



Mining and Forecasting of Big Time-series Data

Yasushi Sakurai (Kumamoto University)

Yasuko Matsubara (Kumamoto University)

Christos Faloutsos (Carnegie Mellon University)



Roadmap



- Motivation
- Similarity search, pattern discovery and summarization
- **Non-linear modeling and forecasting**
- Extension of time-series data: tensor analysis

Part 1

Part 2

Part 3





Part 2

Roadmap



Problem

– Why: “non-linear” modeling

Fundamentals

– Non-linear (“gray-box”) models

Applications

– Epidemics



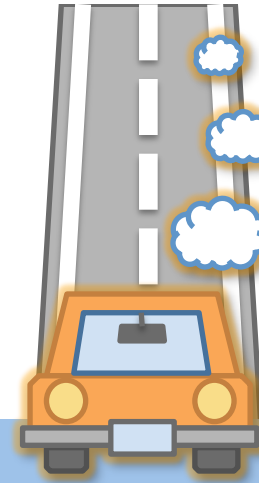
– Information diffusion



– (Online) competition



vs.



Non-linear mining and forecasting

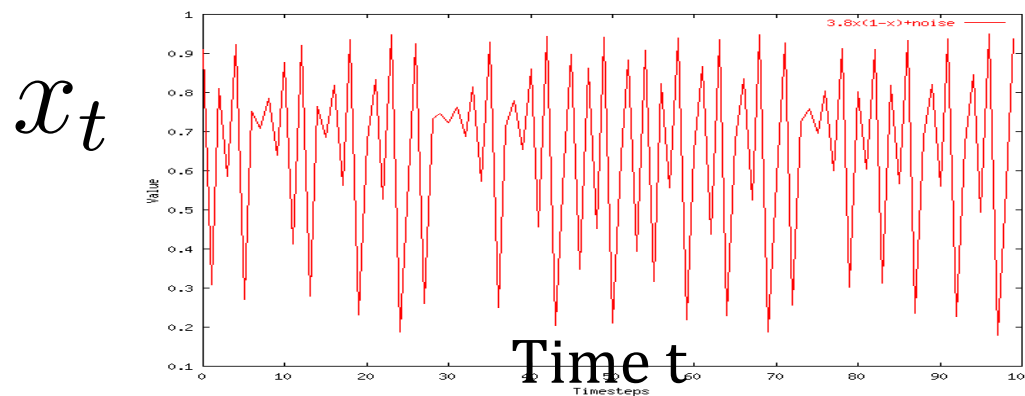
Q. What are “non-linear phenomena”?

Example: logistic parabola

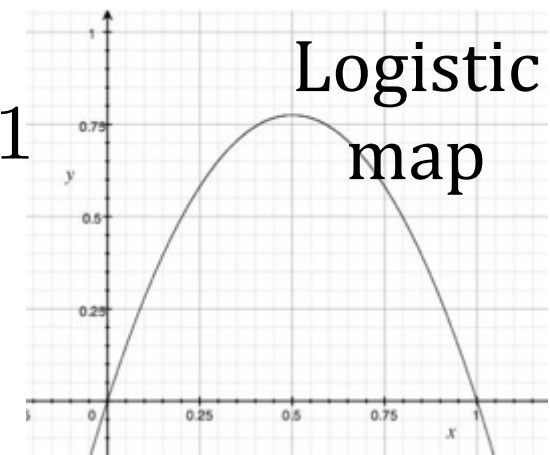
Models population of flies [R. May/1976]

$$x_{t+1} = ax_t \cdot (1 - x_t)$$

Time-series plot



x_{t+1}



x_t

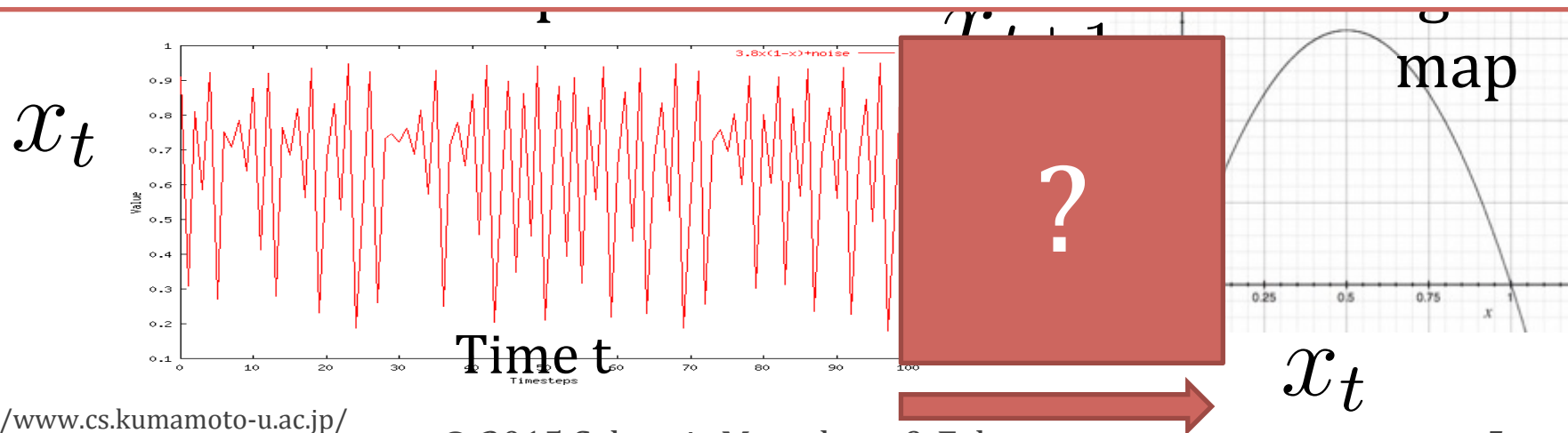
Non-linear mining and forecasting

Q. What are “non-linear phenomena”?

Problem:

Given: a time series x_t

Predict: its future course, i.e., x_{t+1}, x_{t+2}, \dots

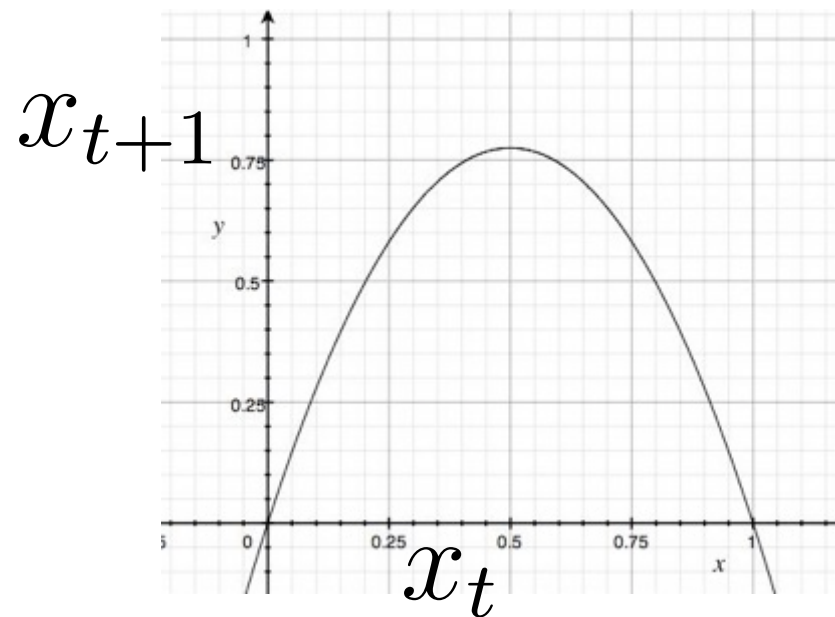




How to forecast?

Solution 1

Linear equations, e.g., AR, ARIMA, ...





How to forecast?

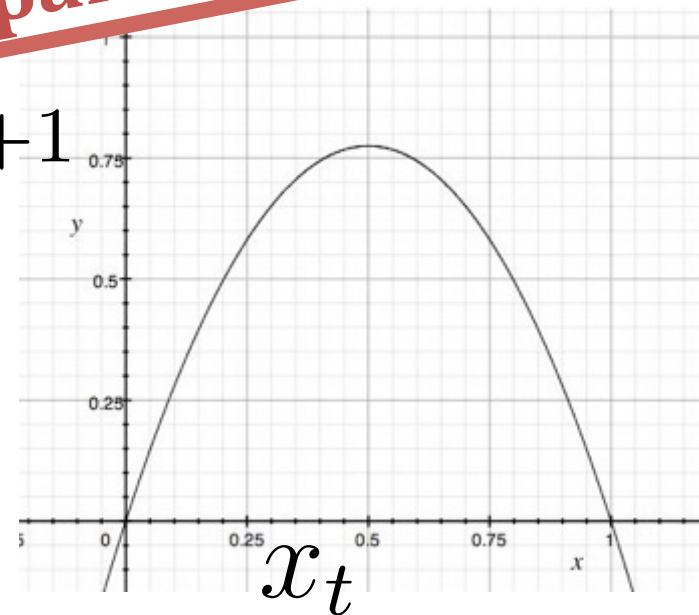
Solution 1

Linear equations, e.g., AR, ARIMA, ...

Details @ part1

e.g., AR(1)

$$x_{t+1} = ax_t + \epsilon$$





How to forecast?

Solution 1

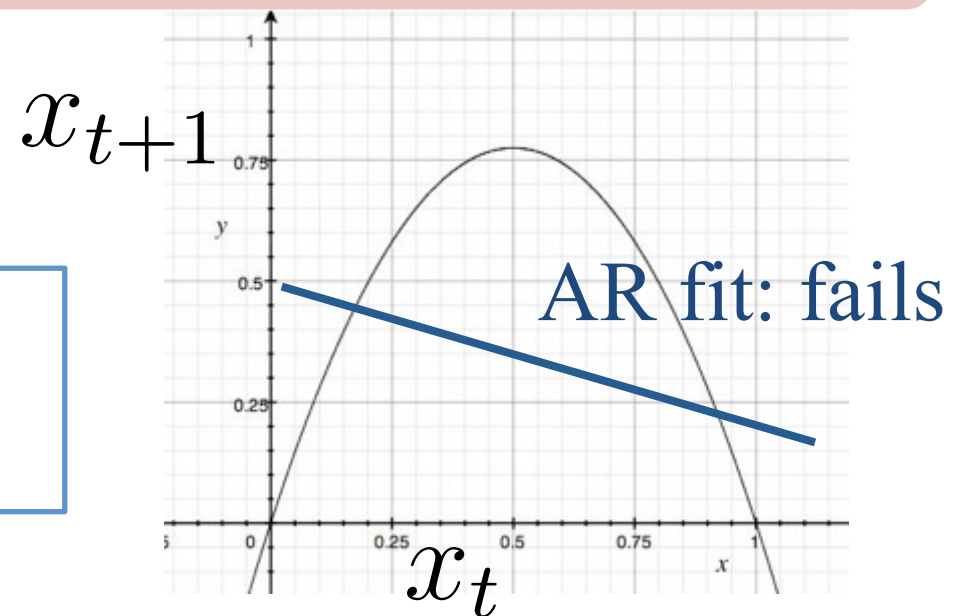
Linear equations, e.g., AR, ARIMA, ...



but: linearity assumption

e.g., AR(1)

$$x_{t+1} = ax_t + \epsilon$$





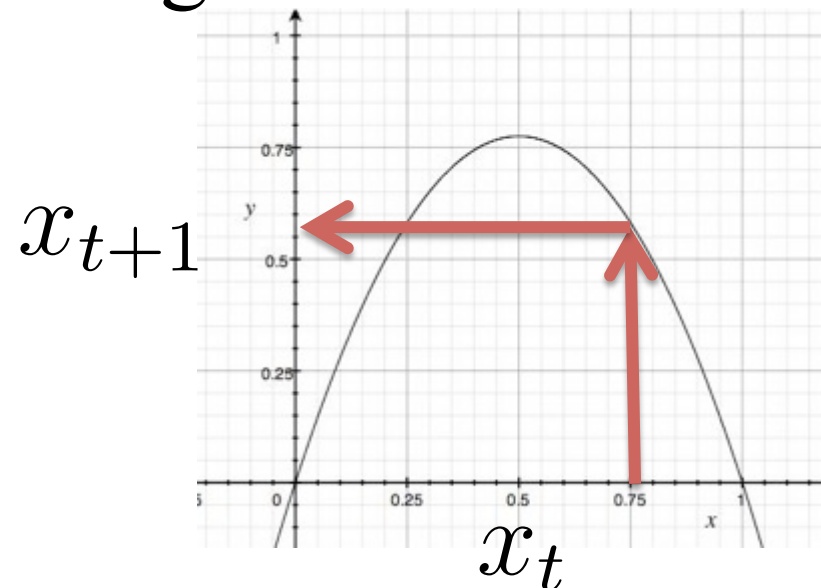
How to forecast?

