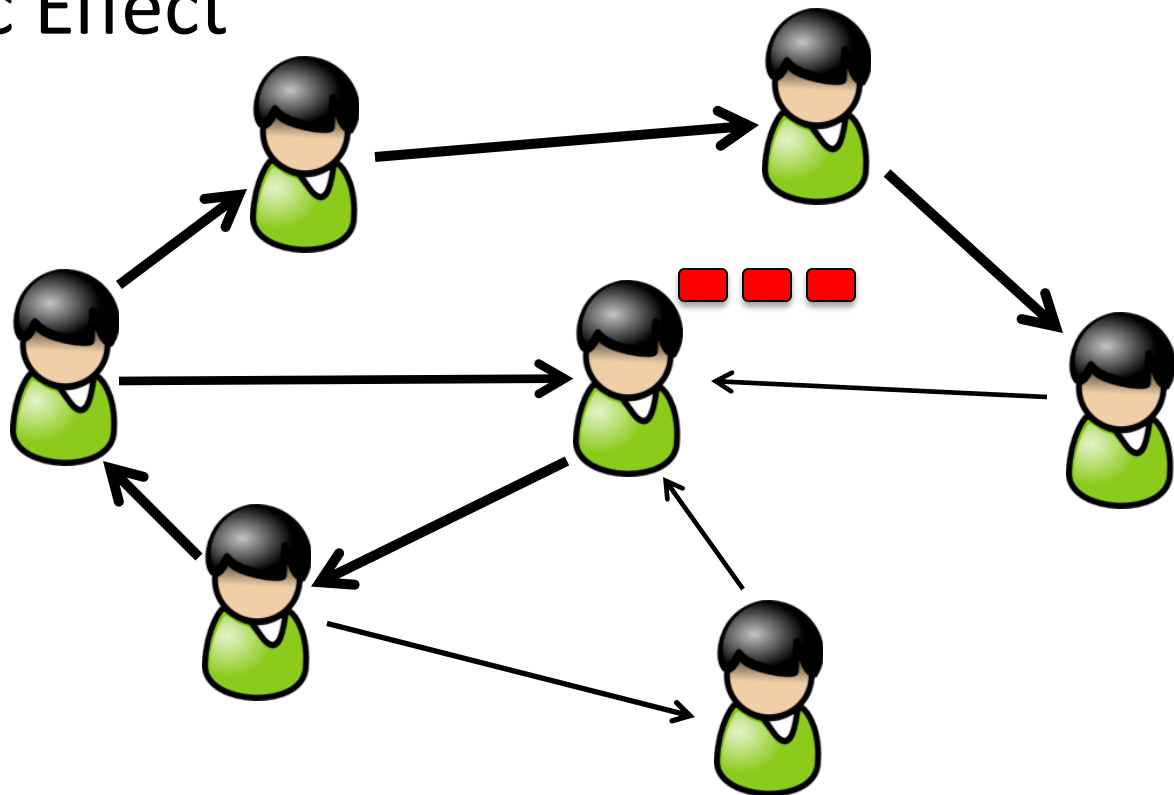


antidotes = 3



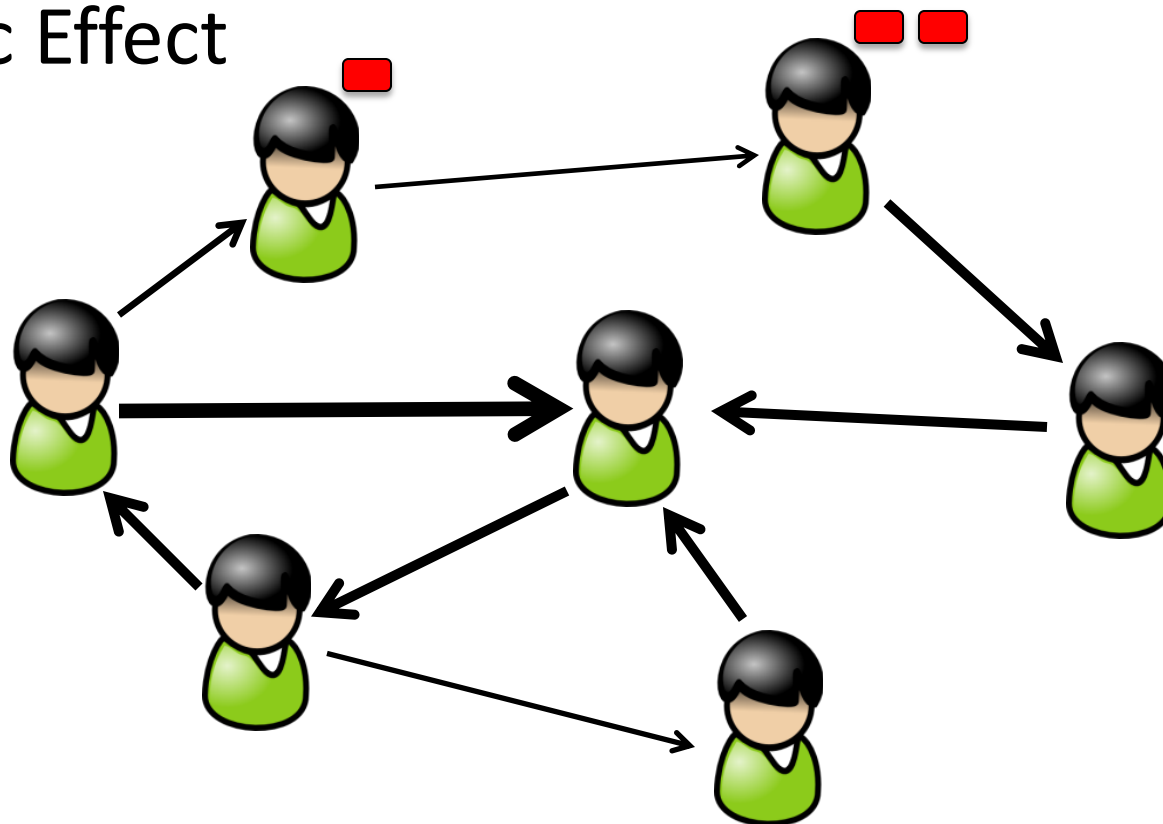
Now: Fractional Asymmetric Immunization

- Fractional Effect [$f(x) = 0.5^x$]
- Asymmetric Effect



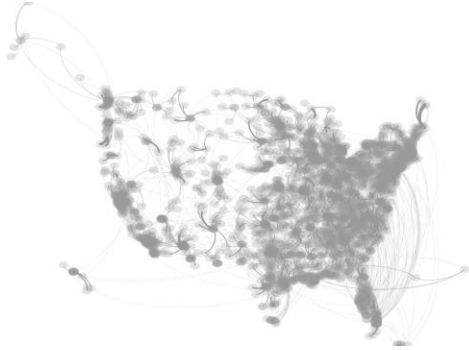
Fractional Asymmetric Immunization

- Fractional Effect [$f(x) = 0.5^x$]
- Asymmetric Effect

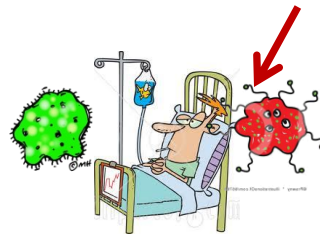


antidotes = 3

Fractional Asymmetric Immunization



Drug-resistant Bacteria
(like XDR-TB)



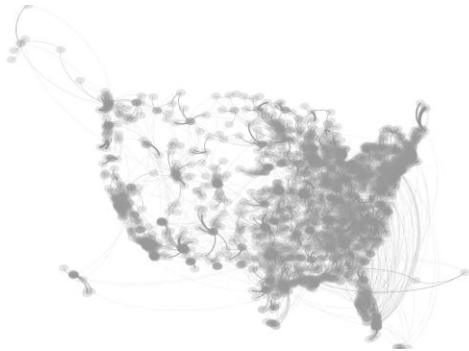
Hospital



Another Hospital



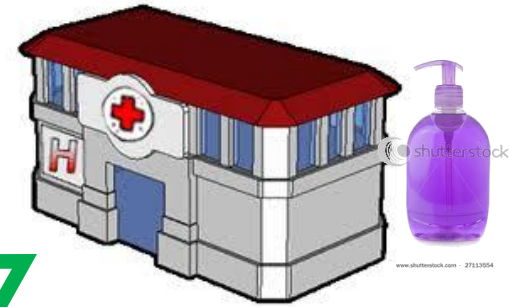
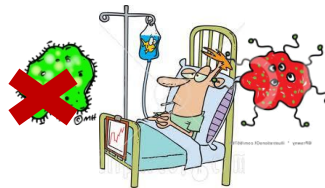
Fractional Asymmetric Immunization



$$\uparrow = f(\uparrow)$$



Hospital



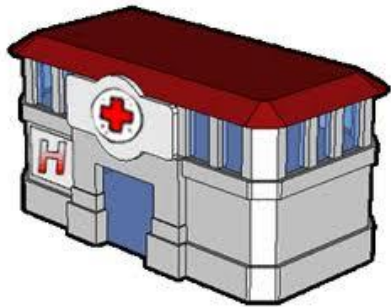
Another Hospital



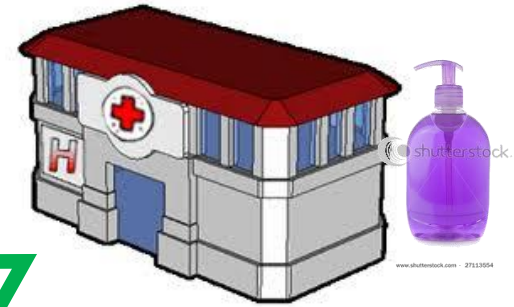
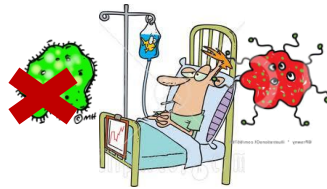
Fractional Asymmetric Immunization



Problem: Given k units of disinfectant, how to distribute them to maximize hospitals saved?



Hospital



Another Hospital

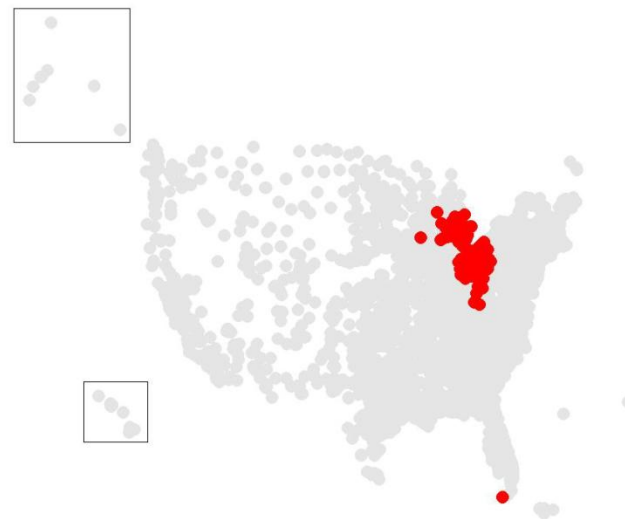
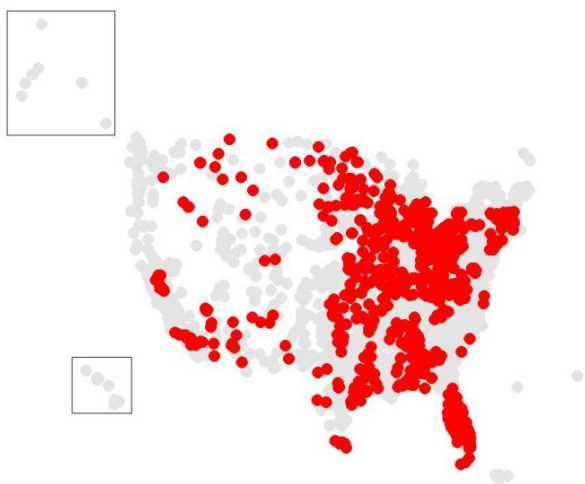


Our Algorithm “SMART-ALLOC”

**~6x
fewer!**

[US-MEDICARE NETWORK 2005]

- Each circle is a hospital, ~3000 hospitals
- More than 30,000 patients transferred



CURRENT PRACTICE

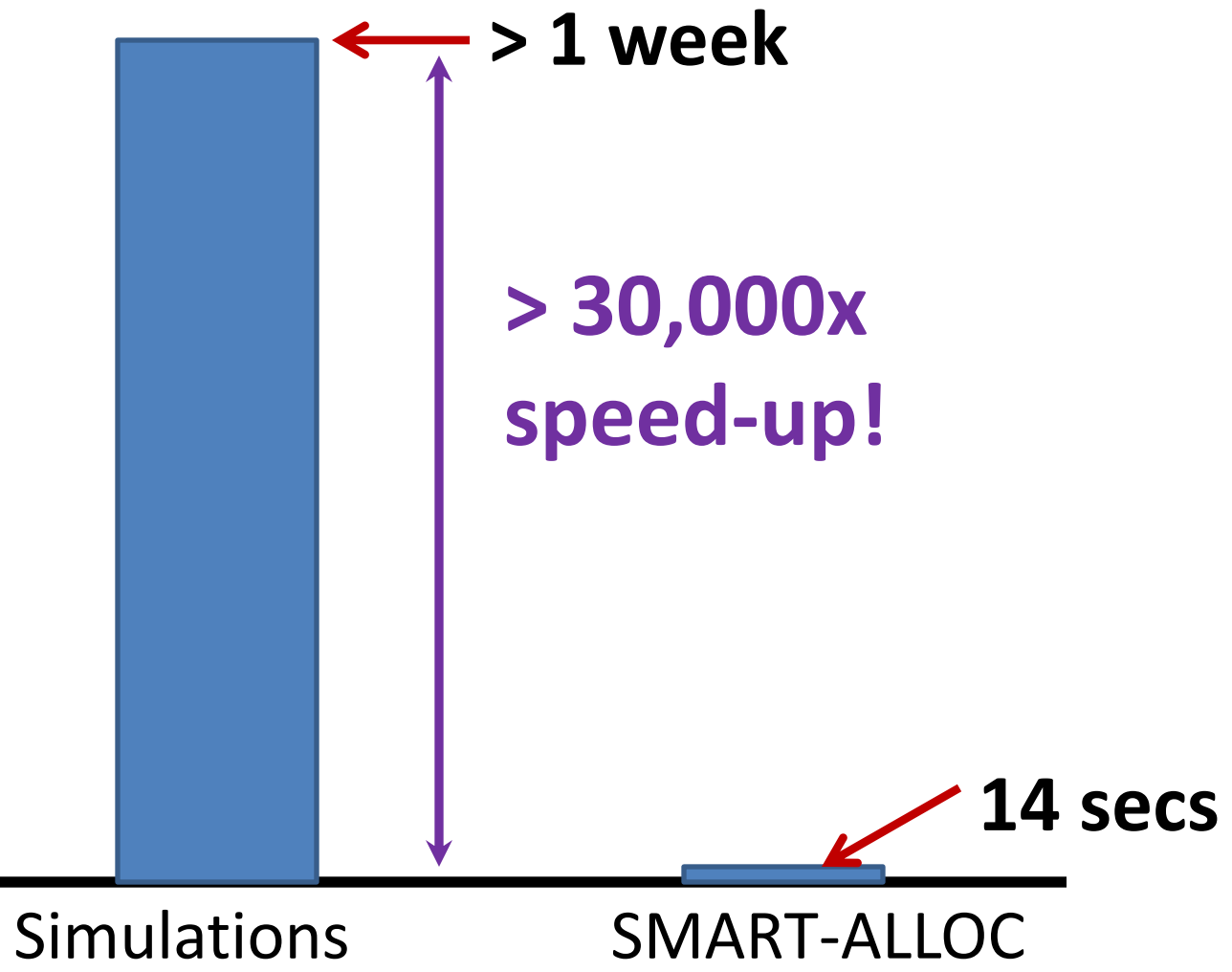
SMART-ALLOC

Running Time

Wall-Clock
Time



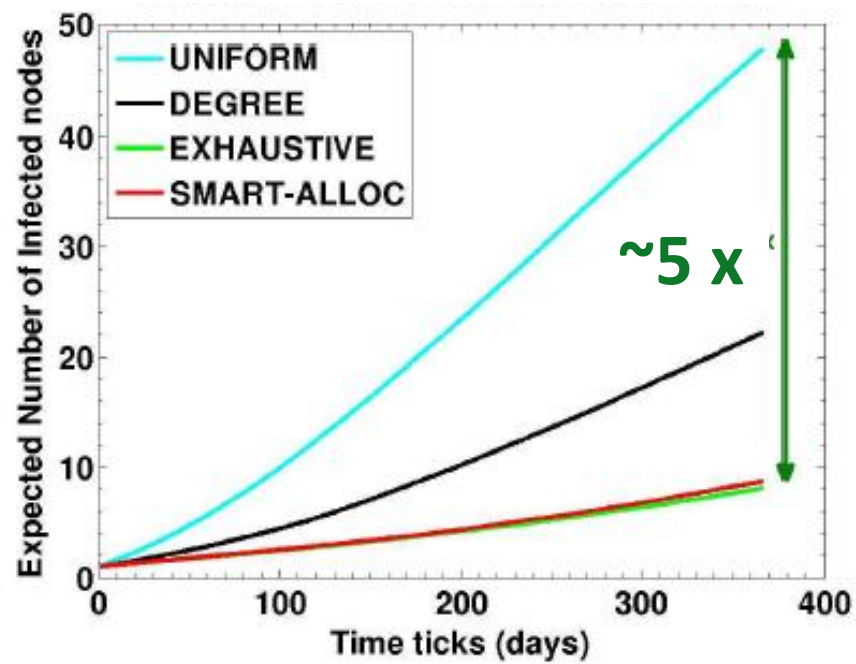
Lower
is
better



Lower
is
better

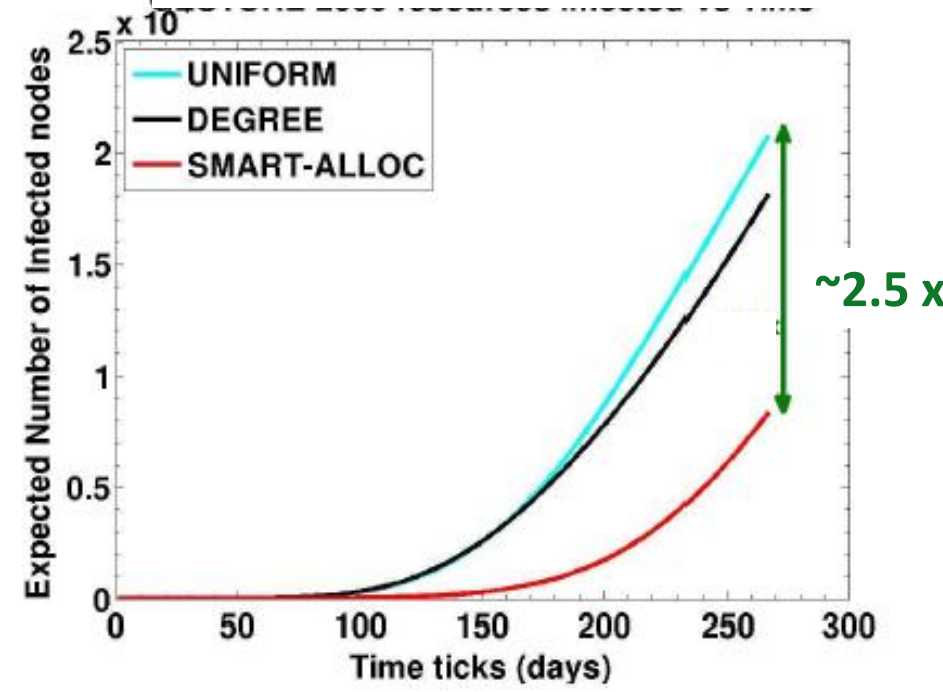
Experiments

PENN-NETWORK



K = 200

SECOND-LIFE



K = 2000

Part 2: Algorithms

- Q3: Whom to immunize?
- **Q4: How to detect outbreaks?**
- Q5: Who are the culprits?

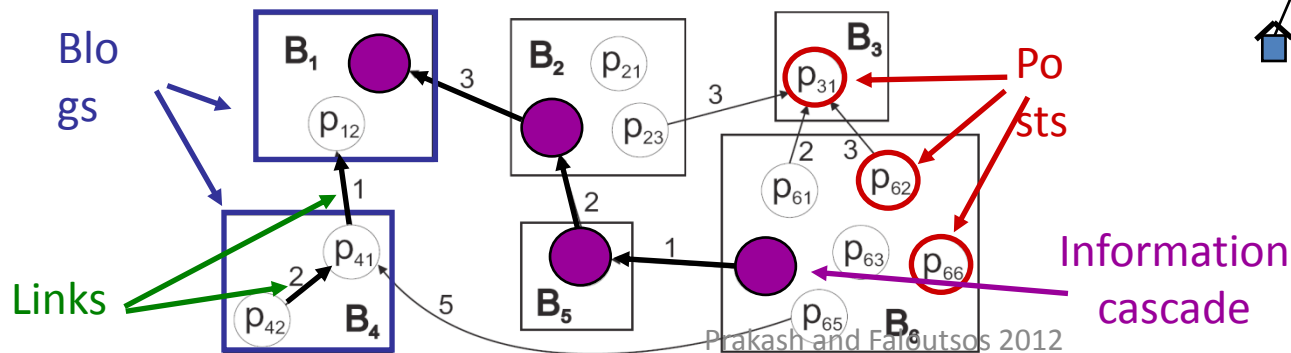
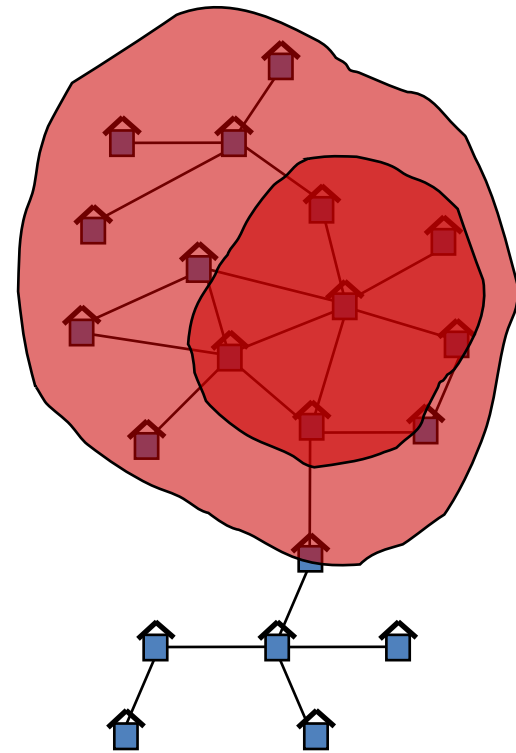
Break!

Part 2: Algorithms

- Q3: Whom to immunize?
- **Q4: How to detect outbreaks?**
- Q5: Who are the culprits?

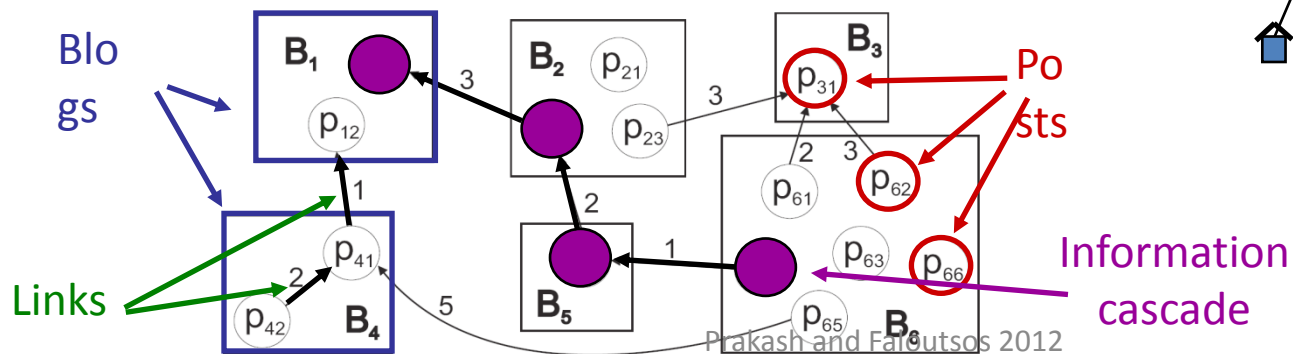
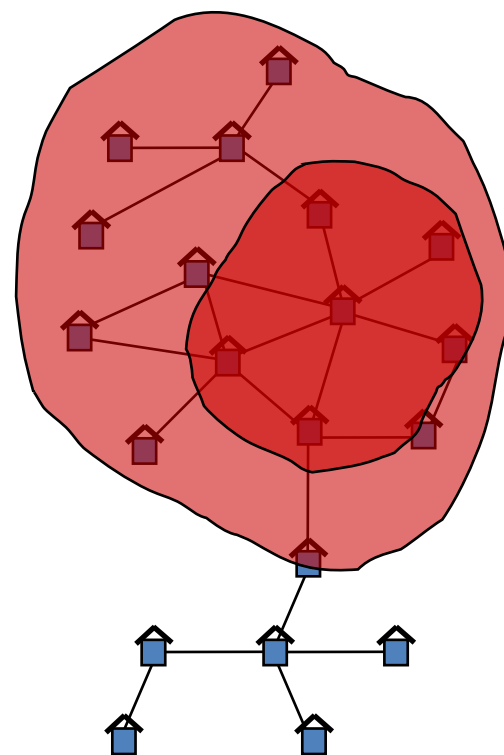
Outbreak detection

- Spot contamination points
 - Minimize time to detection, population affected
 - Maximize probability of detection.
 - Minimize sensor placement cost.

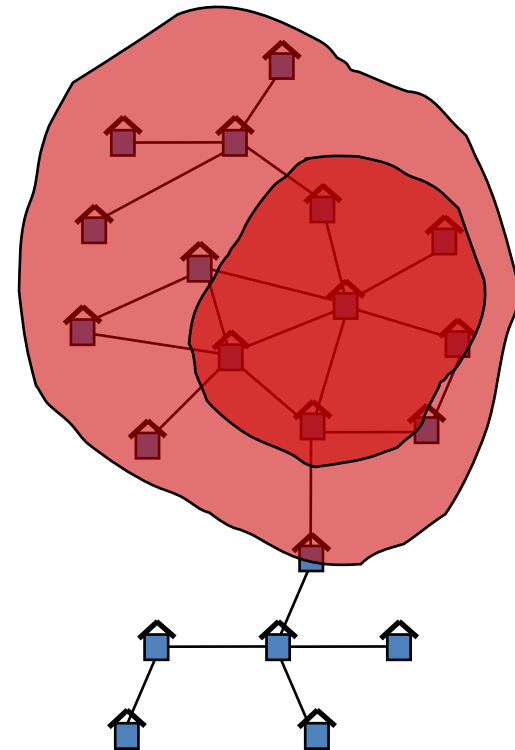


Outbreak detection

- Spot 'hot blogs'
 - Minimize time to detection, population affected
 - Maximize probability of detection.
 - Minimize sensor placement cost.



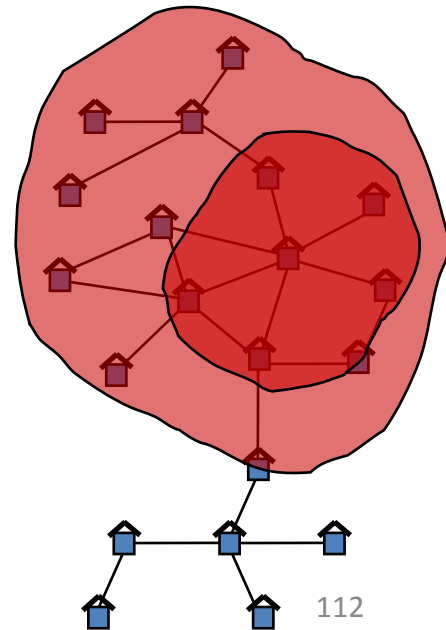
- J. Leskovec, A. Krause, C. Guestrin, C. Faloutsos, J. VanBriesen, N. Glance. "Cost-effective Outbreak Detection in Networks"
KDD 2007



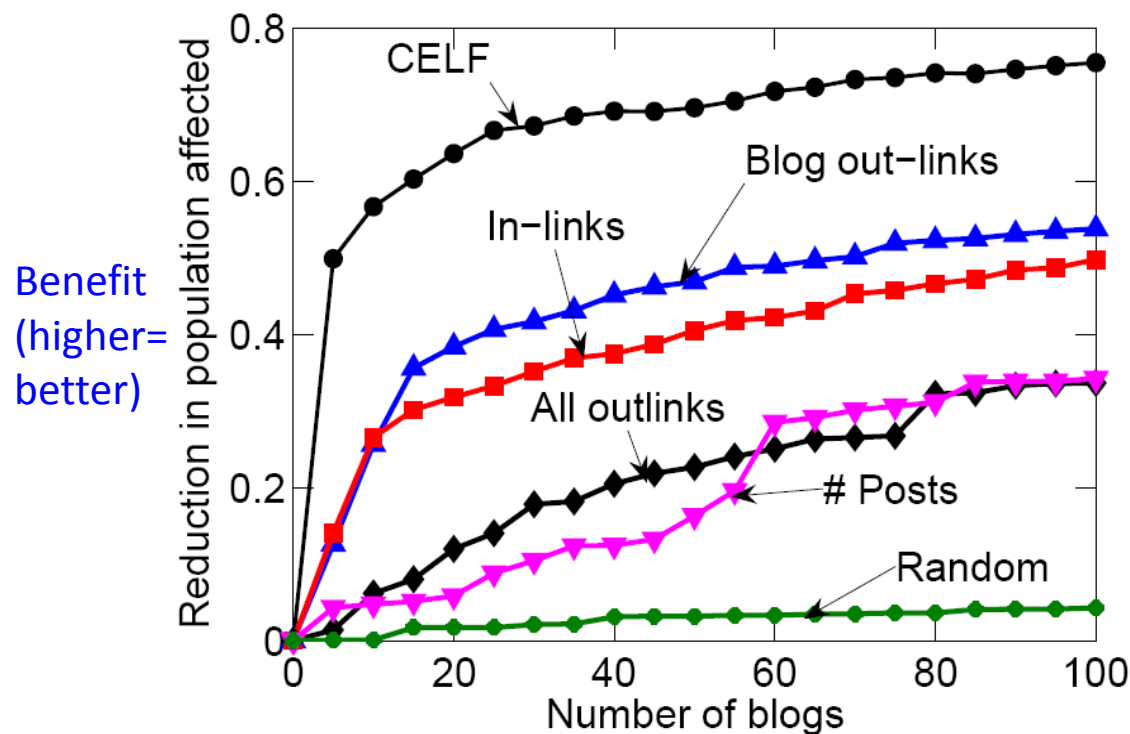
CELF: Main idea

- Given: a graph $G(V,E)$
 - a budget of B sensors
 - data on how contaminations spread over the network:
- Place the sensors
- To minimize time to detect outbreak