Solution 2

“Delayed Coordinate Embedding”

= Lag Plots [Sauer92]

- Based on k-nearest neighbor search

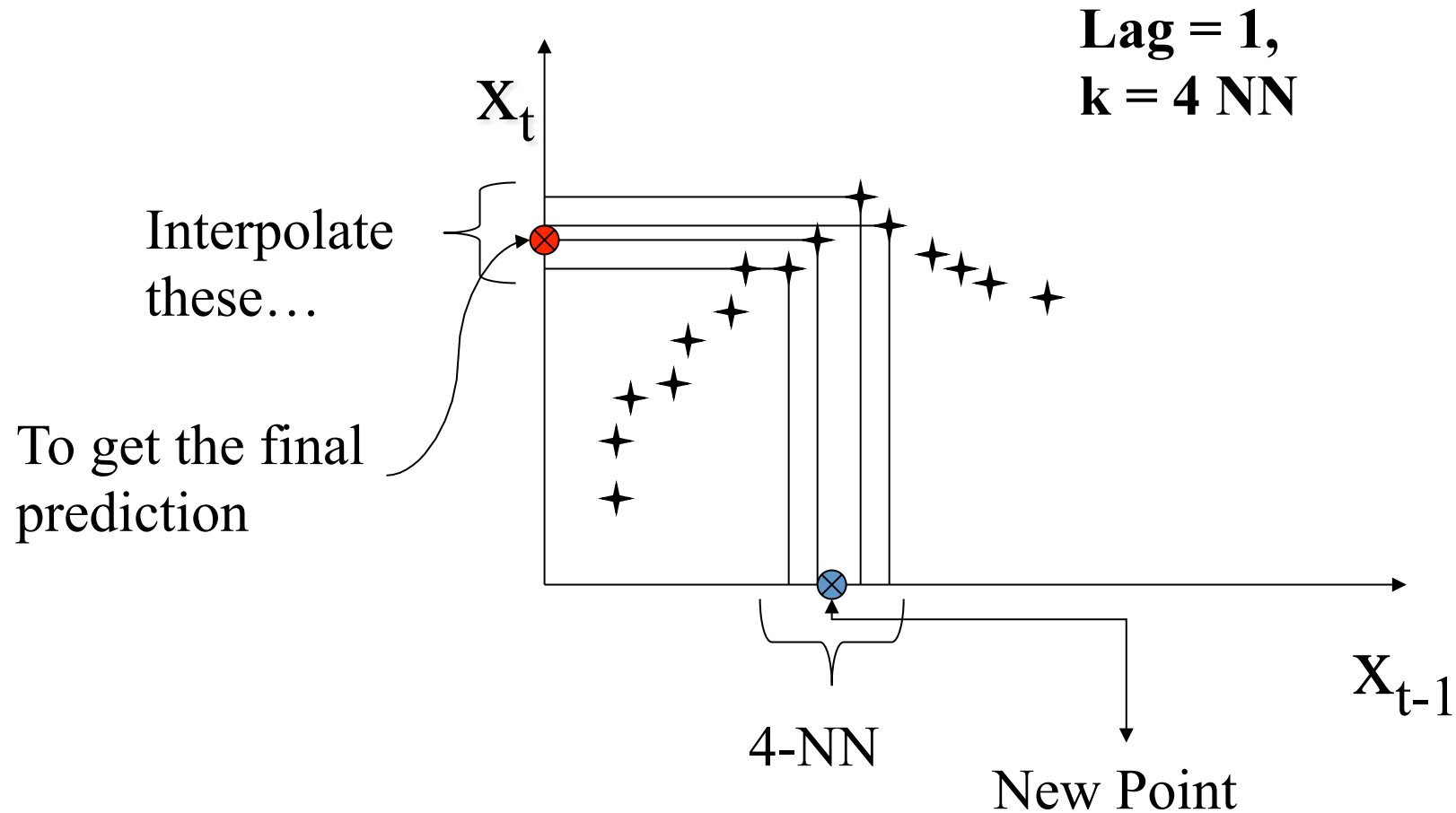




General Intuition (Lag Plot)



Solution 2





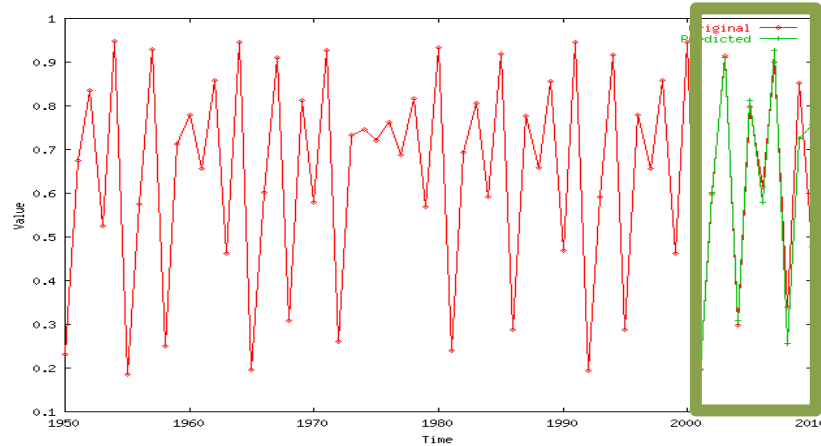
Forecasting results (Lag Plot)



[Chakrabarti+ CIKM'02]

Solution 2

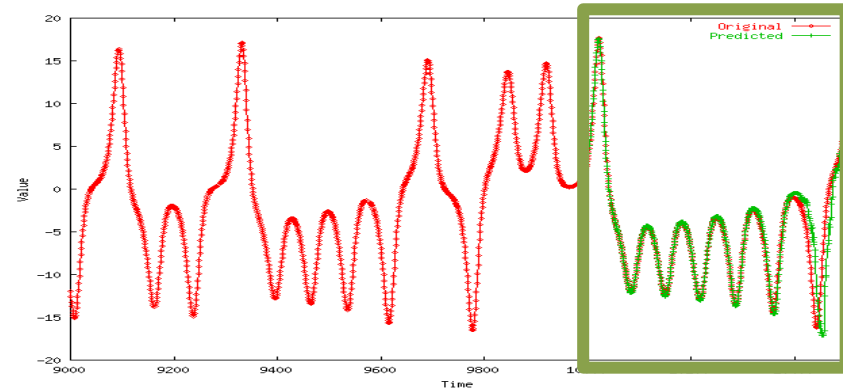
Logistic parabola



Original x_t
(red)

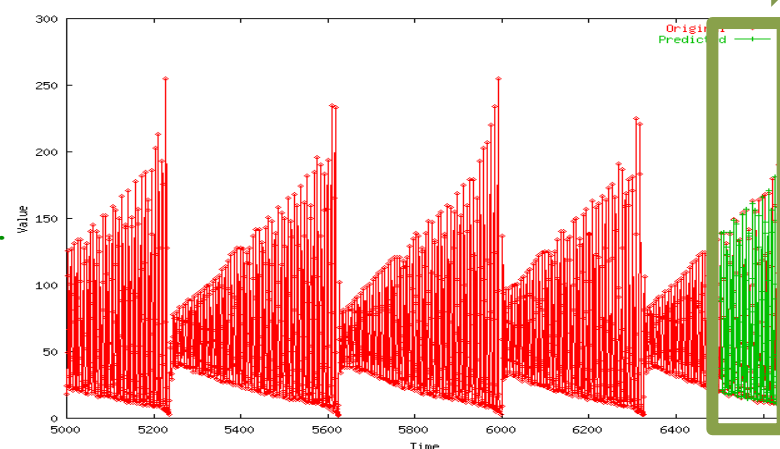
Forecasted $x_{t+1, \dots}$
(green)

LORENZ



Laser

Forecast





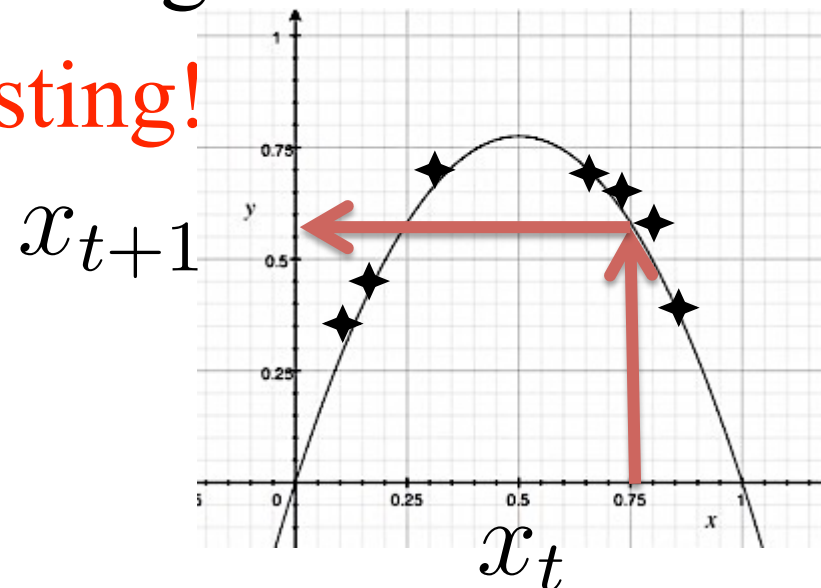
How to forecast?

Solution 2

“Delayed Coordinate Embedding”

= Lag Plots [Sauer92]

- Based on k-nearest neighbor search
- **Non-linear Forecasting!**



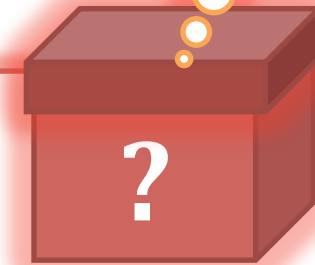


How to forecast?

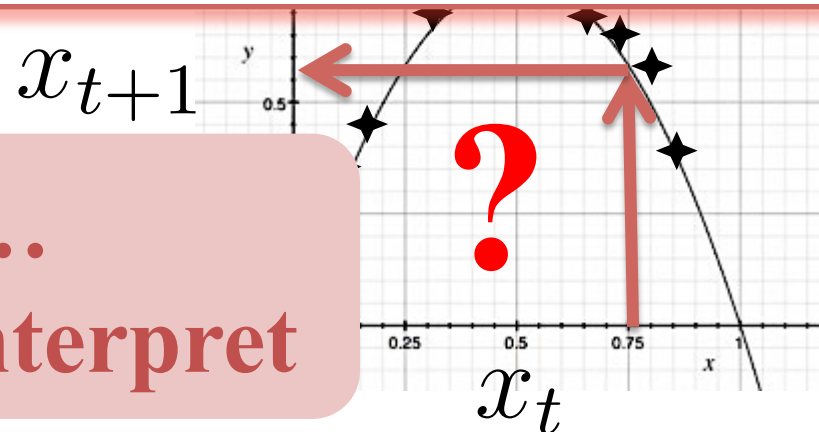
Solution 2

“Delayed Coordinate Embedding”

“Black-box” mining
(we don't know the equations)



But, still, ...
Hard to interpret

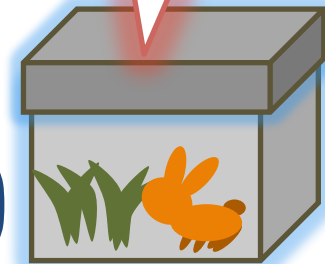
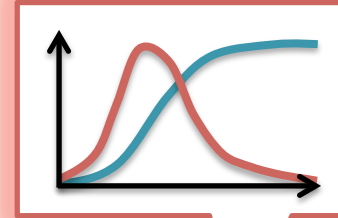




How to forecast?