CELF algorithm uses **submodularity**
and **lazy evaluation**



Blogs: Comparison to heuristics



“Best 10 blogs to read”

NP - number of posts, IL- in-links, OLO- blog out links, OLA- all out links

• k	PA score	Blog	NP	IL	OLO	OLA
• 1	0.1283	http://instapundit.com	4593	4636	1890	5255
• 2	0.1822	http://donsurber.blogspot.com	1534	1206	679	3495
• 3	0.2224	http://sciencepolitics.blogspot.com	924	576	888	2701
• 4	0.2592	http://www.watcherofweasels.com	261	941	1733	3630
• 5	0.2923	http://michellemalkin.com	1839	12642	1179	6323
• 6	0.3152	http://blogometer.nationaljournal.com	189	2313	3669	9272
• 7	0.3353	http://themodulator.org	475	717	1844	4944
• 8	0.3508	http://www.bloggersblog.com	895	247	1244	10201
• 9	0.3654	http://www.boingboing.net	5776	6337	1024	6183
• 10	0.3778	http://atrios.blogspot.com	4682	3205	795	3102

Part 2: Algorithms

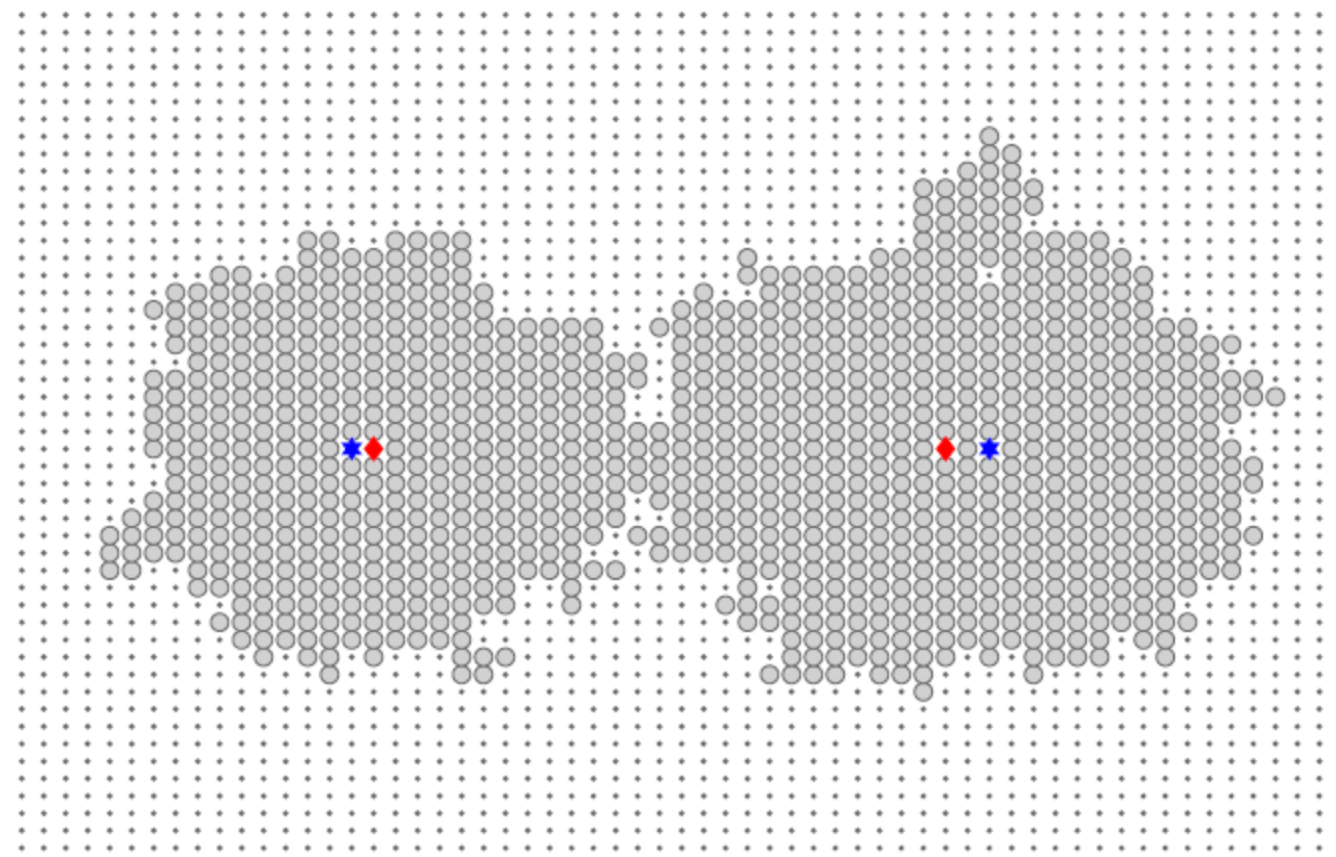
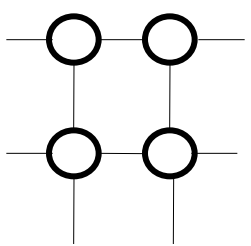
- Q3: Whom to immunize?
- Q4: How to detect outbreaks?
- **Q5: Who are the culprits?**

- B. Aditya Prakash, Jilles Vreeken, Christos Faloutsos ‘Detecting Culprits in Epidemics: Who and How many?’

ICDM 2012, Brussels

Problem definition

2-d grid
 '+' -> infected
 Who started it?

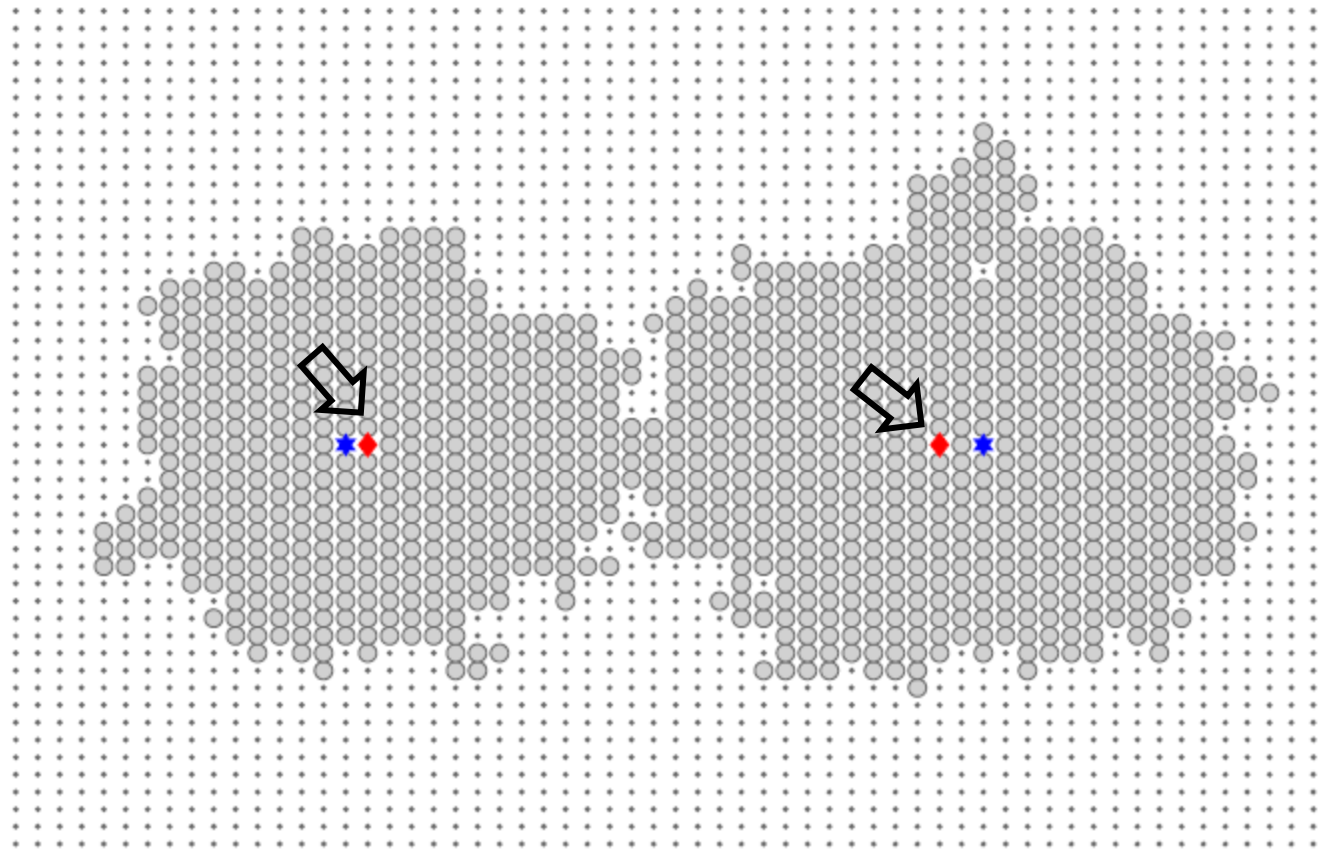
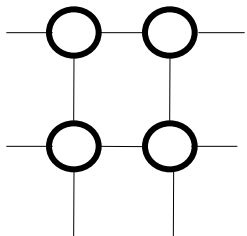


Problem definition

2-d grid

'+' -> infected

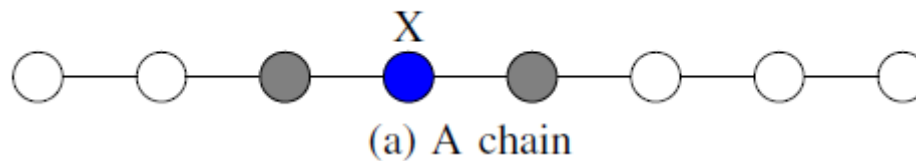
Who started it?



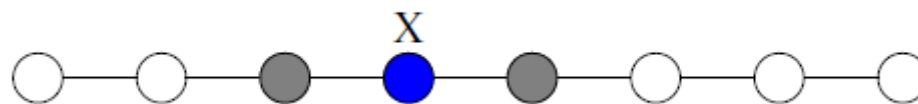
Prior work:

[Lappas et al.
2010, Shah et al.
2011]

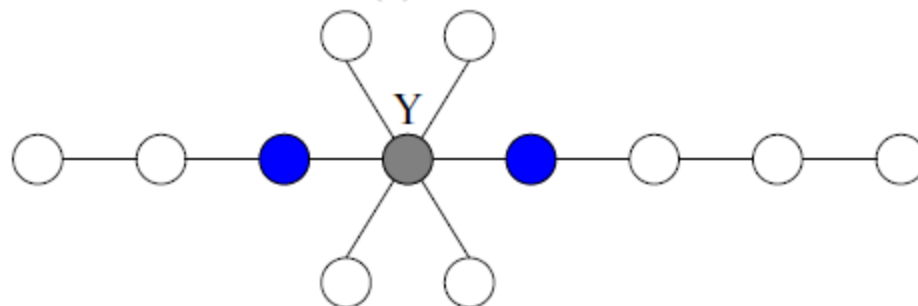
Culprits: Exoneration



Culprits: Exoneration



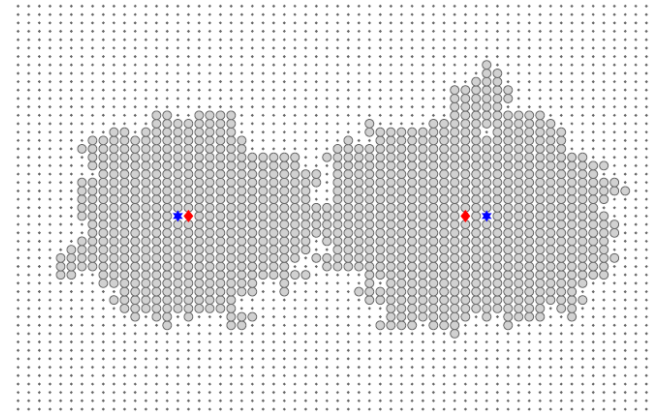
(a) A chain



(b) A chain-star

Who are the culprits

- Two-part solution
 - use MDL for *number* of seeds
 - for a given number:
 - exoneration = centrality + penalty
 - our method uses *smallest* eigenvector of Laplacian submatrix
- Running time = $O(k^*(\mathcal{E}_I + \mathcal{E}_F + \mathcal{V}_I))$
 - linear! (in edges and nodes)



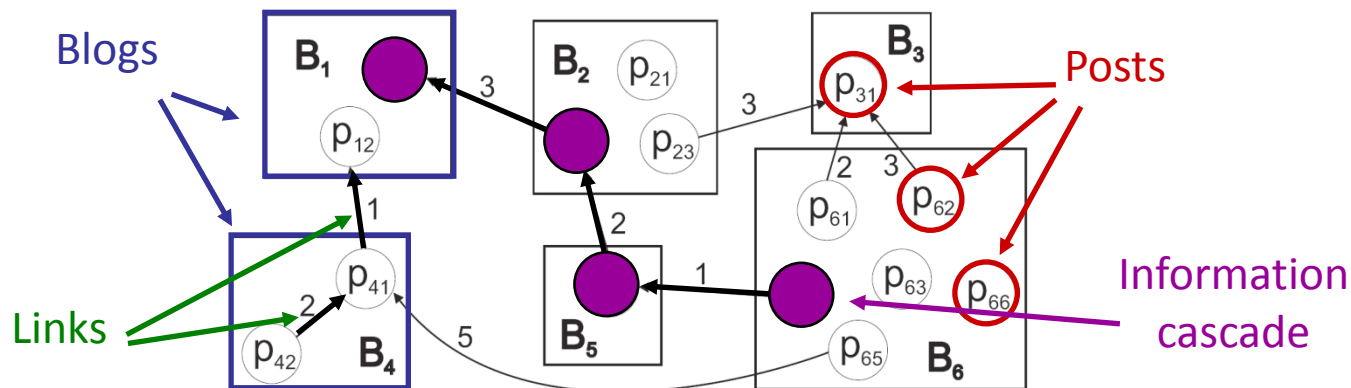
Outline

- Motivation
- Part 1: Understanding Epidemics (Theory)
- Part 2: Policy and Action (Algorithms)
- **Part 3: Learning Models (Empirical Studies)**
- Conclusion

Part 3: Empirical Studies

- **Q6: How do cascades look like?**
- Q7: How does activity evolve over time?
- Q8: How does external influence act?

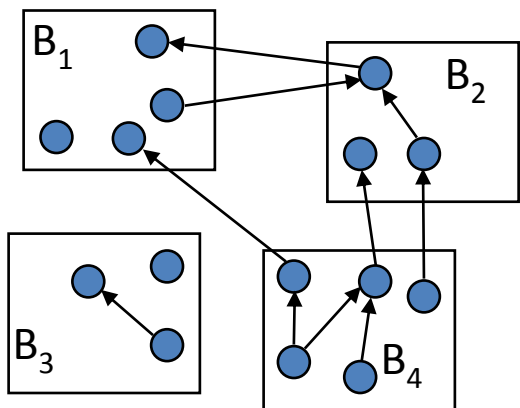
Cascading Behavior in Large Blog Graphs



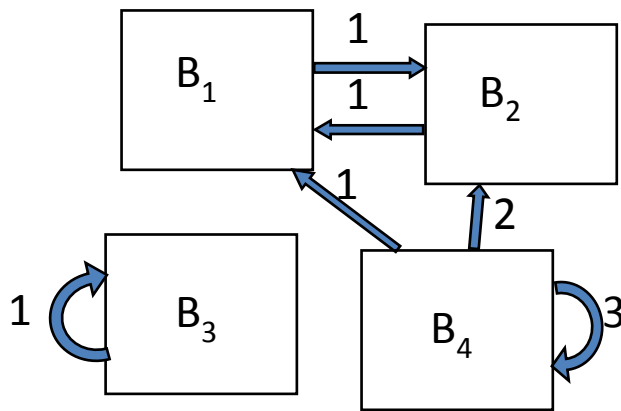
How does information propagate over the blogosphere?