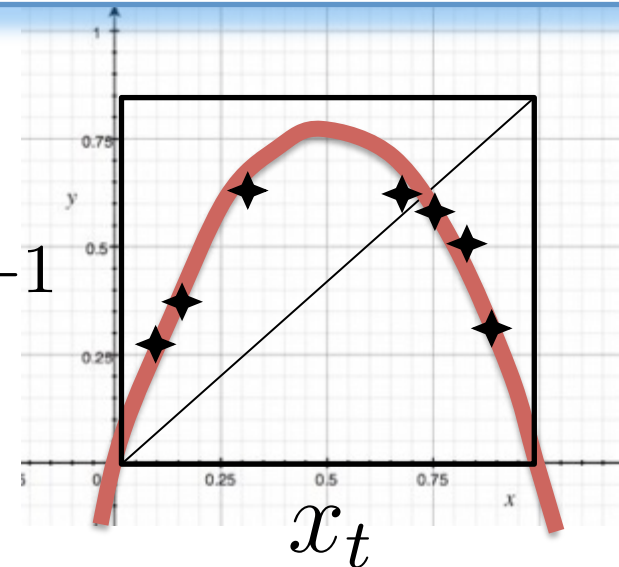
Solution 3

“Gray-box” mining
(if we know the equations)



Non-linear
modeling!

$$x_{t+1} = ax_t \cdot (1 - x_t)$$

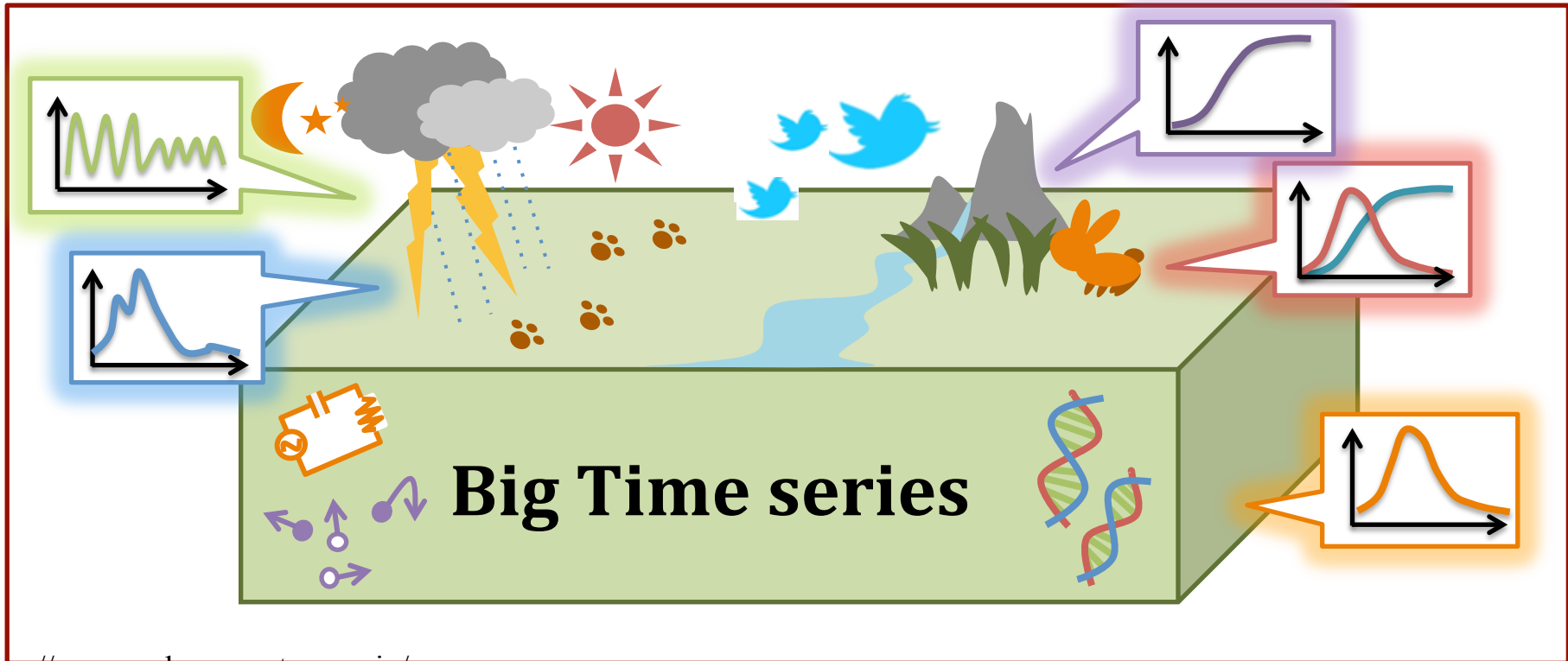
 x_{t+1}

 x_t



How to forecast?

Solution 3

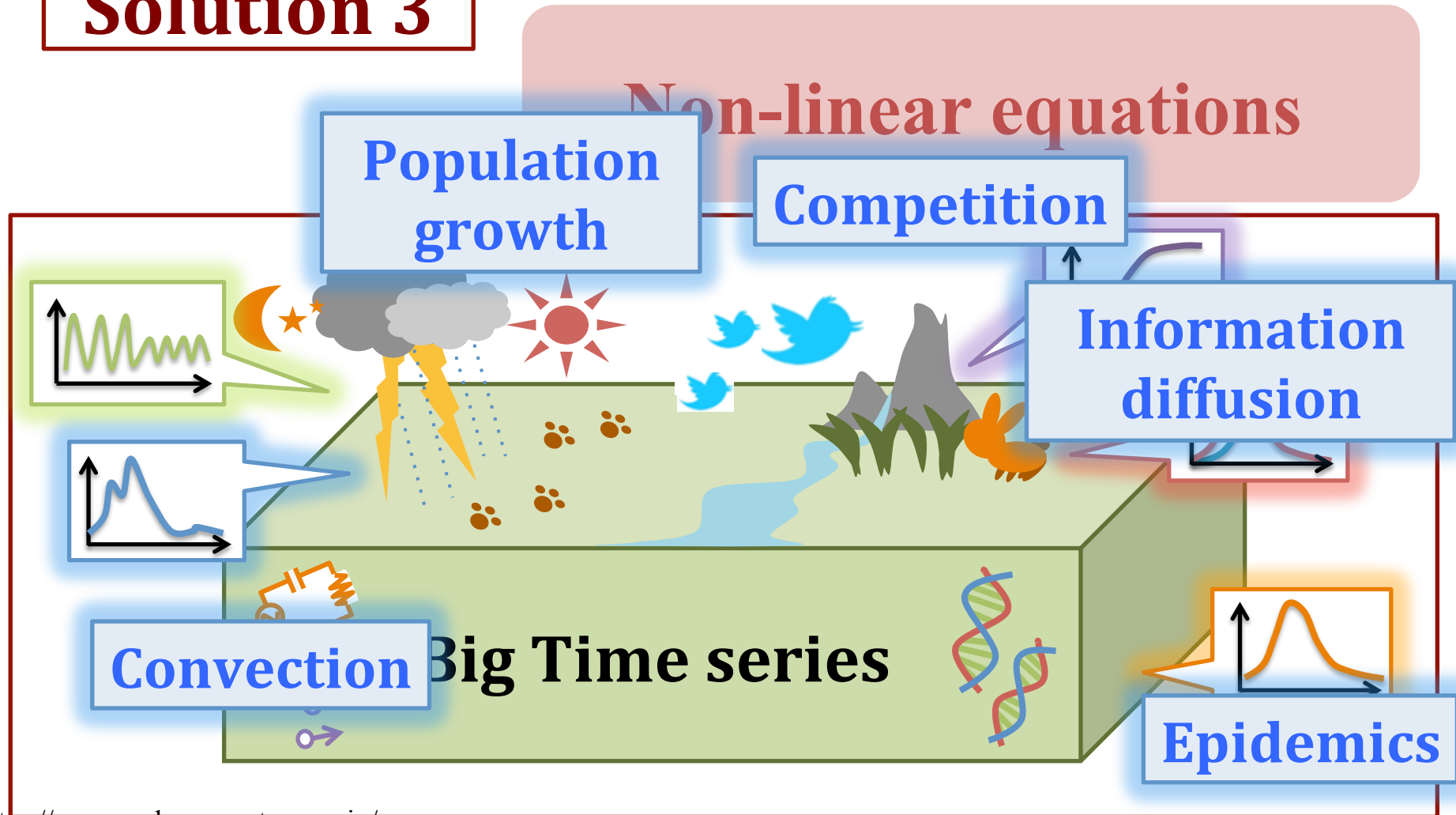
Non-linear equations





How to forecast?

Solution 3





Part 2

Roadmap



Problem

✓ Why: “non-linear” modeling

Fundamentals

– Non-linear (grey-box) models

Applications

– Epidemics



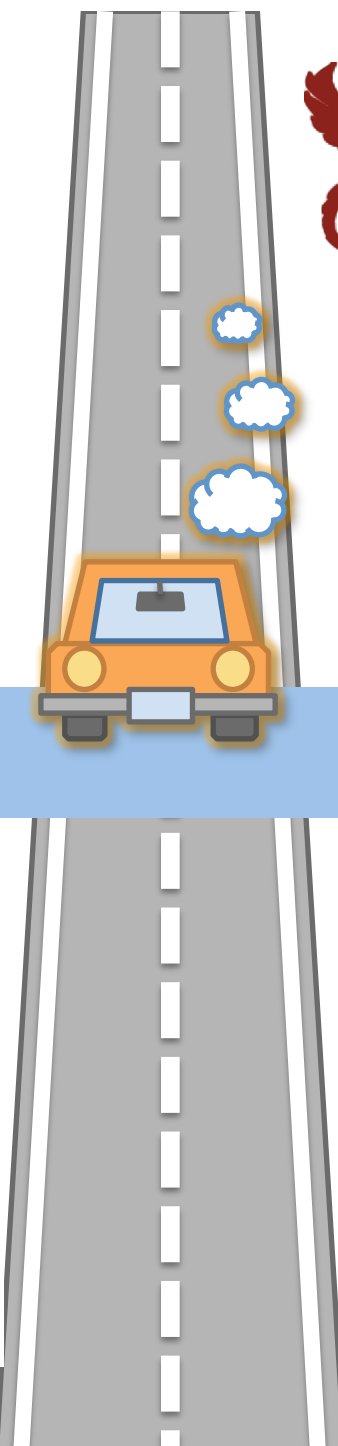
– Information diffusion



– (Online) competition



vs.