J. Leskovec, M. McGlohon, C. Faloutsos, N. Glance, M. Hurst. Cascading Behavior in Large Blog Graphs. SDM 2007.

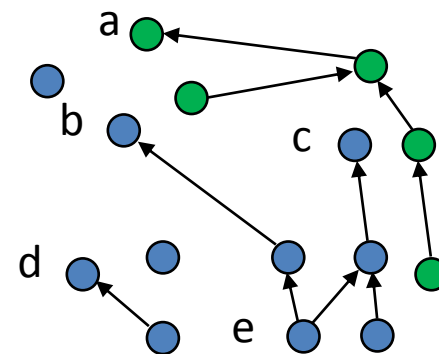
Cascades on the Blogosphere



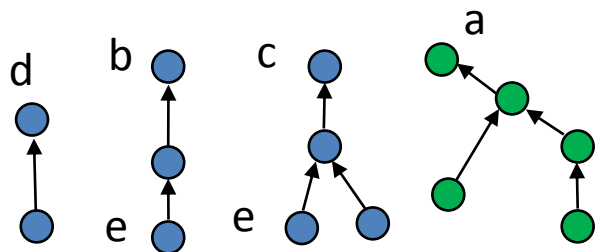
Blogosphere
blogs + posts



Blog network
links among blogs



Post network
links among posts

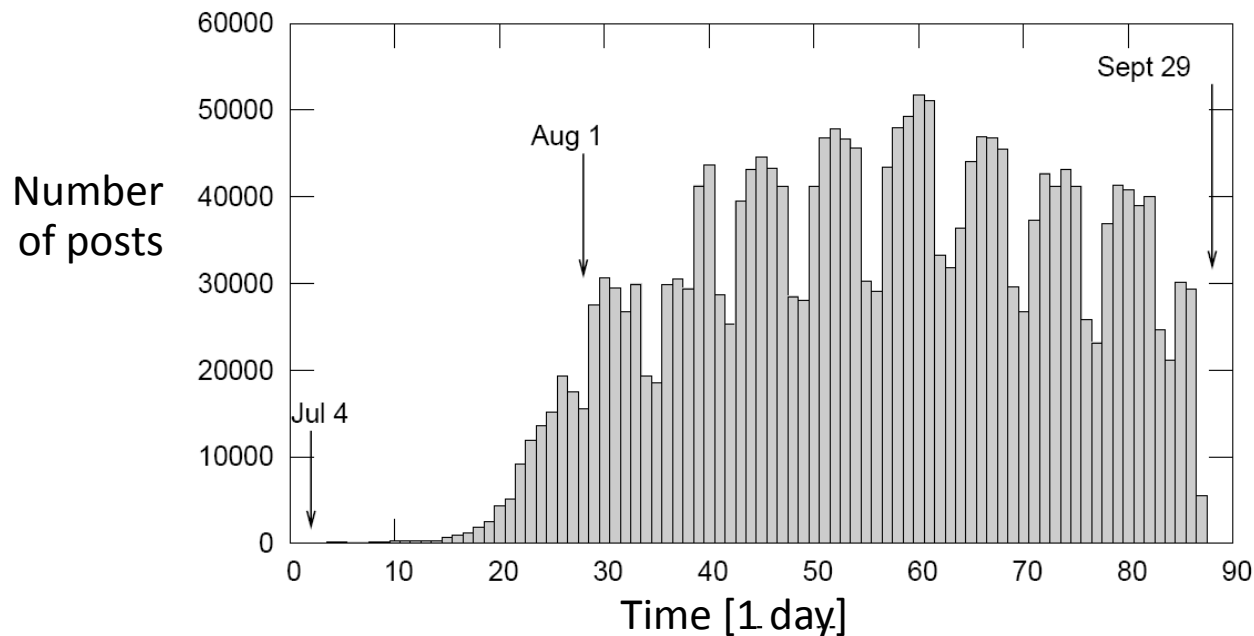


Cascades

Cascade is graph induced by a time ordered propagation of information (edges)

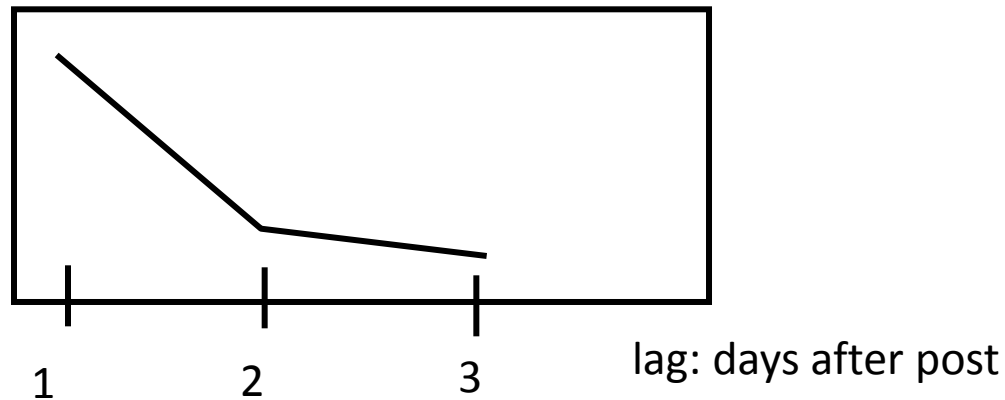
Blog data

- 45,000 blogs participating in cascades
- All their posts for 3 months (Aug-Sept '05)
- 2.4 million posts
- ~5 million links (245,404 inside the dataset)

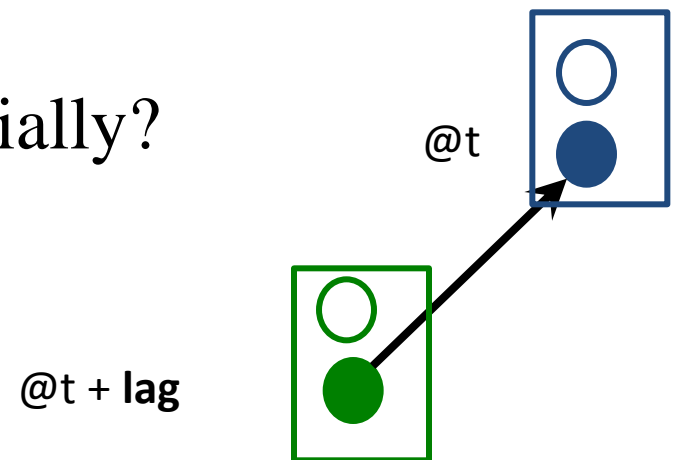


Popularity over time

in links

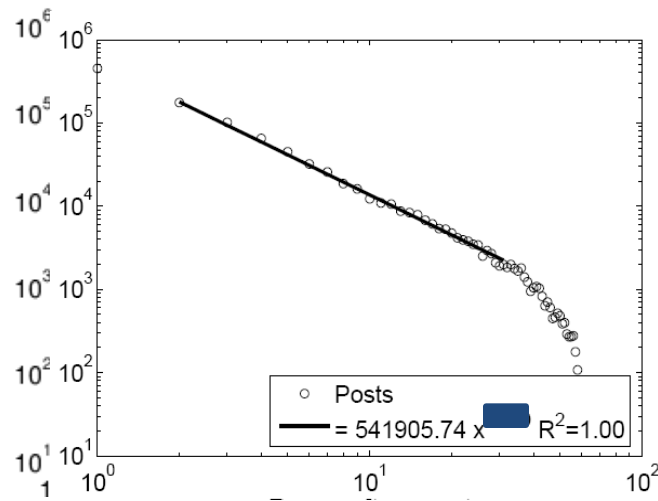


Post popularity drops-off – exponentially?



Popularity over time

in links
(log)



days after post
(log)

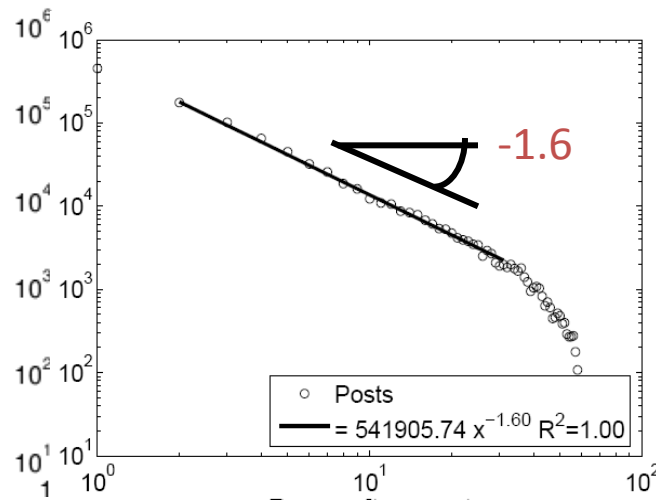
Post popularity drops-off – exponentially?

POWER LAW!

Exponent?

Popularity over time

in links
(log)

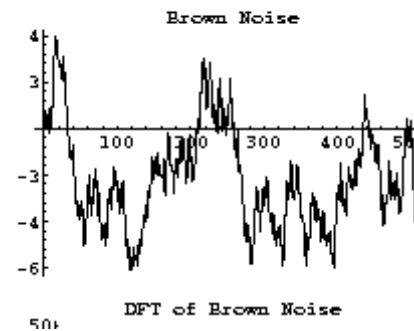


days after post
(log)

Post popularity drops-off – exponentially?
POWER LAW!

Exponent? -1.6

- close to -1.5: Barabasi's stack model
- and like the zero-crossings of a random walk



-1.5 slope

J. G. Oliveira & A.-L. Barabási Human Dynamics: The Correspondence Patterns of Darwin and Einstein. *Nature* **437**, 1251 (2005) . [[PDF](#)]

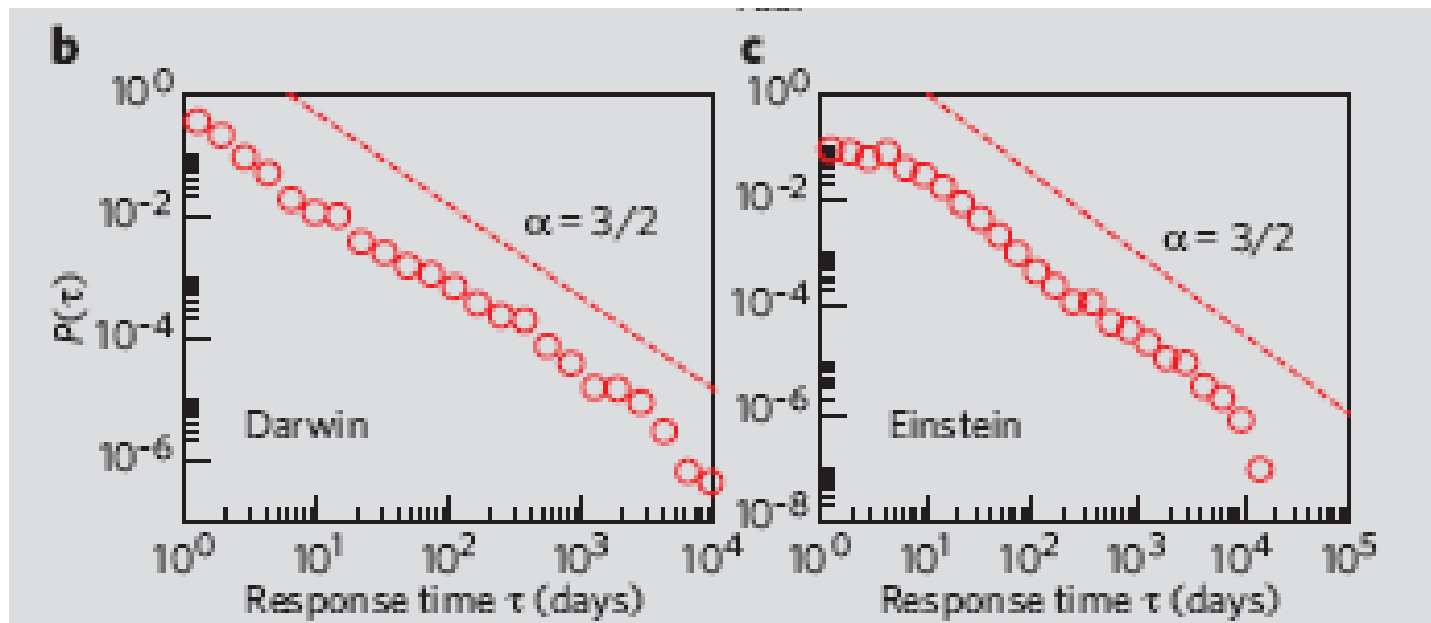
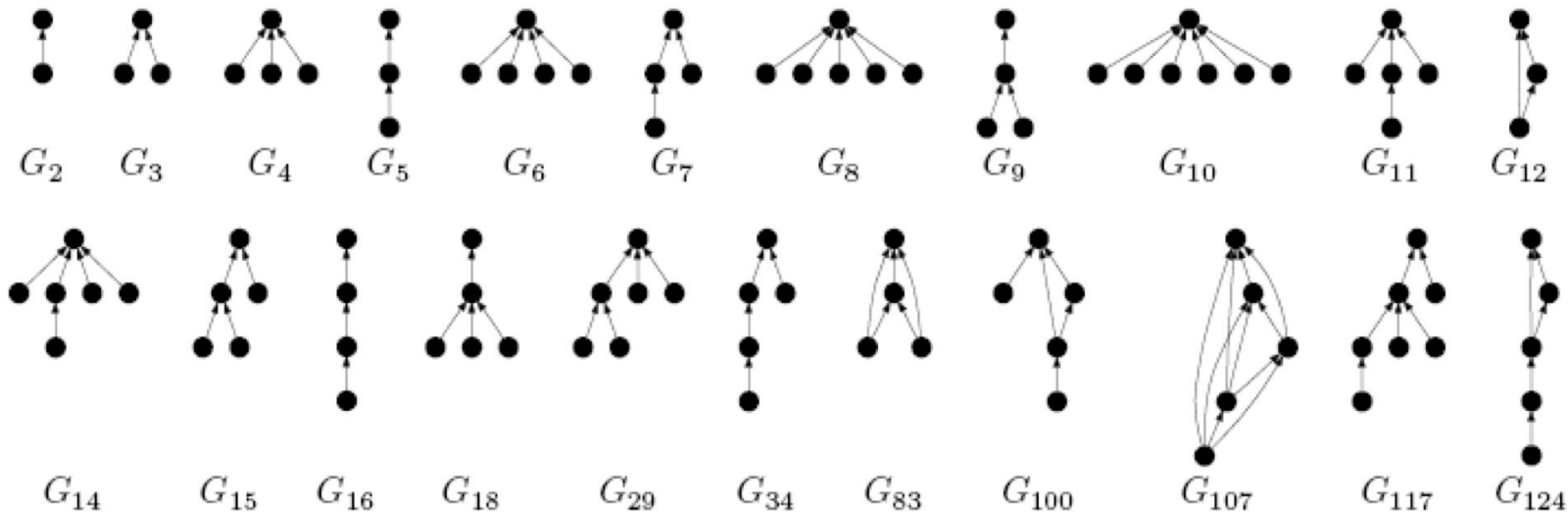


Figure 1 | The correspondence patterns of Darwin and Einstein.

Topological Observations

How do we measure how information flows through the network?

Common cascade shapes extracted using algorithms in [Leskovec, Singh, Kleinberg; PAKDD 2006].