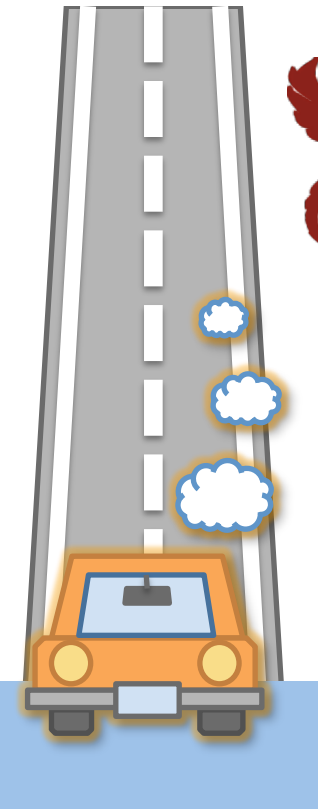
Problem

✓ Why: “non-linear” modeling

Fundamentals

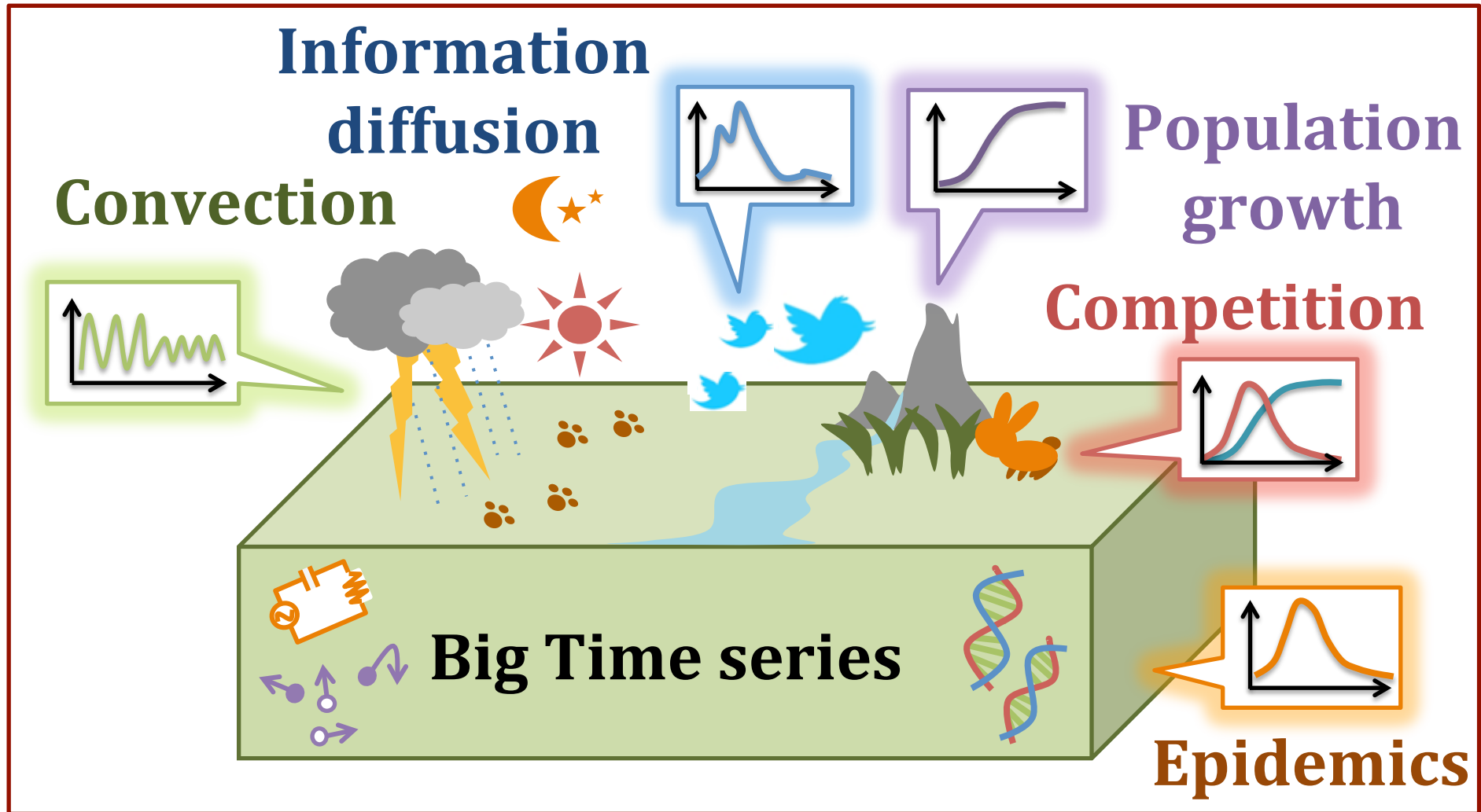
– Non-linear (grey-box) models

- Logistic function
- Lotka-Volterra (prey-predator, competition)
- SI, SIR models, etc.
- Lorenz equations, etc.

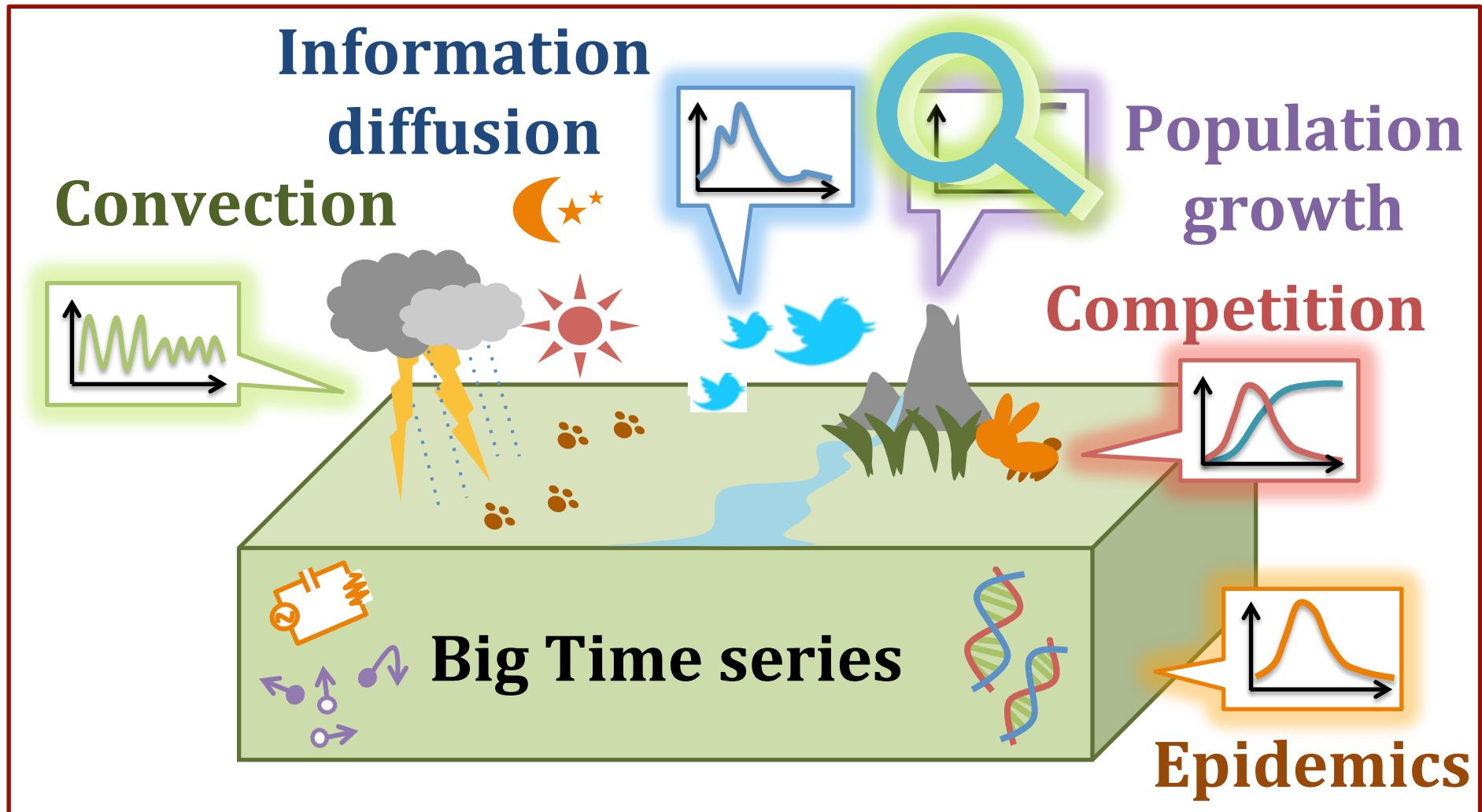




Grey-box mining and non-linear equations



Grey-box mining and non-linear equations

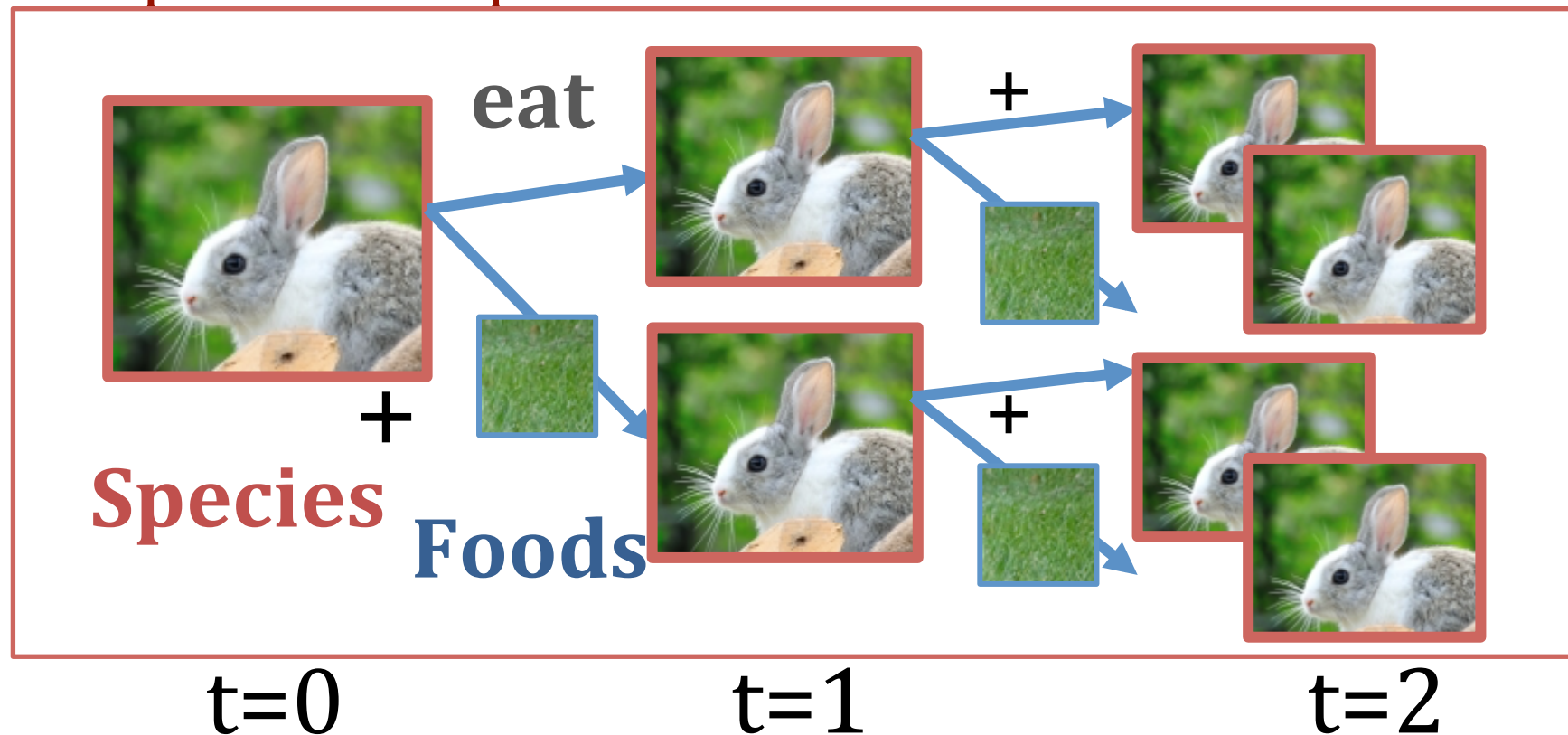




Logistic function

So-called “Verhulst” model (=sigmoid, =Bass)

- Population expansion with limited resources





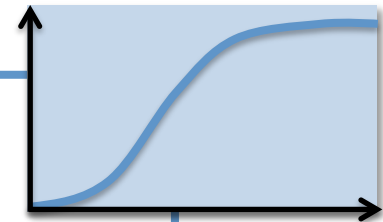
Logistic function

So-called “Verhulst” model (=sigmoid, =Bass)

- Population expansion with limited resources

P : Population size

$$\frac{dP}{dt} = rP \left(1 - \frac{P}{K} \right)$$



p – Initial condition (i.e., $P(0) = p$)

r – Growth rate, reproductively

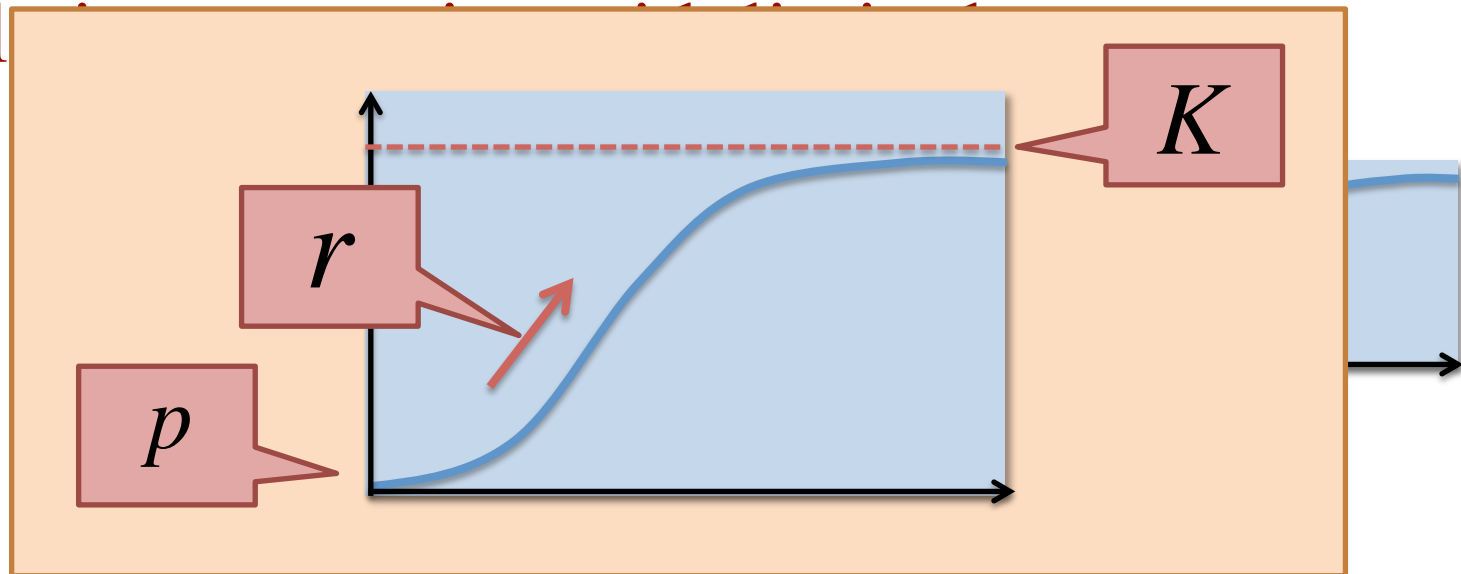
K – Carrying capacity (=available resources)



Logistic function

So-called “Verhulst” model (=sigmoid, =Bass)

- Popul



p – Initial condition (i.e., $P(0) = p$)

r – Growth rate, reproductively

K – Carrying capacity (=available resources)



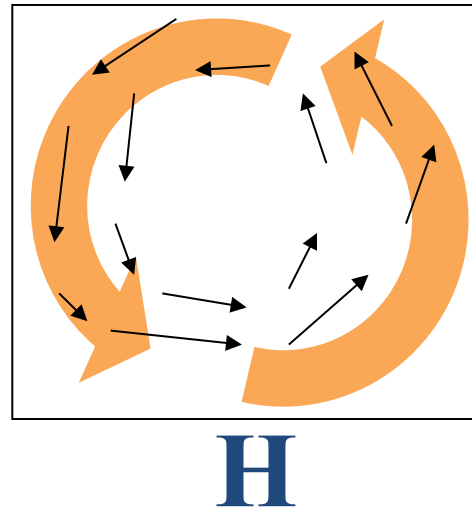
Lotka-Volterra equations



So-called “prey-predator” model



Prey (H)



Predator (P)

- **H : count of prey (e.g., hare)**
- **P : count of predators (e.g., lynx)**



Lotka-Volterra equations



So-called “prey-predator” model



Prey (H)

$$\frac{dH}{dt} = rH - aHP$$

$$\frac{dP}{dt} = bHP - mP$$



Predator (P)

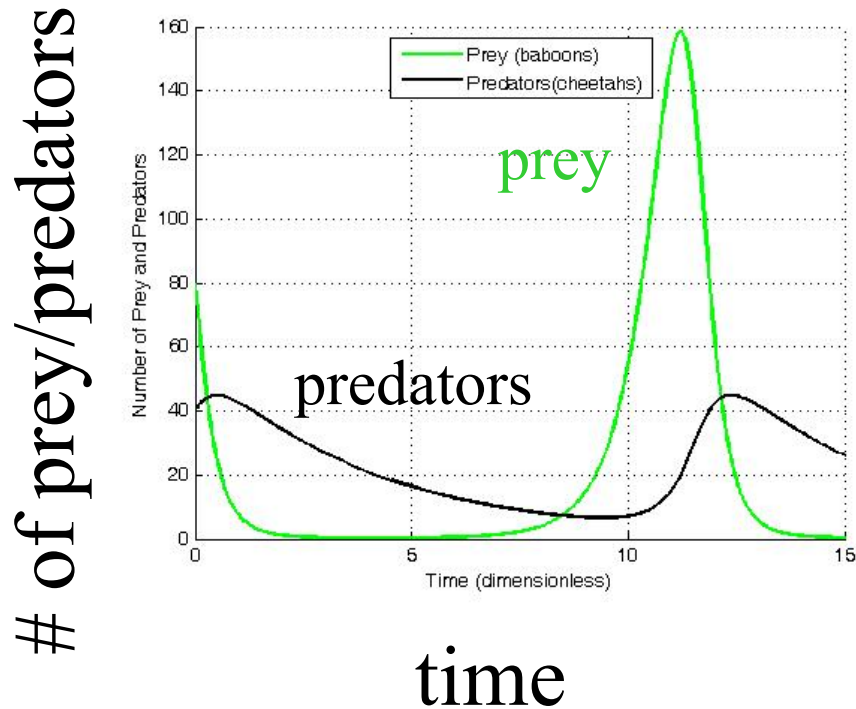
- **H : count of prey (e.g., hare)**
- **P : count of predators (e.g., lynx)**



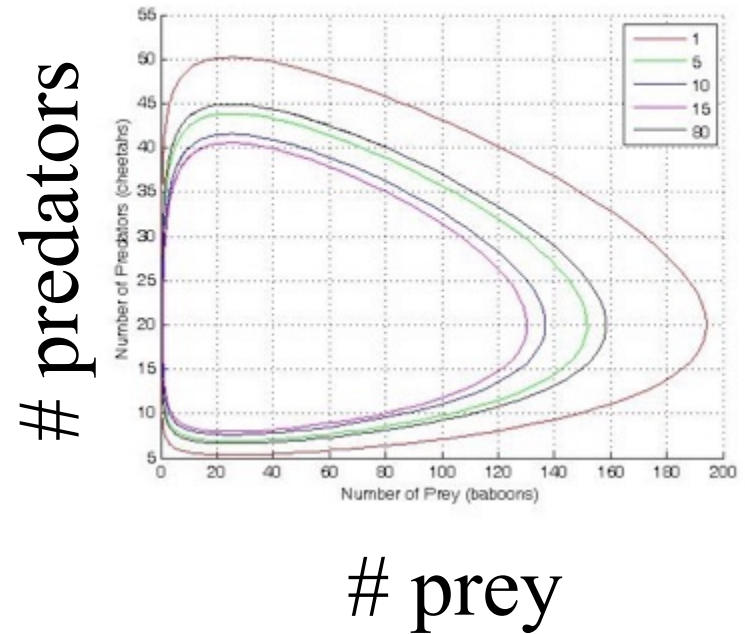
Solution to the Lotka-Volterra equations.



Frequency Plot



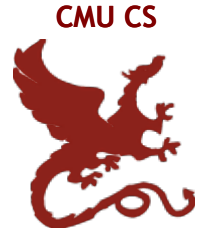
Phase Space Plot



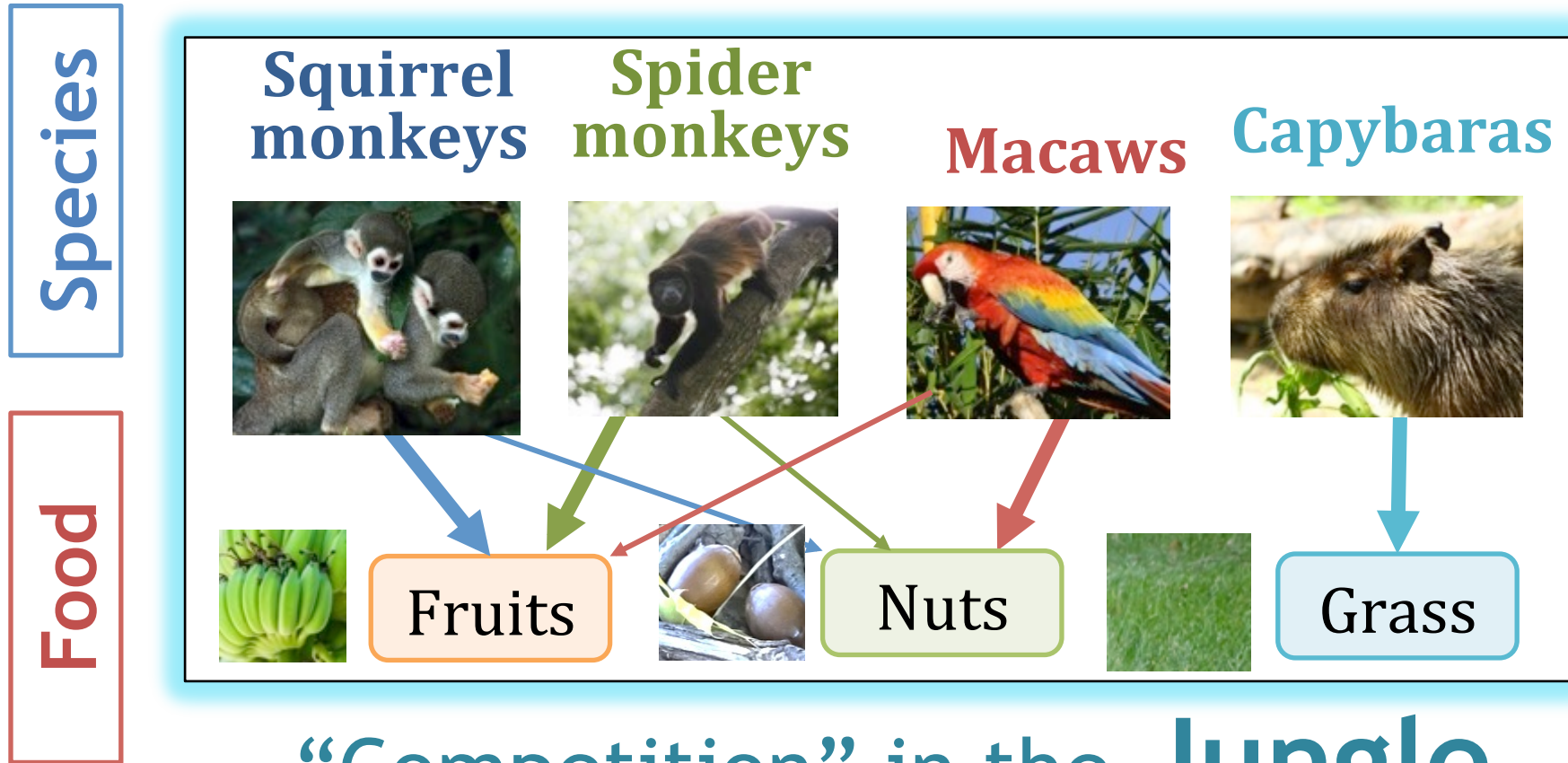
From Wikipedia



Extension: “Competitive” Lotka-Volterra equations



Competition between multiple (d) species



“Competition” in the Jungle

Image courtesy of Tina Phillips and amenic181 at FreeDigitalPhotos.net.