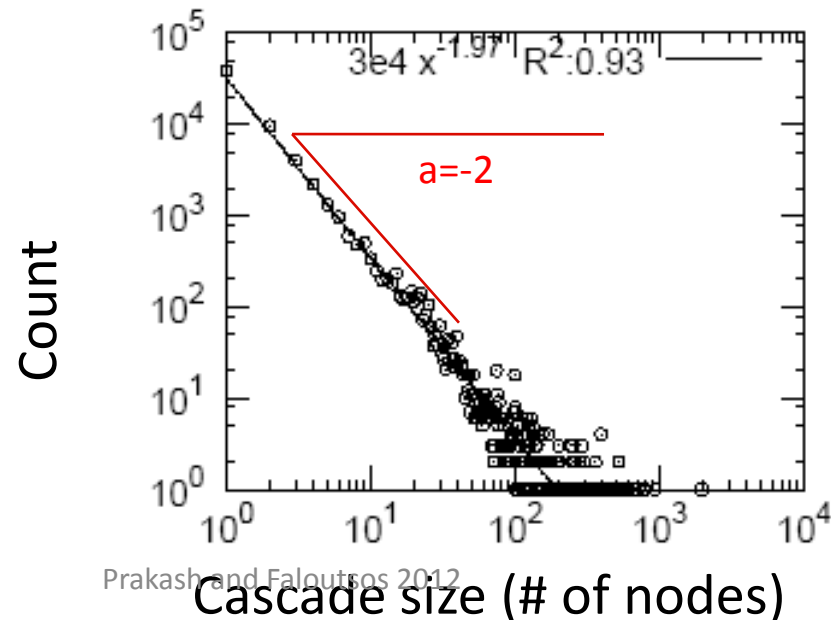
Topological Observations

What graph properties do cascades exhibit?

Cascade size distributions also follow power law.

Observation 2: *The probability of observing a cascade on n nodes follows a Zipf distribution:*

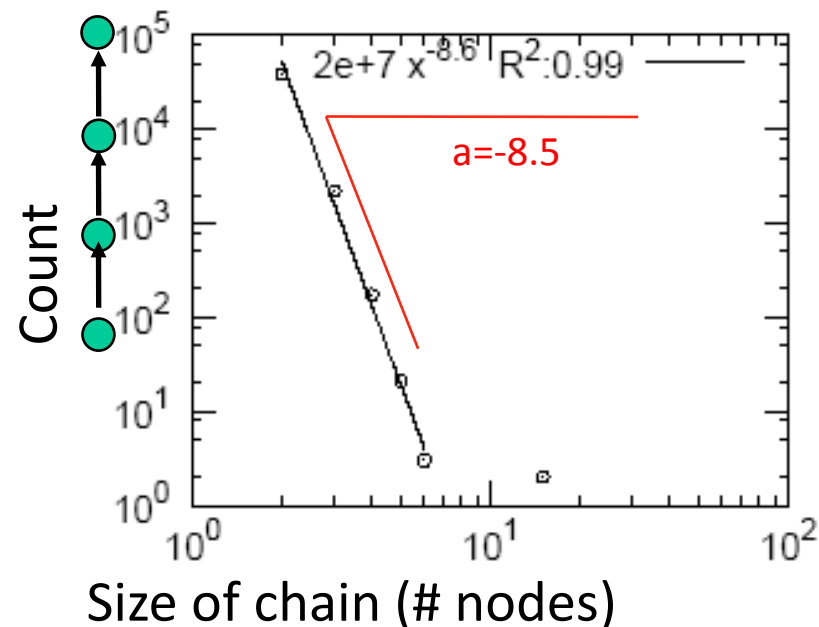
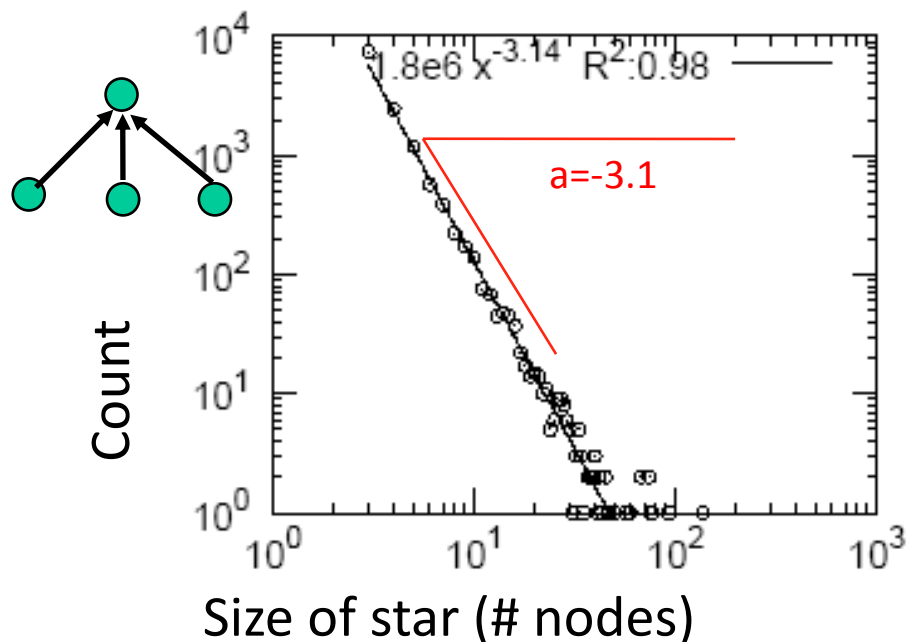
$$p(n) \propto n^{-2}$$



Topological Observations

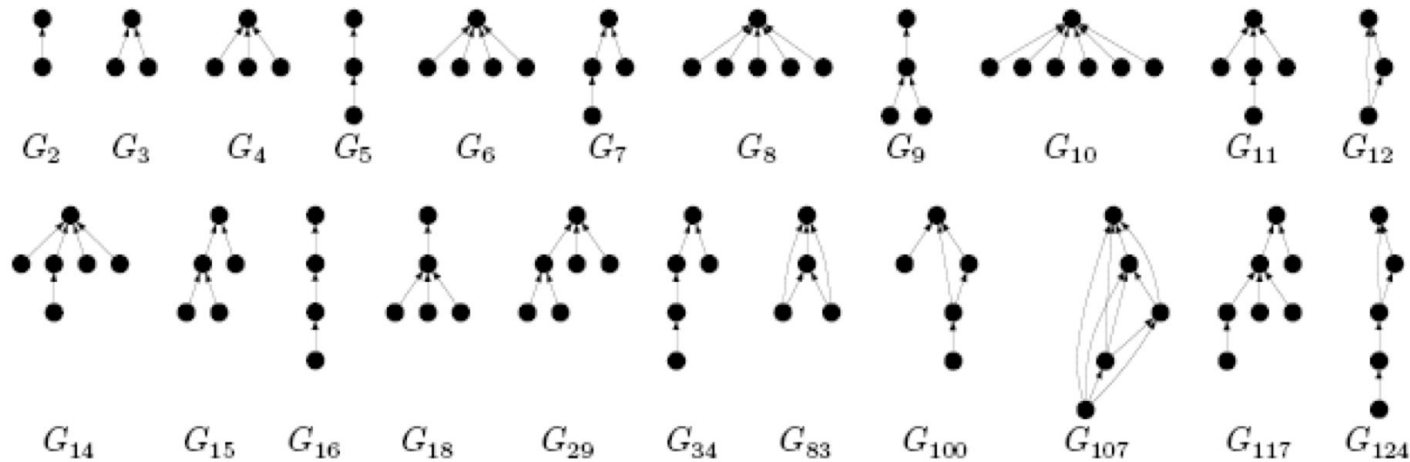
What graph properties do cascades exhibit?

Stars and chains also follow a power law, with different exponents (star -3.1, chain -8.5).



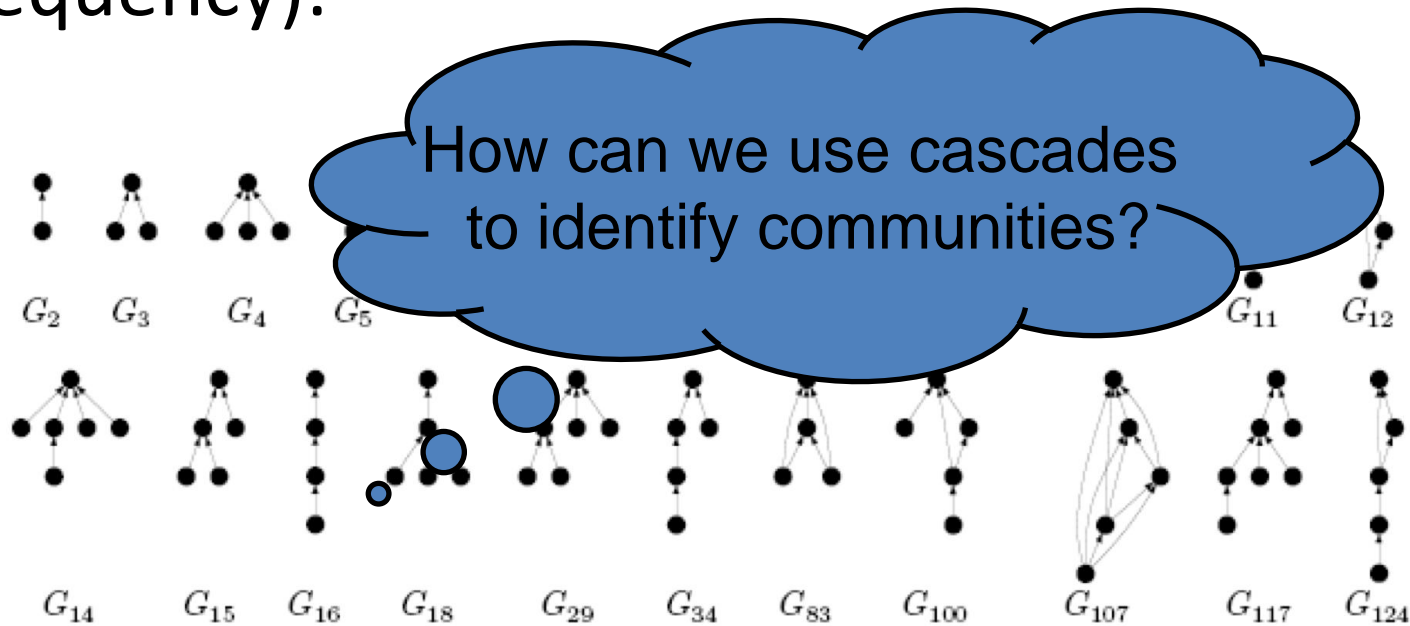
Blogs and structure

- Cascades take on different shapes (sorted by frequency):



Blogs and structure

- Cascades take on different shapes (sorted by frequency):



PCA on cascade types

- Perform PCA on sparse matrix.
- Use $\log(\text{count}+1)$
- Project onto 2 PC...

~9,000 cascade types

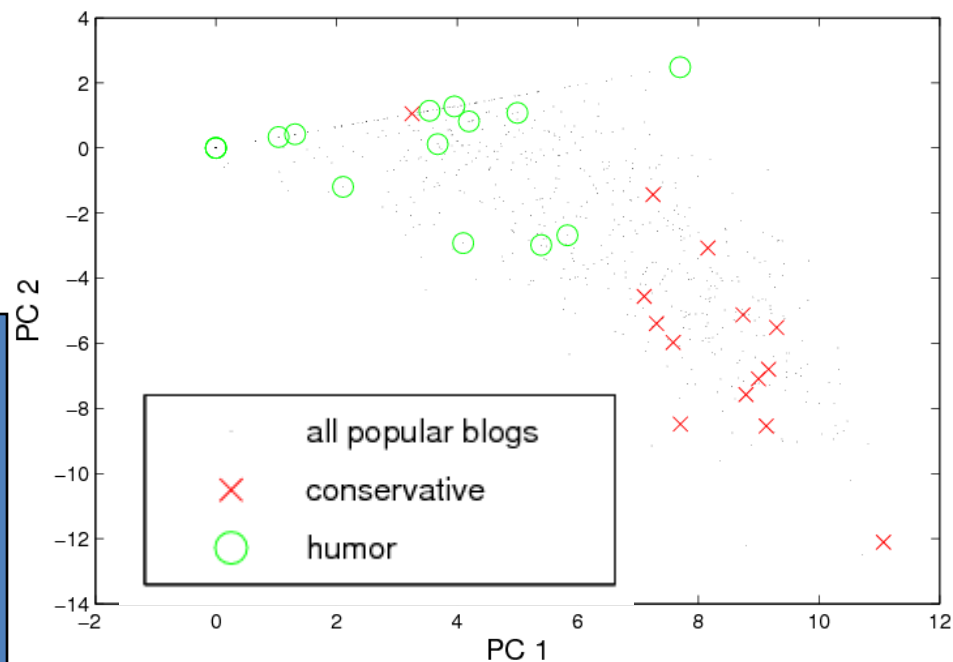


~44,000 blogs

<i>slashdot</i>	4.6	2.1	.09			
<i>boingboing</i>	3.2	1.1		3.4	.07	
...	4.2					
...	5.1					
...	2.1		1.1			
...	.67			.07		
...	.01					

PCA on cascade types

- Observation: Content of blogs and cascade behavior are often related.
- Distinct clusters for “conservative” and “humorous” blogs (hand-labeling).

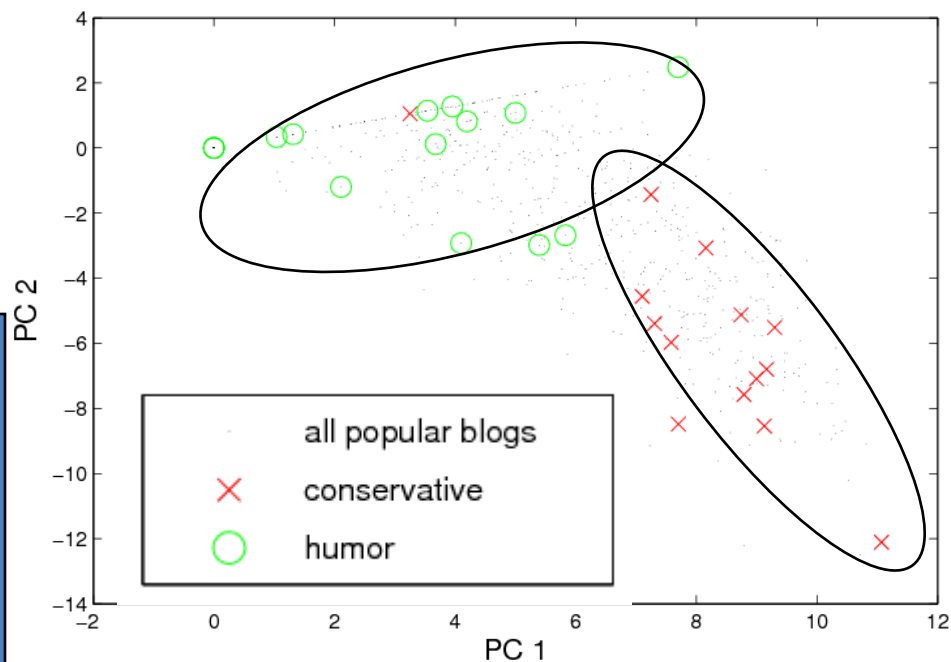


M. McGlohon, J. Leskovec, C. Faloutsos, M. Hurst, N. Glance. Finding Patterns in Blog Shapes and Blog Evolution. ICWSM 2007.

PCA on cascade types

- Observation: Content of blogs and cascade behavior are often related.
- Distinct clusters for “conservative” and “humorous” blogs (hand-labeling).