“Competitive”



Lotka-Volterra equations

Competition between multiple (d) species

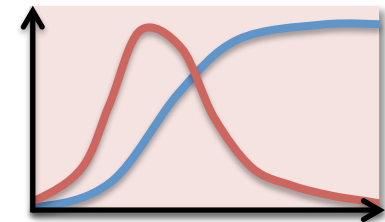
Population of species i

Population of j

$$\frac{dP_i}{dt} = r_i P_i \left(1 - \frac{\sum_{j=1}^d a_{ij} P_j}{K_i} \right)$$

$(i = 1, \dots, d)$

a_{ij} : Interaction coefficient
i.e., effect rate of species j on i



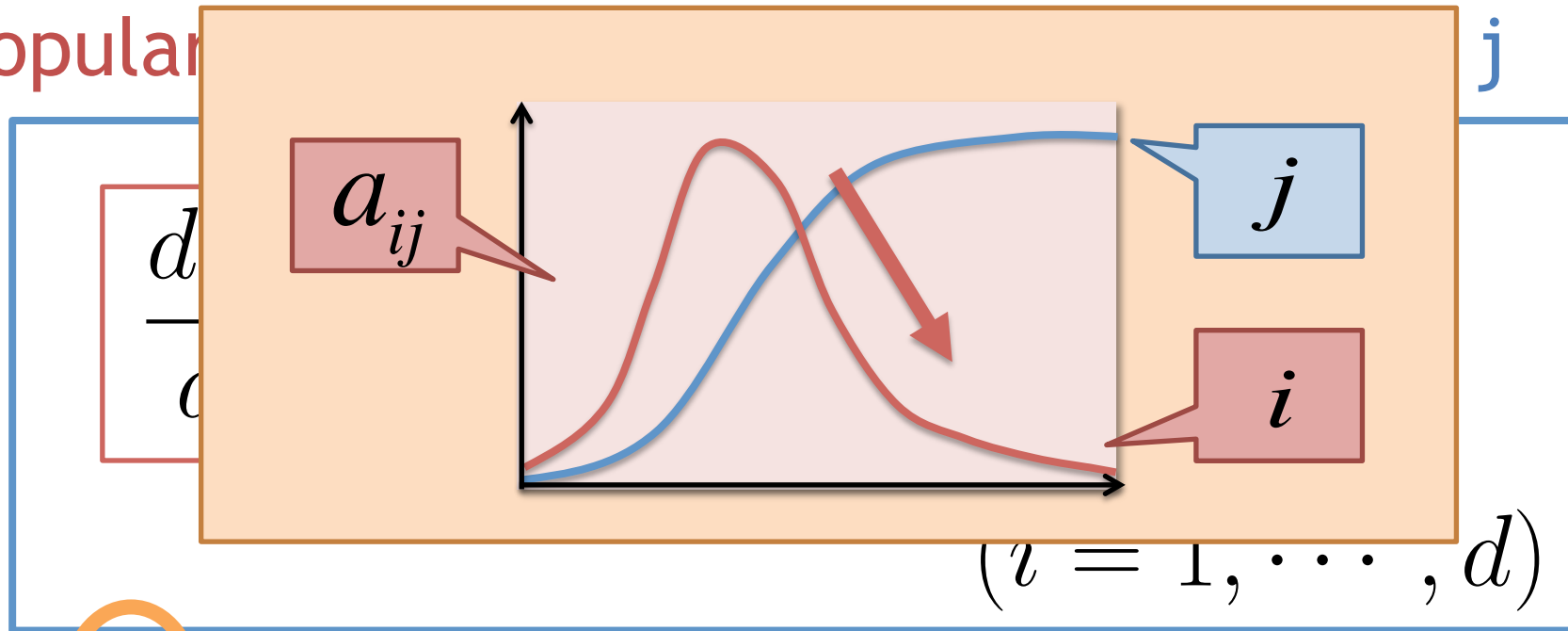
“Competitive”



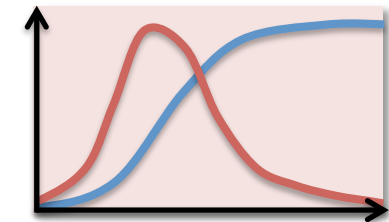
Lotka-Volterra equations

Competition between multiple (d) species

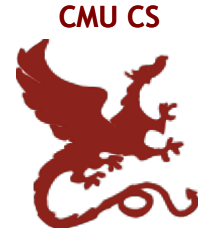
Popular



a_{ij} : Interaction coefficient
i.e., effect rate of species j on i



“Competitive”



Lotka-Volterra equations

- Biological interaction

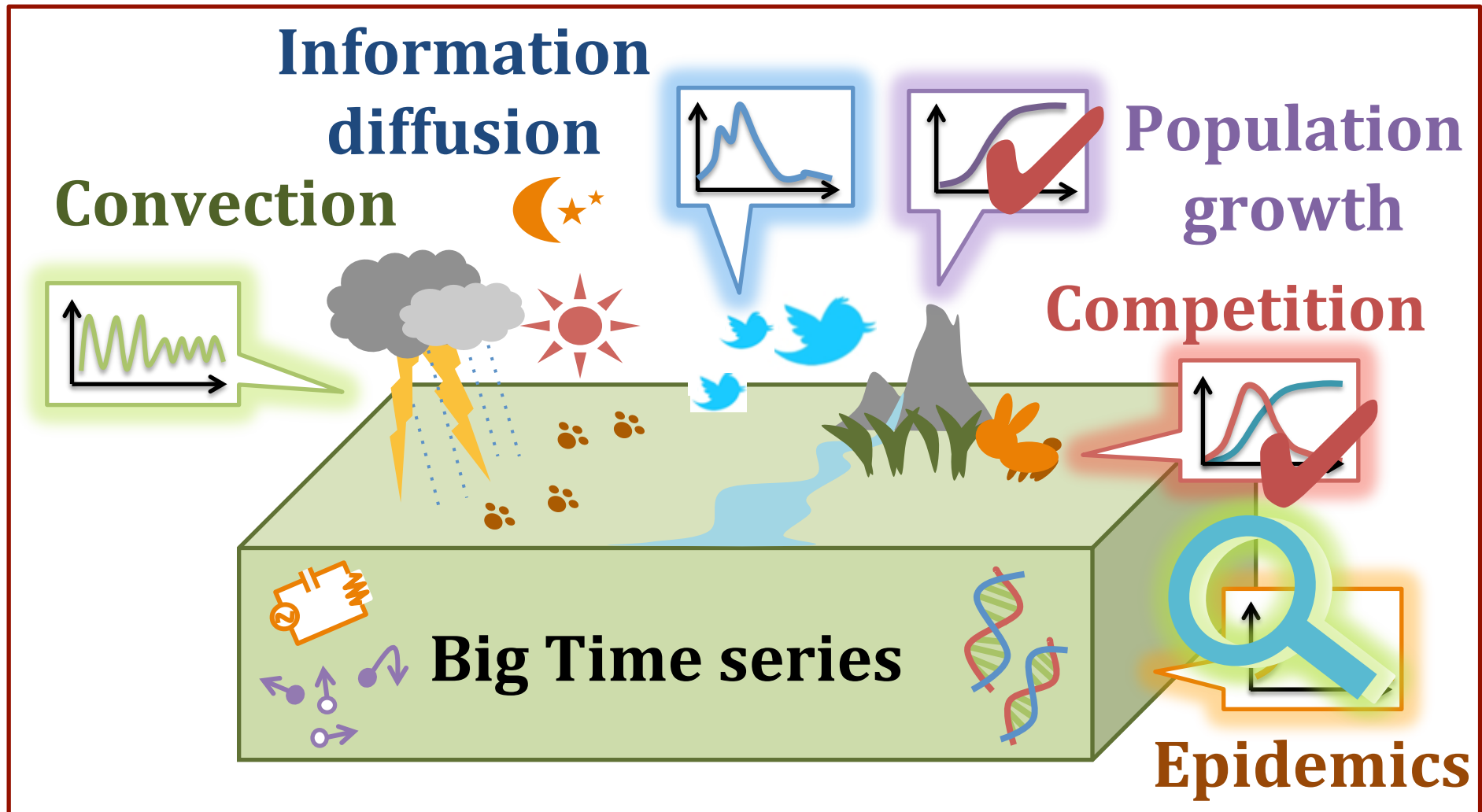
– Table: Type of interaction

0 : no effect
 - : detrimental
 + : beneficial

		Species B		
		+	0	-
Species A	+	Mutualism		
	0	Commensalism	Neutralism	
	-	Antagonism	Amensalism	Competition



Grey-box mining and non-linear equations

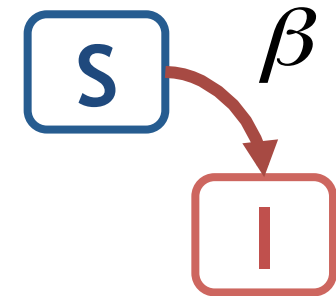




Epidemics: Susceptible-Infected (SI) model



Each node is in one of two states

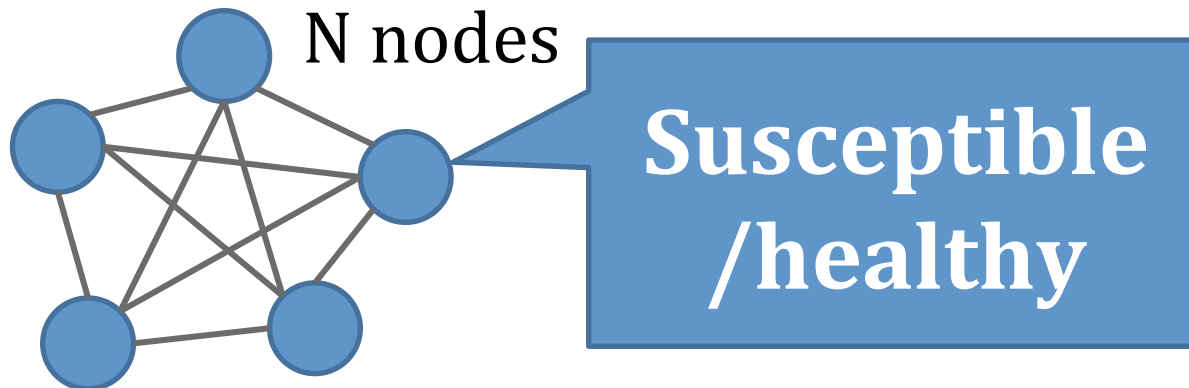
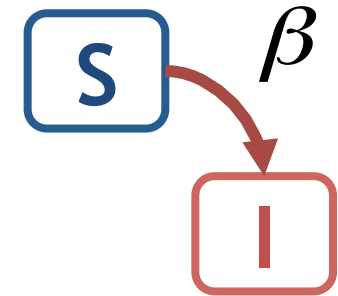




Epidemics: Susceptible-Infected (SI) model



Each node is in one of two states



Time $t=0$



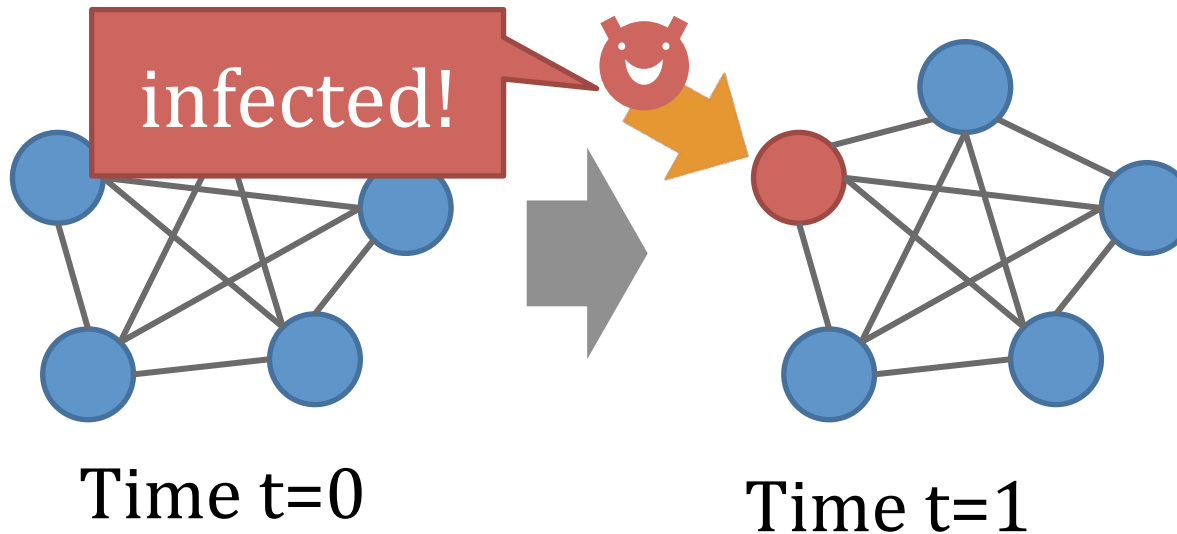
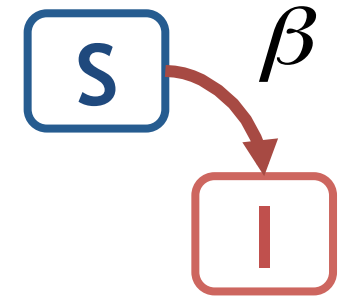
Epidemics: Susceptible-Infected (SI) model



Each node is in one of two states

S – Susceptible (healthy)

I – Infected





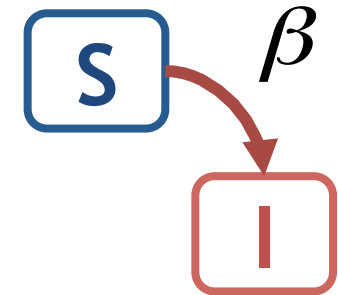
Epidemics: Susceptible-Infected (SI) model



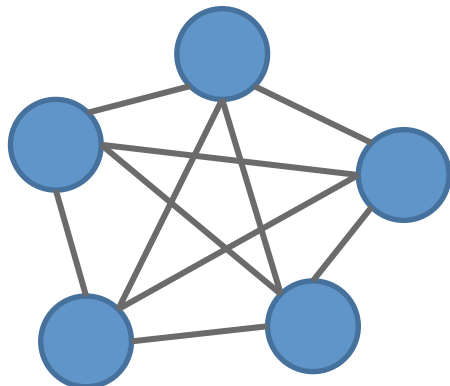
Each node is in one of two states

S – Susceptible (healthy)

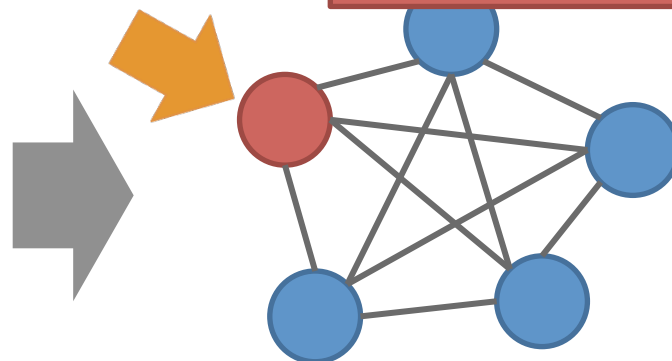
I – Infected β : infection rate



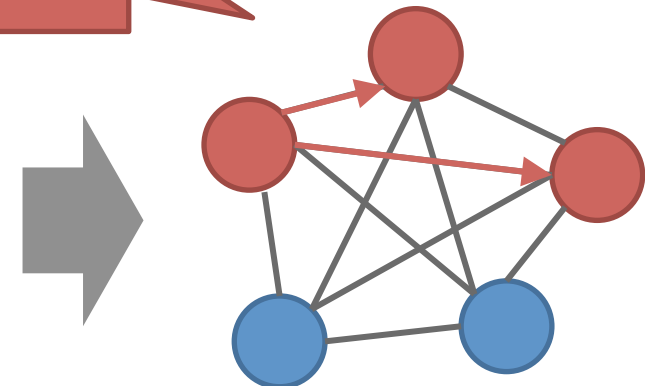
Prob. β



Time $t=0$



Time $t=1$



Time $t=2$



Epidemics: Susceptible-Infected (SI) model



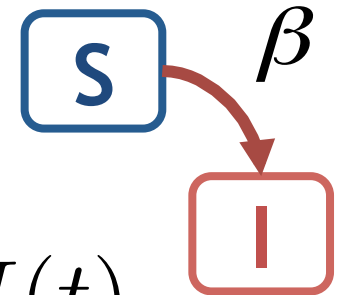
Each node is in one of two states

$$\frac{dS}{dt} = -\beta SI$$

$$\frac{dI}{dt} = +\beta SI$$

$$N = S(t) + I(t)$$

β : Infection strength
 N : Population size



i.e.,
$$\frac{dI}{dt} = \beta(N - I)I$$

Epidemics: Susceptible-Infected (SI) model



Each node is in one of two states

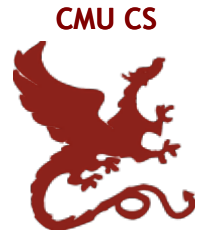
$\frac{dP}{dt} = rP(1 - \frac{P}{K})$

$\frac{dI}{dt} = \beta N \cdot I(1 - \frac{I}{N})$

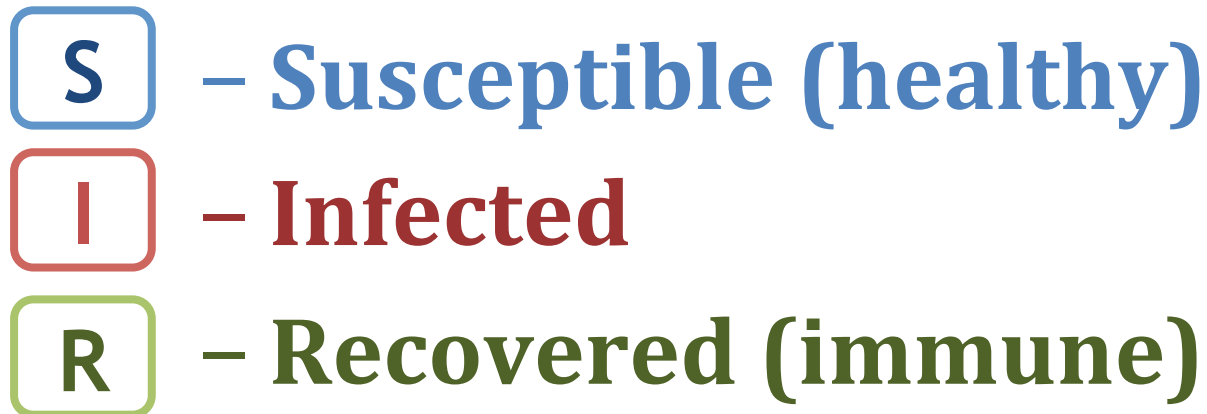
i.e.,
$$\frac{dI}{dt} = \beta(N - I)I$$



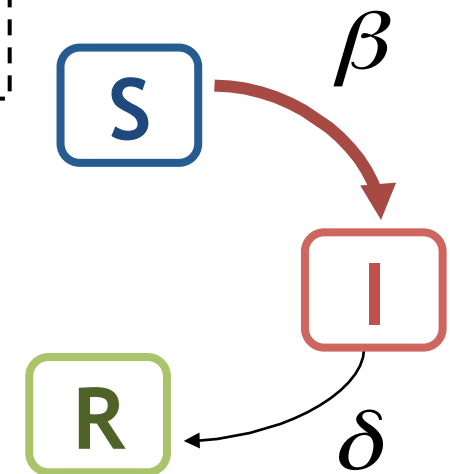
Susceptible-Infected-recovered (SIR) model