M. McGlohon, J. Leskovec, C. Faloutsos, M. Hurst, N. Glance. Finding Patterns in Blog Shapes and Blog Evolution. ICWSM 2007.



Part 3: Empirical Studies

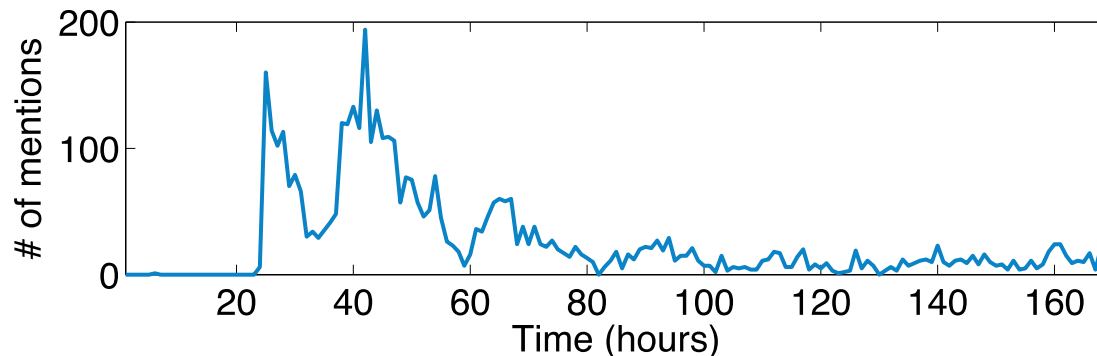
- Q6: How do cascades look like?
- **Q7: How does activity evolve over time?**
- Q8: How does external influence act?

Rise and fall patterns in social media

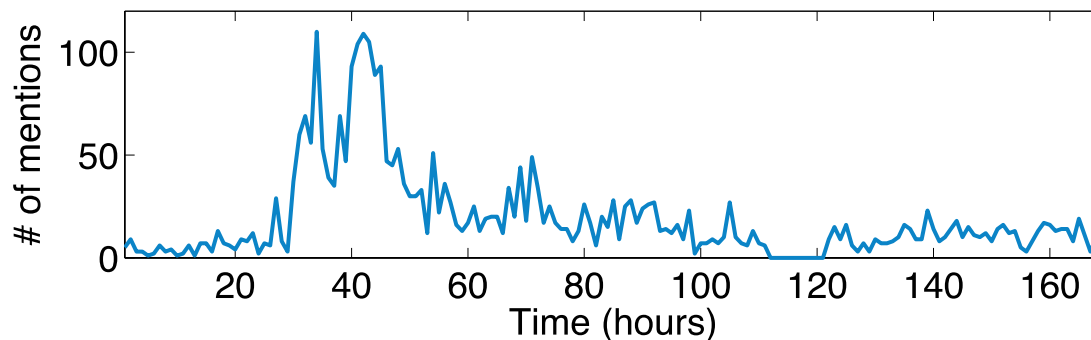
- Meme (# of mentions in blogs)

- short phrases Sourced from U.S. politics in 2008

“you can put lipstick on a pig”



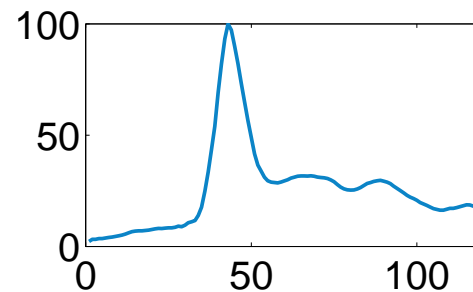
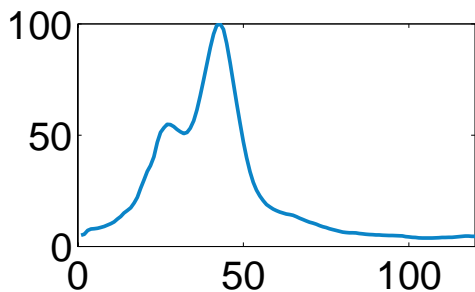
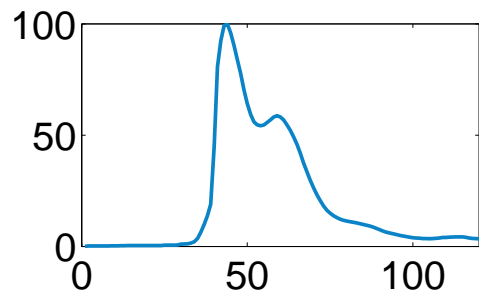
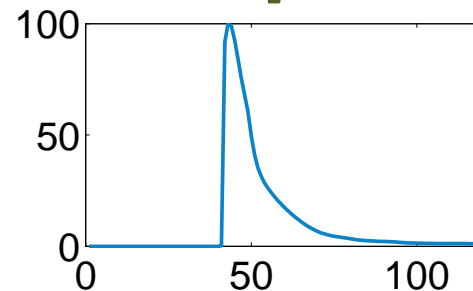
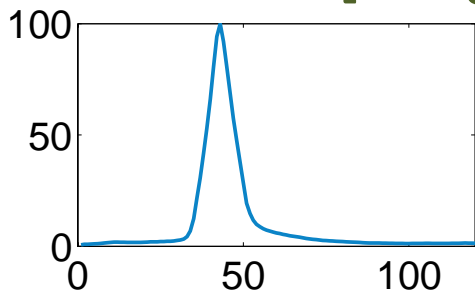
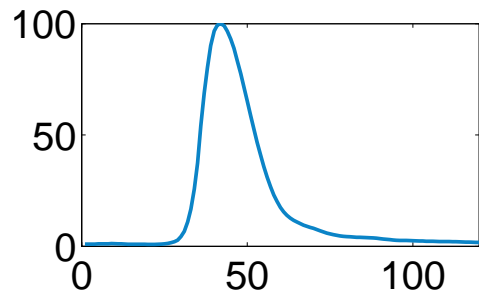
“yes we can”



Rise and fall patterns in social media

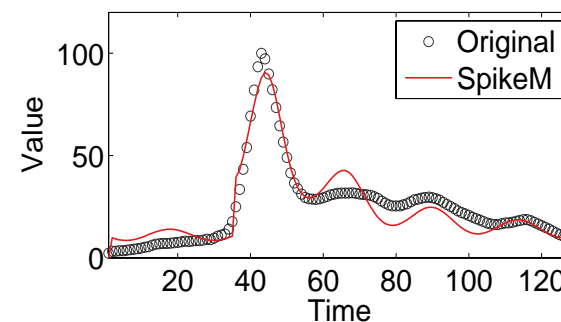
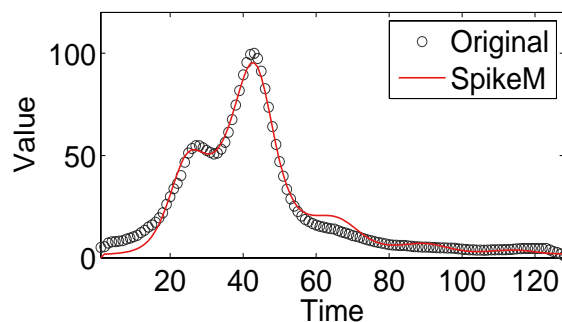
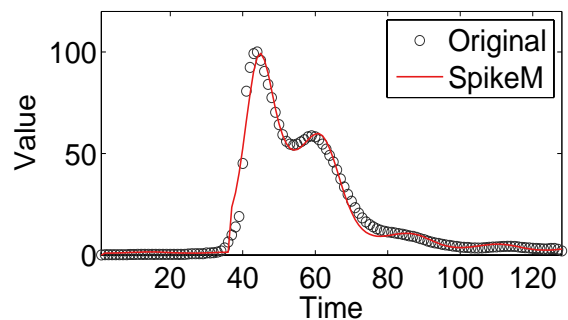
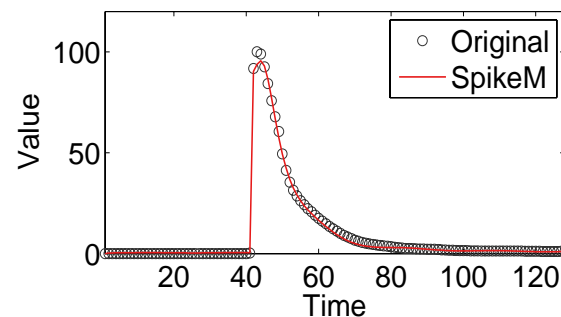
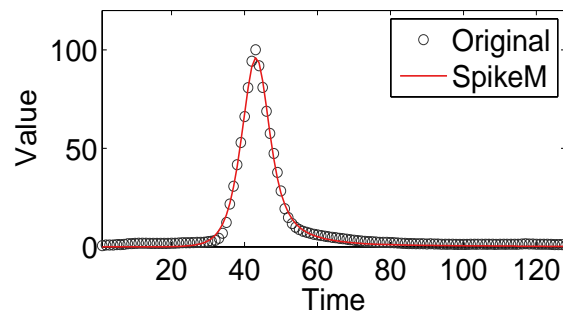
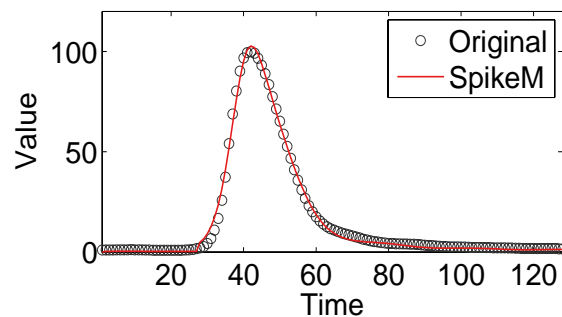
- Can we find a unifying model, which includes these patterns?
 - **four** classes on YouTube [Crane et al. '08]

- **six** classes on Meme [Yang et al. '11]

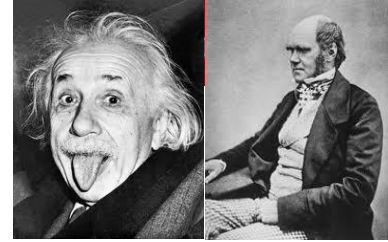


Rise and fall patterns in social media

- Answer: YES!

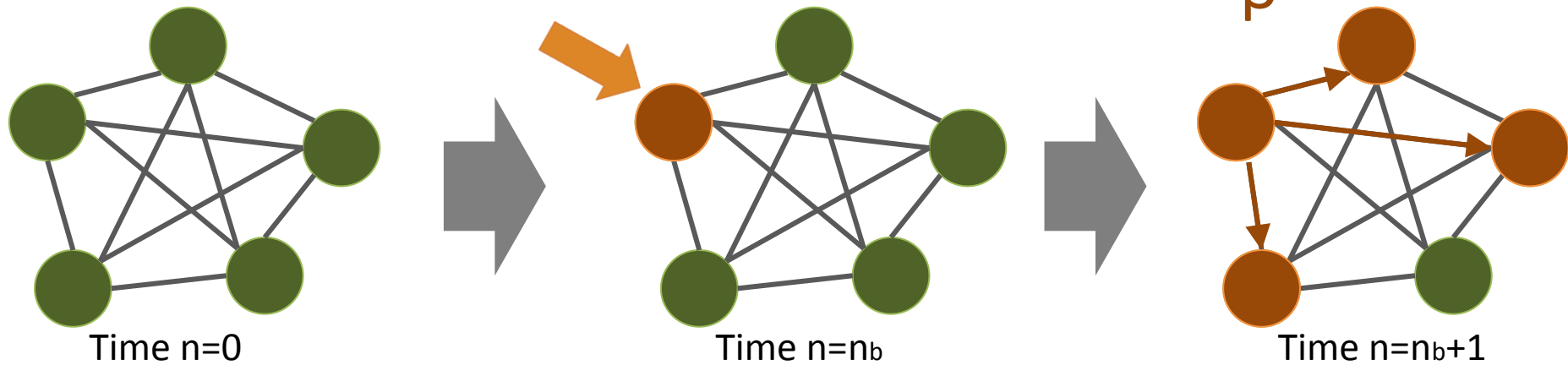


- We can represent all patterns by **single model**



Main idea - SpikeM

- 1. **Un-informed bloggers** (uninformed about rumor)
- 2. **External shock** at time n_b (e.g, breaking news)
- 3. **Infection** (word-of-mouth)

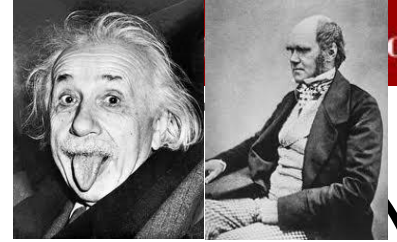


Infectiveness of a blog-post at age n :

$$f(n) = b * n^{-1.5}$$

- b - Strength of infection (quality of news)
- $f(n)$ - Decay function (how infective a blog posting is)

Power Law



-1.5 slope

J. G. Oliveira & A.-L. Barabási Human Dynamics: The Correspondence Patterns of Darwin and Einstein. *Nature* **437**, 1251 (2005) . [[PDF](#)]

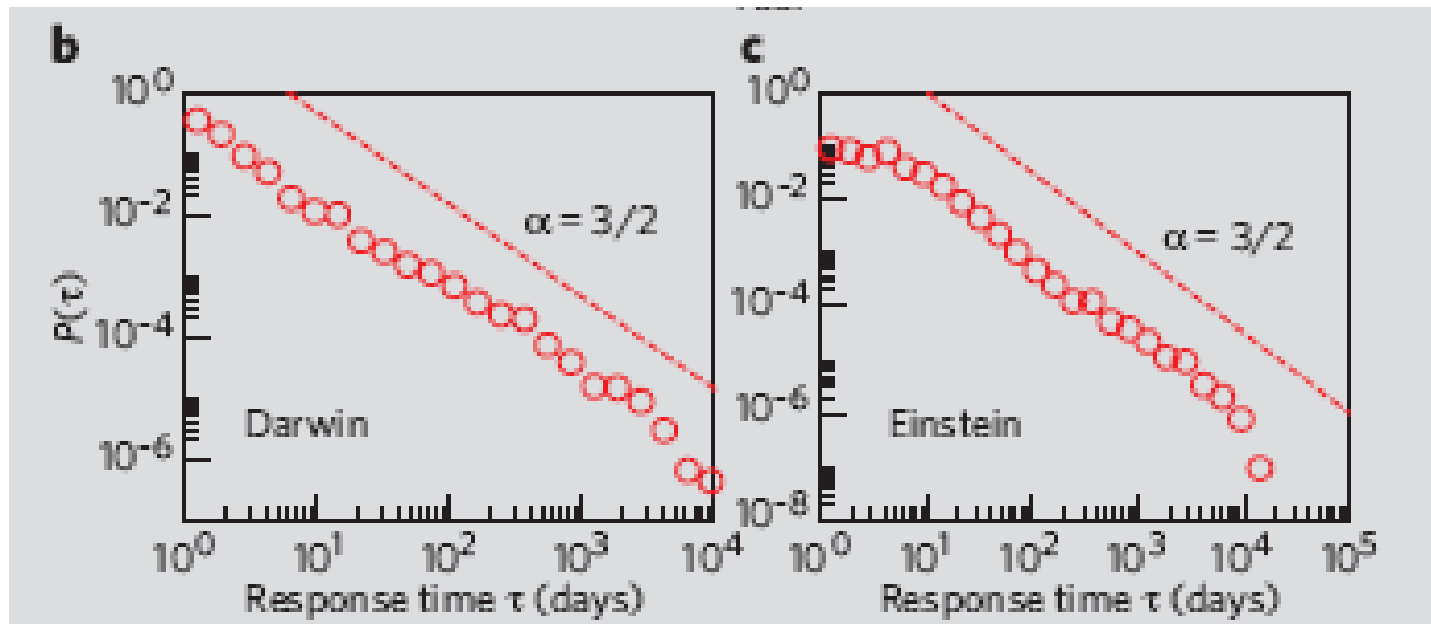


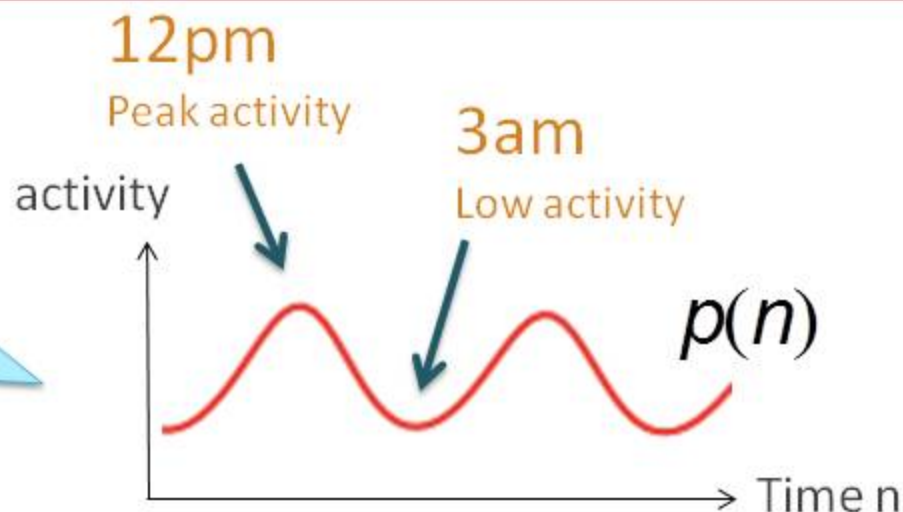
Figure 1 | The correspondence patterns of Darwin and Einstein.

SpikeM - with periodicity

- Full equation of SpikeM

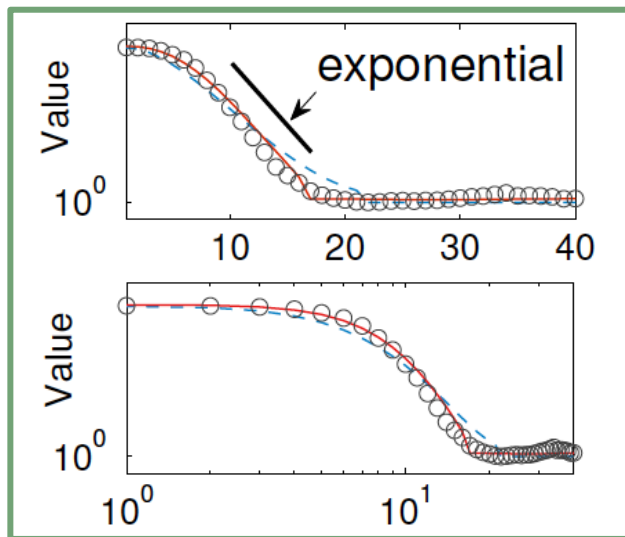
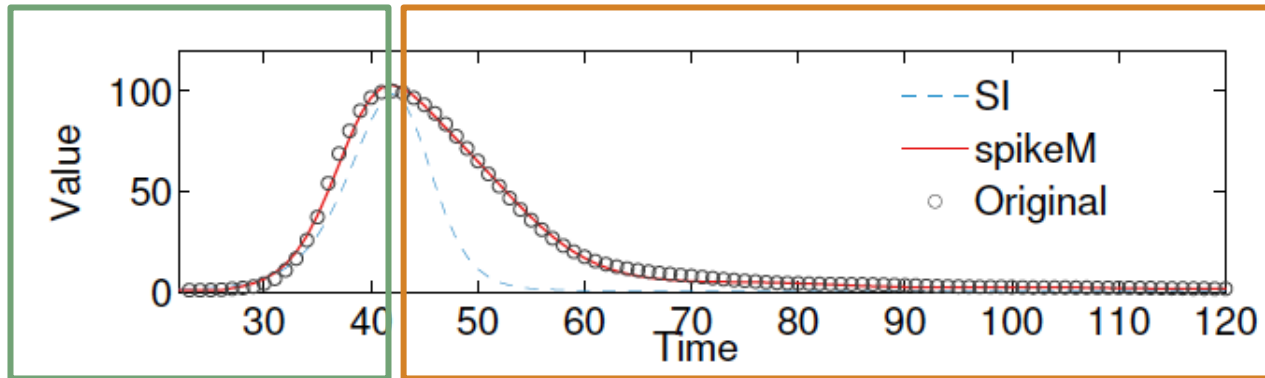
$$\Delta B(n+1) = \underbrace{p(n+1)}_{\text{Periodicity}} \cdot \left[U(n) \cdot \sum_{t=n_b}^n (\Delta B(t) + S(t)) \cdot f(n+1-t) + \varepsilon \right]$$

Bloggers change their activity over time
(e.g., daily, weekly, yearly)



Details

- Analysis – exponential rise and power-law fall



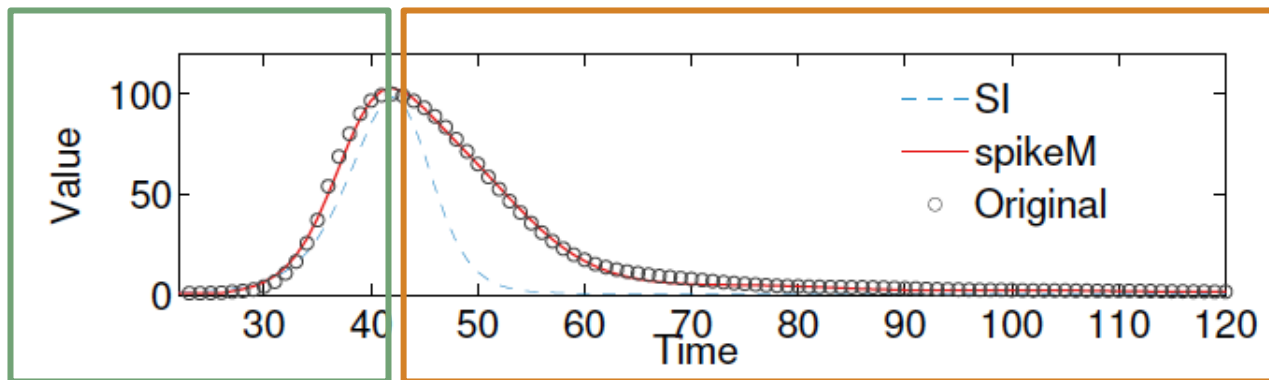
Rise-part

SI -> exponential

SpikeM -> exponential

Details

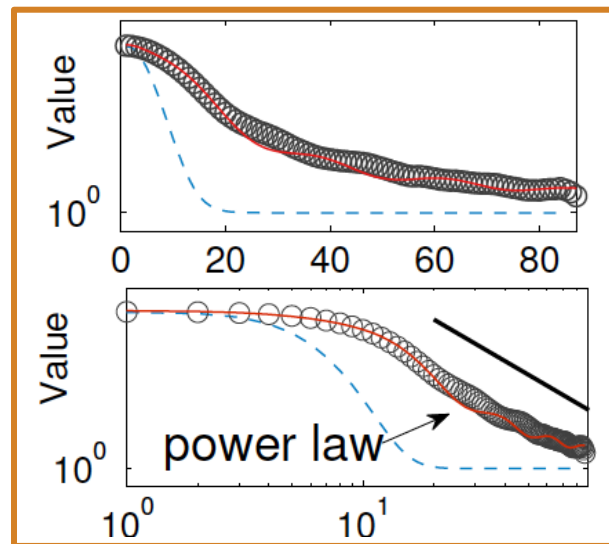
- Analysis – exponential rise and power-law fall



Fall-part

✗ SI → exponential

SpikeM → power law

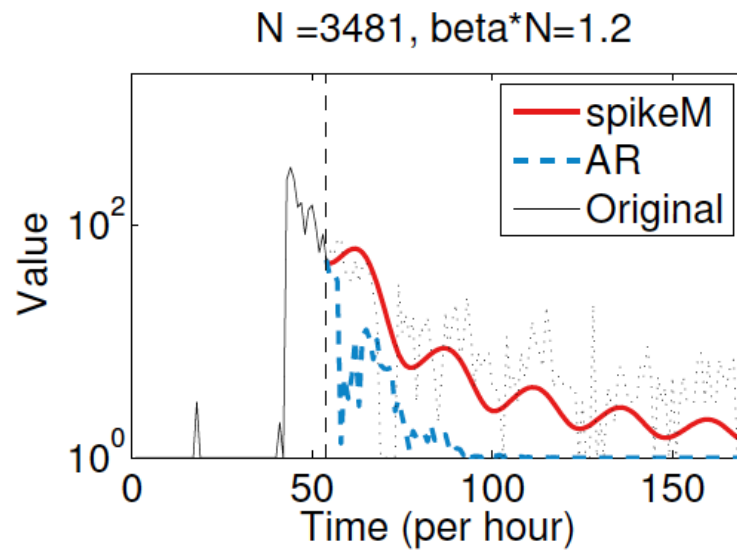
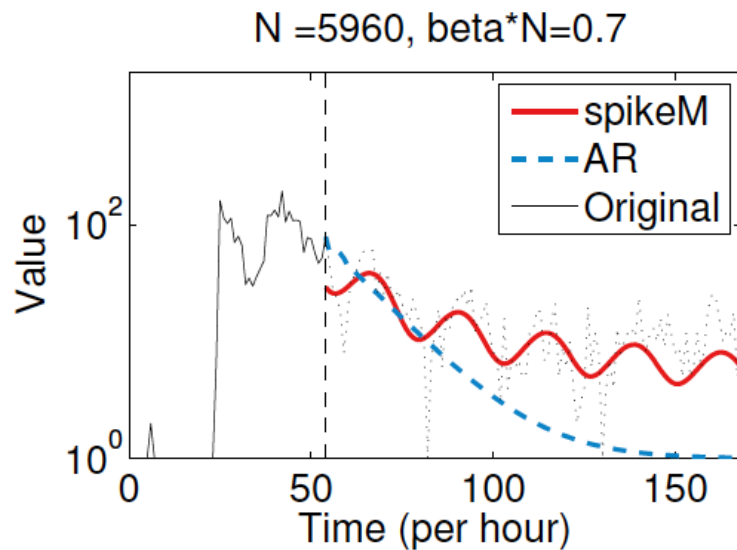


Liner-log

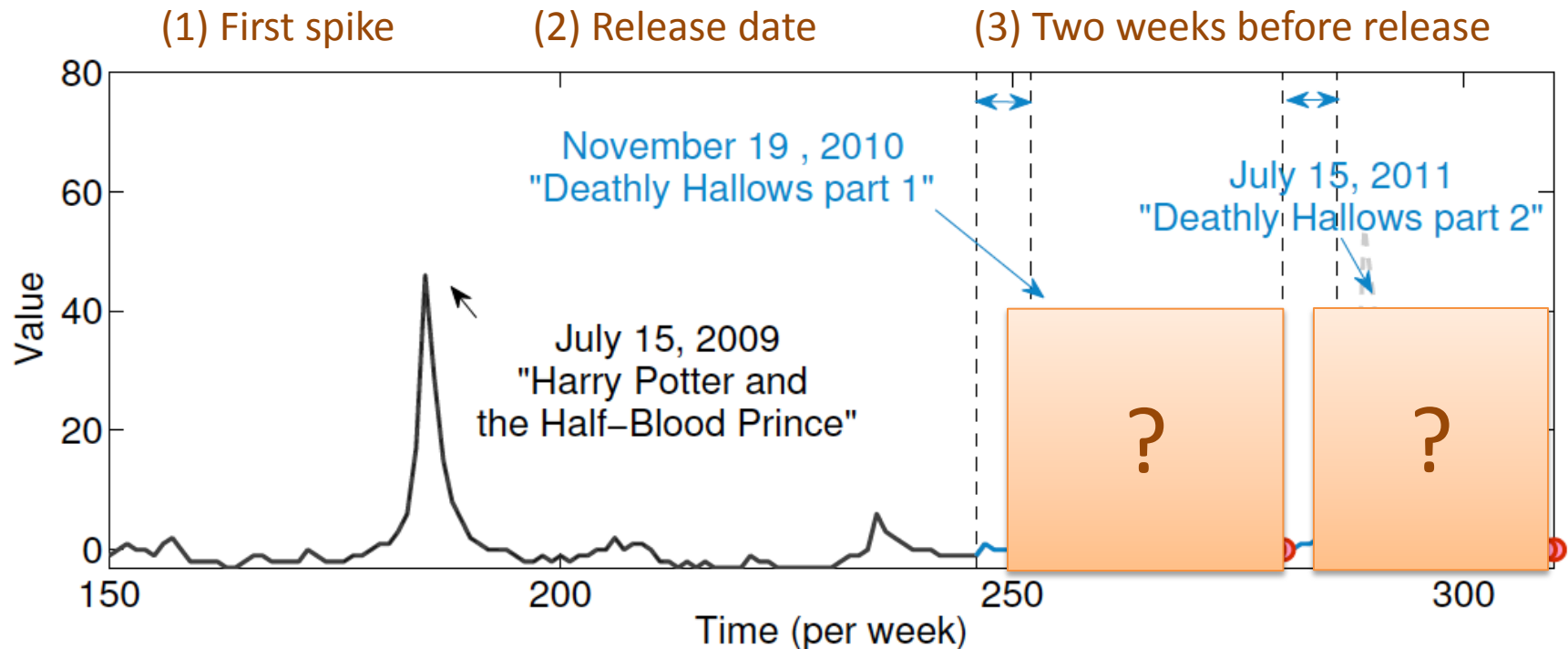
Log-log

Tail-part forecasts

- **SpikeM** can capture tail part



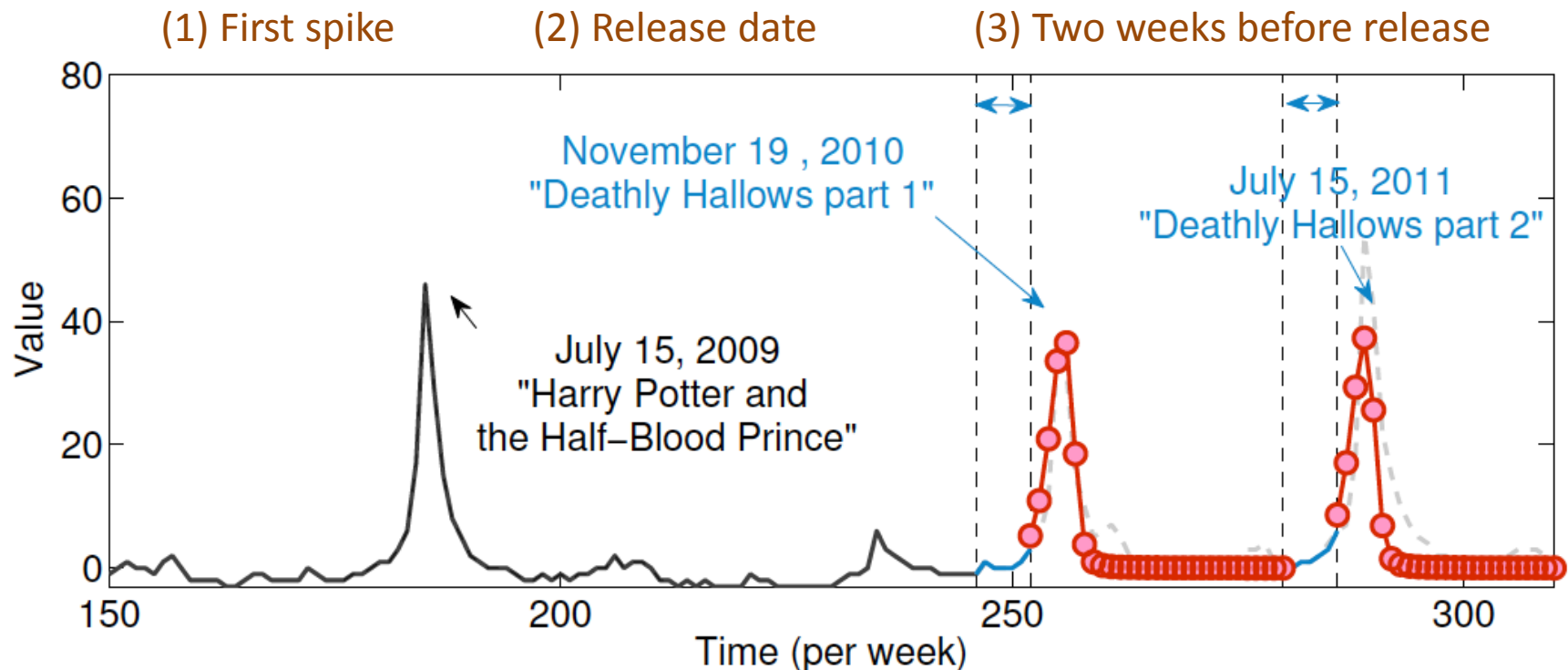
“What-if” forecasting



- e.g., given
- (1) first spike,
 - (2) release date of two sequel movies
 - (3) access volume before the release date

“What-if” forecasting

–SpikeM can forecast not only tail-part, but also **rise-part!**





- SpikeM can forecast **upcoming spikes**

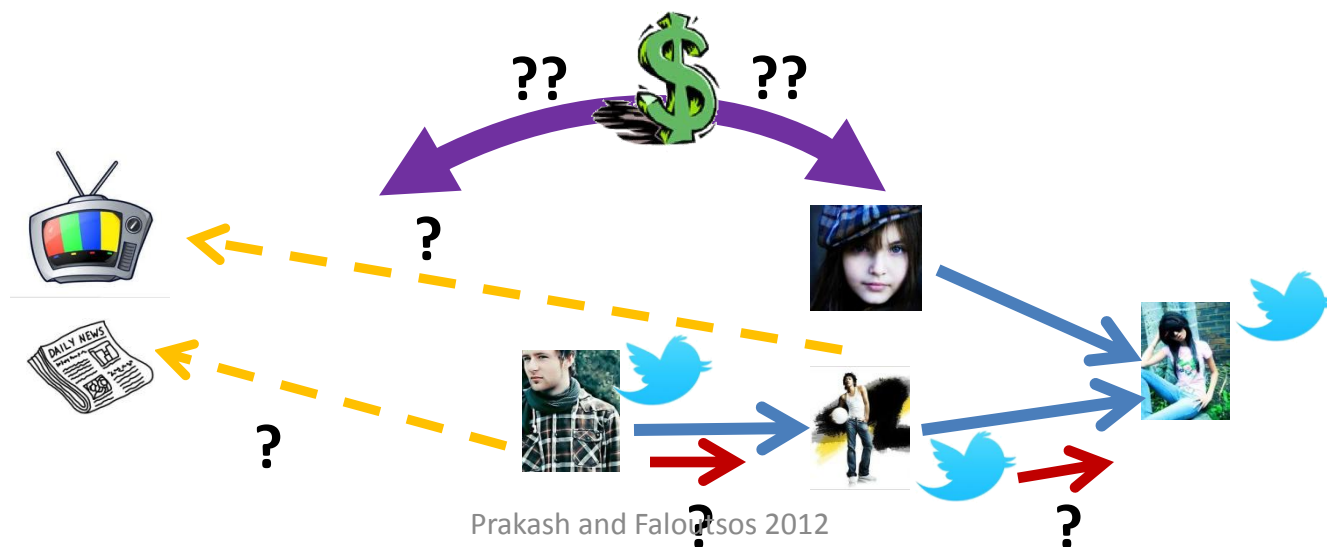
Part 3: Empirical Studies

- Q6: How do cascades look like?
- Q7: How does activity evolve over time?
- **Q8: How does external influence act?**

Tweets Diffusion: Problem Definition

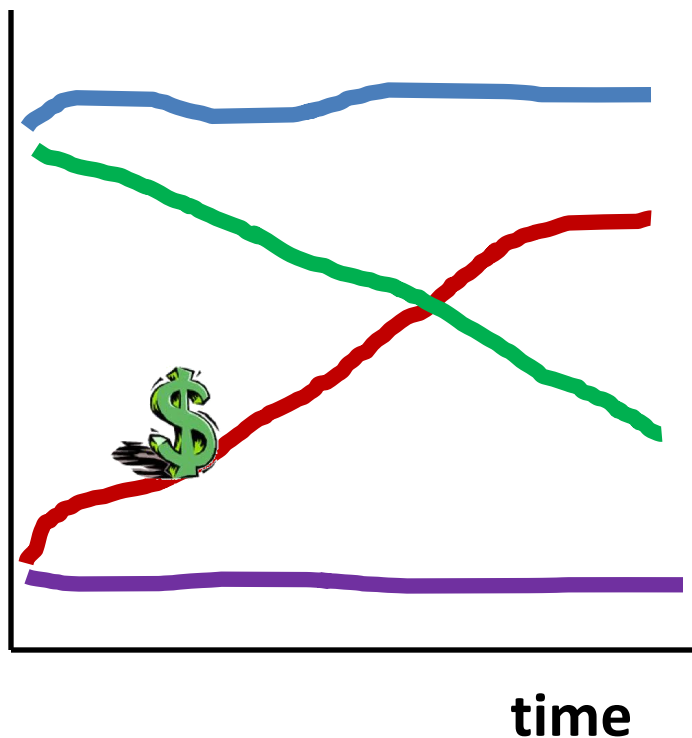


- Given:
 - Action log of people tweeting a #hashtag ()
 - A network of users ()
- Find:
 - How external influence varies with #hashtags?



Results: External Influence vs Time

“External Effects”



Long-running tags
 #nowwatching, #nowplaying, #epictweets

“Word-of-mouth”
 #purpleglasses, #brits, #famouslies

Bursty, external events
 #oscar, #25jan

“Word-of-mouth”
 Not trending
 #openfollow, #ihatequotes, #tweetmyjobs

Can also use for Forecasting, Anomaly Detection!

Outline

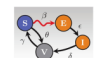
- Motivation
- Part 1: Understanding Epidemics (**Theory**)
- Part 2: Policy and Action (**Algorithms**)
- Part 3: Learning Models (**Empirical Studies**)
- **Conclusion**

Conclusions

- Epidemic Threshold
 - It's the Eigenvalue
- Fast Immunization
 - Max. drop in eigenvalue, linear-time near-optimal algorithm
- Bursts: SpikeM model
 - Exponential growth, Power-law decay

Our thresholds for some models

- $s = \text{effective strength}$
- $s < 1$: below threshold




Models	Effective Strength (s)	Threshold (tipping point)
SIS, SIR, SIRS, SEIR	$s = \lambda \cdot \left(\frac{\beta}{\delta} \right)$	$s = 1$
SIV, SEIV	$s = \lambda \cdot \left(\frac{\beta \gamma}{\delta(\gamma + \theta)} \right)$	
SI, I ₂ , V ₁ , V ₂ (H.I.V.)	$s = \lambda \cdot \left(\frac{\beta_1 v_2 + \beta_2 \epsilon}{v_1(\epsilon + v_1)} \right)$	

Our Algorithm "SMART-ALLOC"

~6x fewer!

[US-MEDICARE NETWORK 2005]

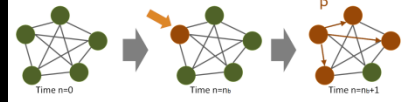
- Each circle is a hospital, ~3000 hospitals
- More than 30,000 patients transferred



CURRENT PRACTICE SMART-ALLOC

Main idea - SpikeM

1. Un-informed bloggers (uninformed about rumor)
2. External shock at time n_0 (e.g. breaking news)
3. Infection (word-of-mouth)

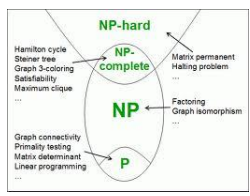


Infectiveness of a blog-post at age n : $f(n) = \beta * n^{-1.5}$

β - Strength of infection (quality of news)

$f(n)$ - Decay function (how infective a blog posting is)

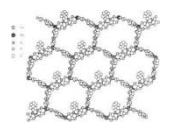
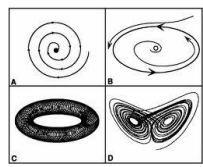
Power Law



Theory & Algo.

Biology

Physics



Comp. Systems

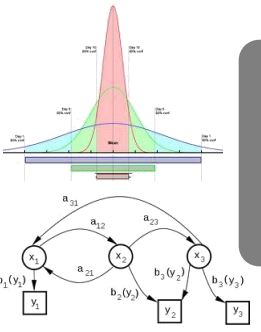


Social Science



Propagation on Networks

ML & Stats.



Econ.



References

1. *Winner-takes-all: Competing Viruses or Ideas on fair-play networks* (**B. Aditya Prakash**, Alex Beutel, Roni Rosenfeld, Christos Faloutsos) – In WWW 2012, Lyon
2. *Threshold Conditions for Arbitrary Cascade Models on Arbitrary Networks* (**B. Aditya Prakash**, Deepayan Chakrabarti, Michalis Faloutsos, Nicholas Valler, Christos Faloutsos) - In IEEE ICDM 2011, Vancouver (*Invited to KAIS Journal **Best Papers of ICDM.***)
3. *Times Series Clustering: Complex is Simpler!* (Lei Li, **B. Aditya Prakash**) - In ICML 2011, Bellevue
4. *Epidemic Spreading on Mobile Ad Hoc Networks: Determining the Tipping Point* (Nicholas Valler, **B. Aditya Prakash**, Hanghang Tong, Michalis Faloutsos and Christos Faloutsos) – In IEEE NETWORKING 2011, Valencia, Spain
5. *Formalizing the BGP stability problem: patterns and a chaotic model* (**B. Aditya Prakash**, Michalis Faloutsos and Christos Faloutsos) – In IEEE INFOCOM NetSciCom Workshop, 2011.
6. *On the Vulnerability of Large Graphs* (Hanghang Tong, **B. Aditya Prakash**, Tina Eliassi-Rad and Christos Faloutsos) – In IEEE ICDM 2010, Sydney, Australia
7. *Virus Propagation on Time-Varying Networks: Theory and Immunization Algorithms* (**B. Aditya Prakash**, Hanghang Tong, Nicholas Valler, Michalis Faloutsos and Christos Faloutsos) – In ECML-PKDD 2010, Barcelona, Spain
8. *MetricForensics: A Multi-Level Approach for Mining Volatile Graphs* (Keith Henderson, Tina Eliassi-Rad, Christos Faloutsos, Leman Akoglu, Lei Li, Koji Maruhashi, **B. Aditya Prakash** and Hanghang Tong) - In SIGKDD 2010, Washington D.C.
9. *Parsimonious Linear Fingerprinting for Time Series* (Lei Li, **B. Aditya Prakash** and Christos Faloutsos) - In VLDB 2010, Singapore
10. *EigenSpokes: Surprising Patterns and Scalable Community Chipping in Large Graphs* (**B. Aditya Prakash**, Ashwin Sridharan, Mukund Seshadri, Sridhar Machiraju and Christos Faloutsos) – In PAKDD 2010, Hyderabad, India
11. *BGP-lens: Patterns and Anomalies in Internet-Routing Updates* (**B. Aditya Prakash**, Nicholas Valler, David Andersen, Michalis Faloutsos and Christos Faloutsos) – In ACM SIGKDD 2009, Paris, France.
12. *Surprising Patterns and Scalable Community Detection in Large Graphs* (**B. Aditya Prakash**, Ashwin Sridharan, Mukund Seshadri, Sridhar Machiraju and Christos Faloutsos) – In IEEE ICDM Large Data Workshop 2009, Miami
13. *FRAPP: A Framework for high-Accuracy Privacy-Preserving Mining* (Shipra Agarwal, Jayant R. Haritsa and **B. Aditya Prakash**) – In Intl. Journal on Data Mining and Knowledge Discovery (DKMD), Springer, vol. 18, no. 1, February 2009, Ed: Johannes Gehrke.
14. *Complex Group-By Queries For XML* (C. Gokhale, N. Gupta, P. Kumar, L. V. S. Lakshmanan, R. Ng and **B. Aditya Prakash**) – In IEEE ICDE 2007, Istanbul, Turkey.

Acknowledgements

Collaborators

Christos Faloutsos 

Roni Rosenfeld, 

Michalis Faloutsos, 

Lada Adamic, 

Theodore Iwashyna (M.D.), 

Dave Andersen, 

Tina Eliassi-Rad, 

Iulian Neamtiu, 

Varun Gupta, 
The University of Chicago Booth School of Business

Jilles Vreeken, 

Deepayan Chakrabarti, 

Hanghang Tong, 


Kunal Punera, 

Ashwin Sridharan, Sprint 

Sridhar Machiraju, 

Mukund Seshadri, 

Alice Zheng, **Microsoft**

Lei Li, 

Polo Chau, 

Nicholas Valler, 

Alex Beutel, 

Xuetao Wei 

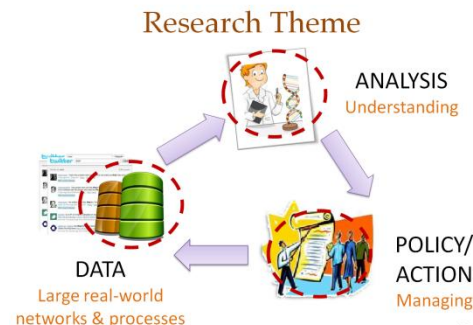
Acknowledgements

Funding



Dynamical Processes on Large Networks

B. Aditya Prakash
Christos Faloutsos



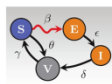
Analysis

Policy/Action

Data

Our thresholds for some models

- $s = \text{effective strength}$
- $s < 1$: below threshold

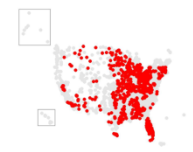


Models	Effective Strength (s)	Threshold (tipping point)
SIS, SIR, SIRS, SEIR	$s = \lambda \cdot \left(\frac{\beta}{\delta} \right)$	$s = 1$
SIV, SEIV	$s = \lambda \cdot \left(\frac{\beta\gamma}{\delta(\gamma + \theta)} \right)$	
SI ₁ I ₂ V ₁ V ₂ (H.I.V.)	$s = \lambda \cdot \left(\frac{\beta_1 v_2 + \beta_2 \epsilon}{v_2(\epsilon + v_1)} \right)$	

Our Algorithm "SMART-ALLOC"

~6x fewer!

- [US-MEDICARE NETWORK 2005]
- Each circle is a hospital, ~3000 hospitals
 - More than 30,000 patients transferred

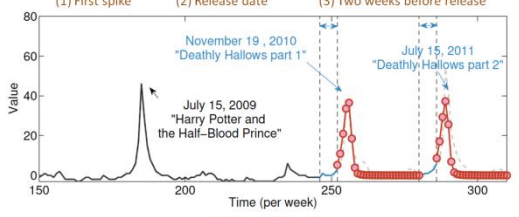


CURRENT PRACTICE

SMART-ALLOC

"What-if" forecasting

-SpikeM can forecast not only tail-part, but also rise-part!



- SpikeM can forecast upcoming spikes