Recovered with immunity



β : Infection rate
 δ : Recovery rate

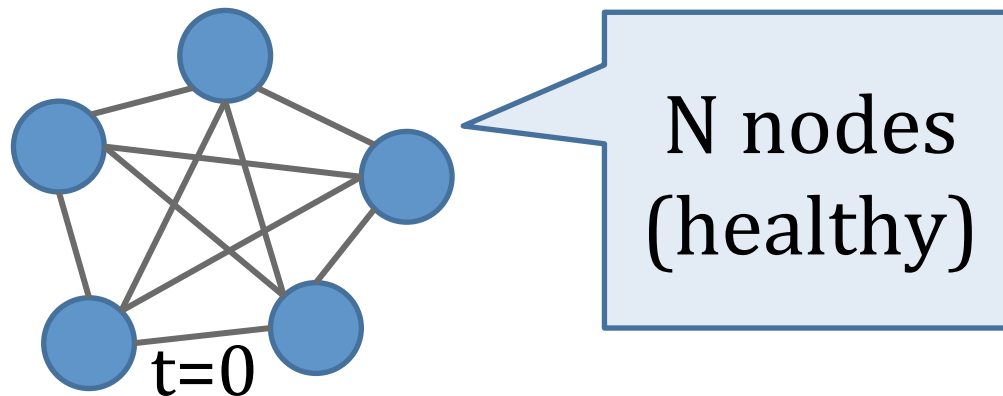


Details

Susceptible-Infected-recovered (SIR) model



Recovered with immunity

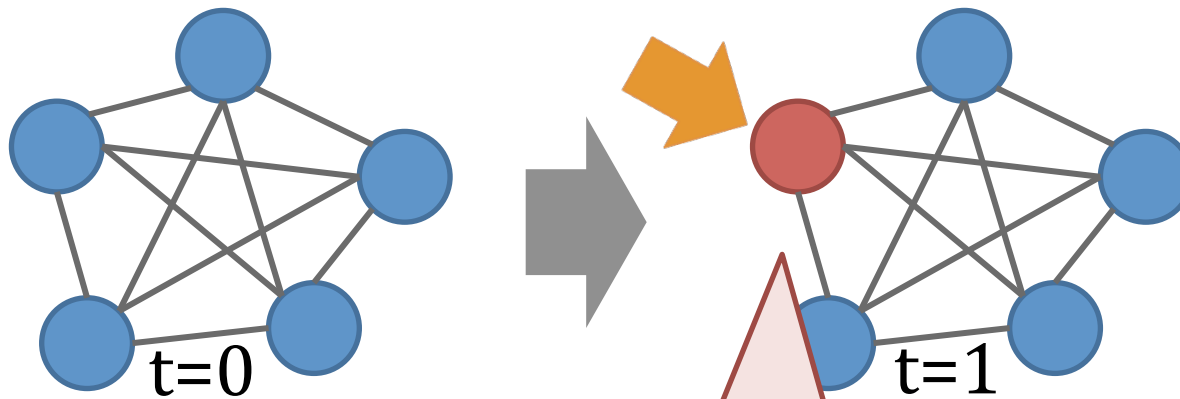



Details

Susceptible-Infected-recovered (SIR) model



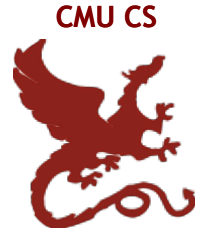
Recovered with immunity



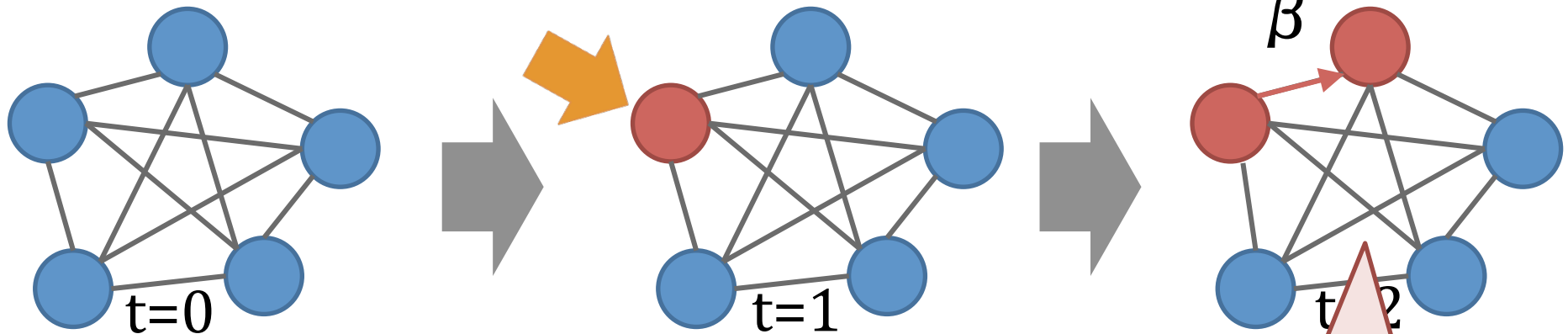
infection 


Details

Susceptible-Infected-recovered (SIR) model



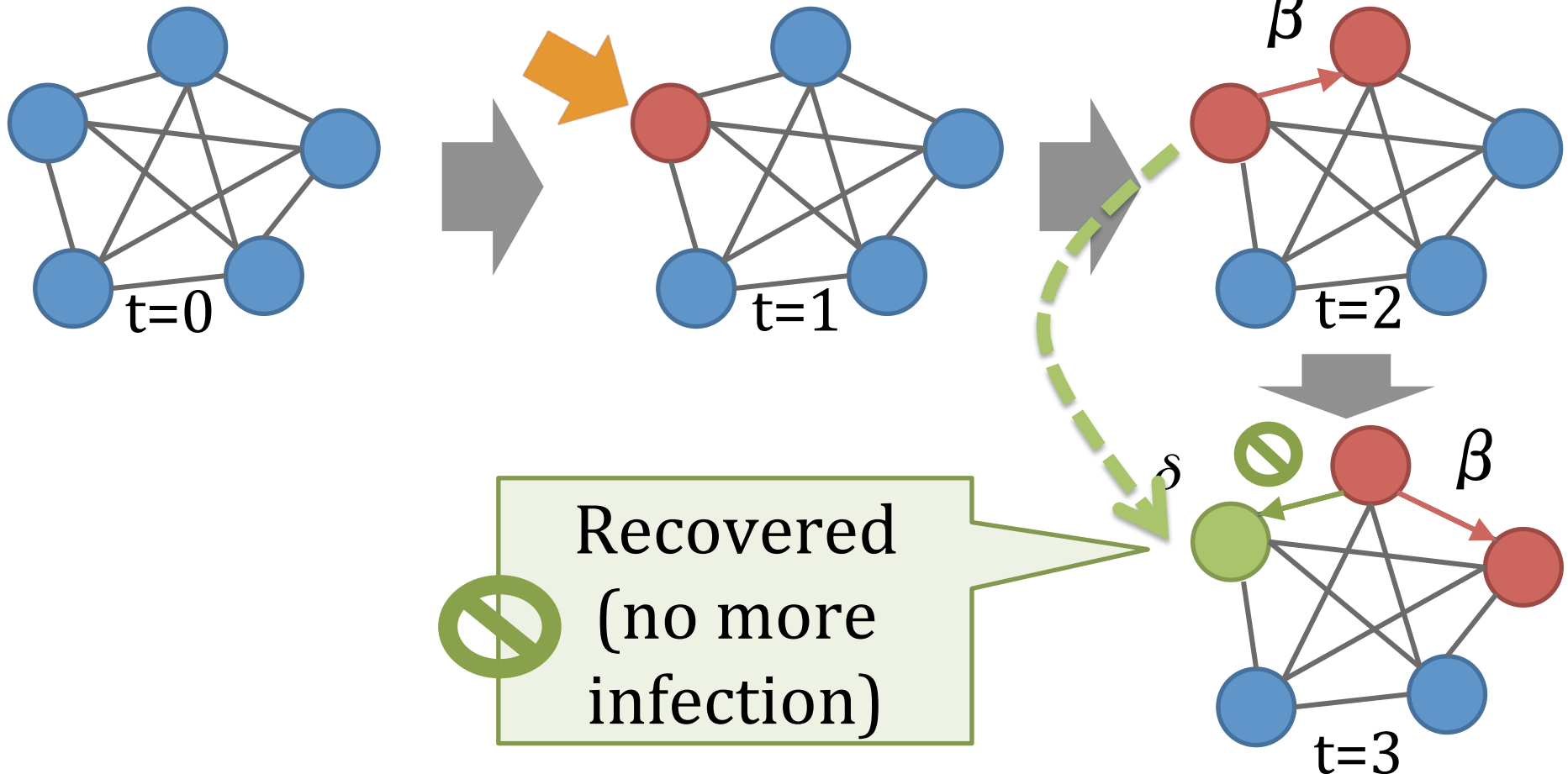
Recovered with immunity



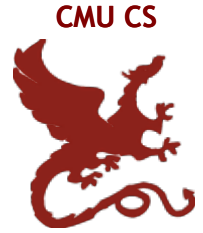
Propagation 

Susceptible-Infected-recovered (SIR) model

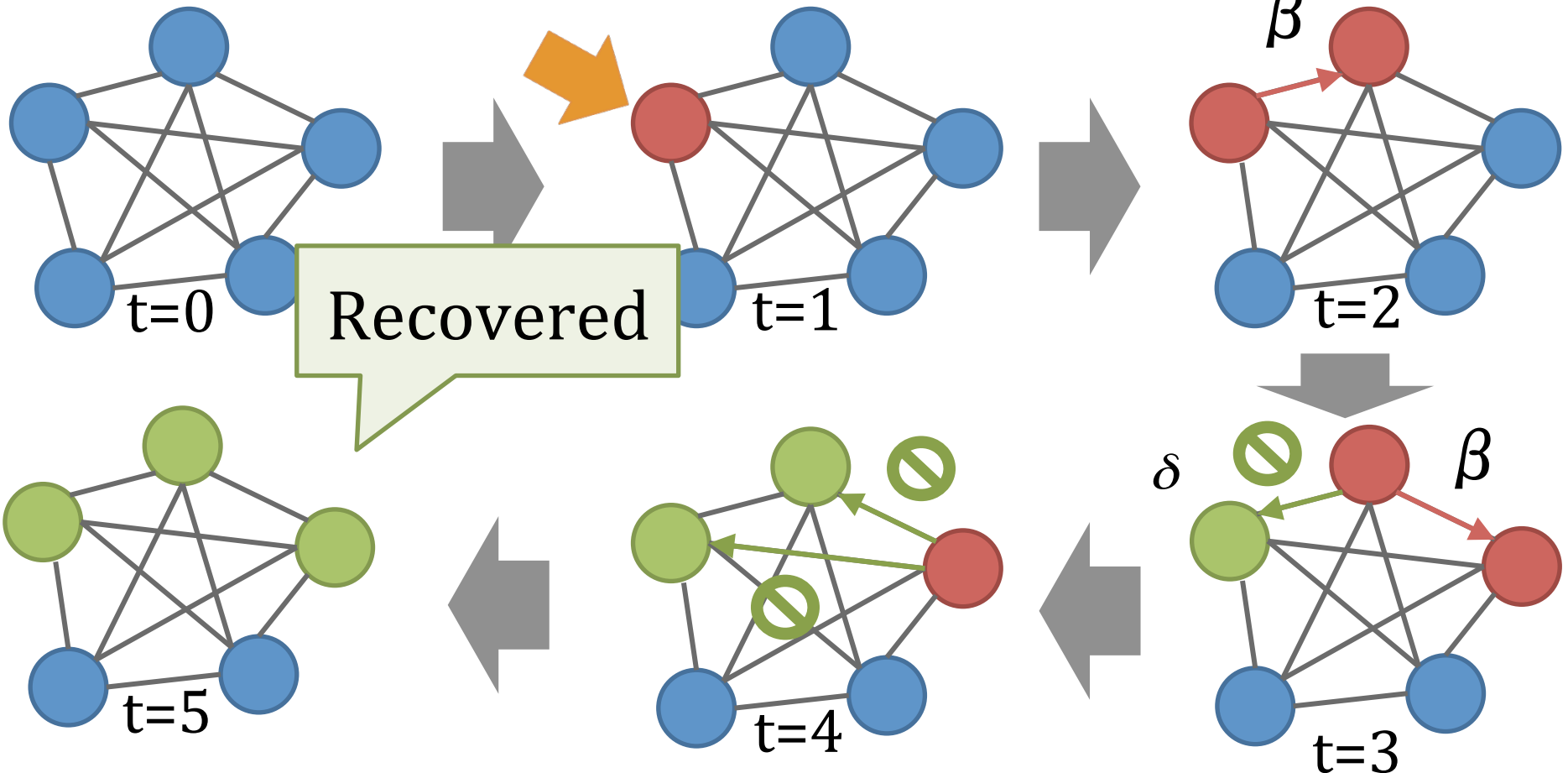
Recovered with immunity



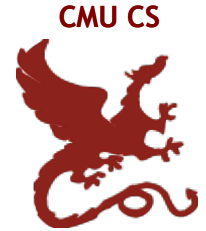
Susceptible-Infected-recovered (SIR) model



Recovered with immunity



Susceptible-Infected-recovered (SIR) model



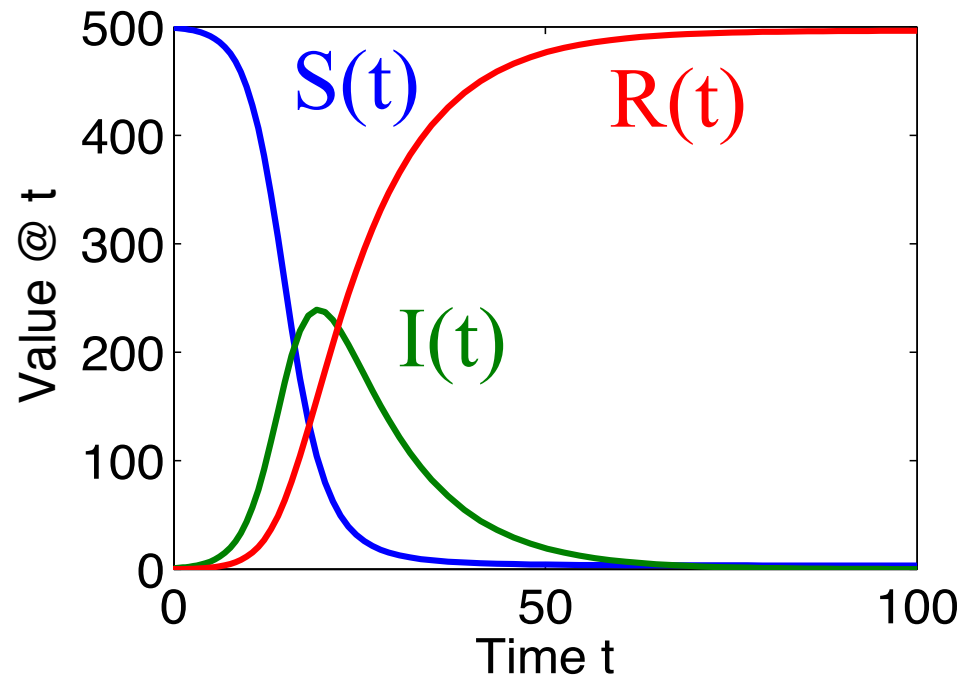
Recovered with immunity

$$\frac{dS}{dt} = -\frac{\beta SI}{N}$$

$$\frac{dI}{dt} = \frac{\beta SI}{N} - \delta I$$

$$\frac{dR}{dt} = \delta I$$

$$S(t) + I(t) + R(t) = N$$

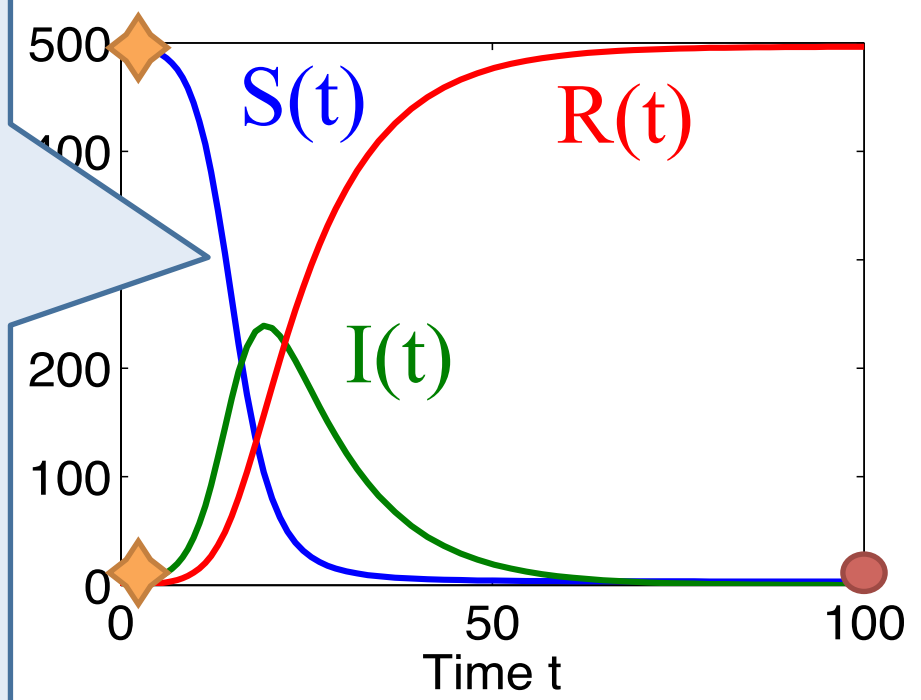
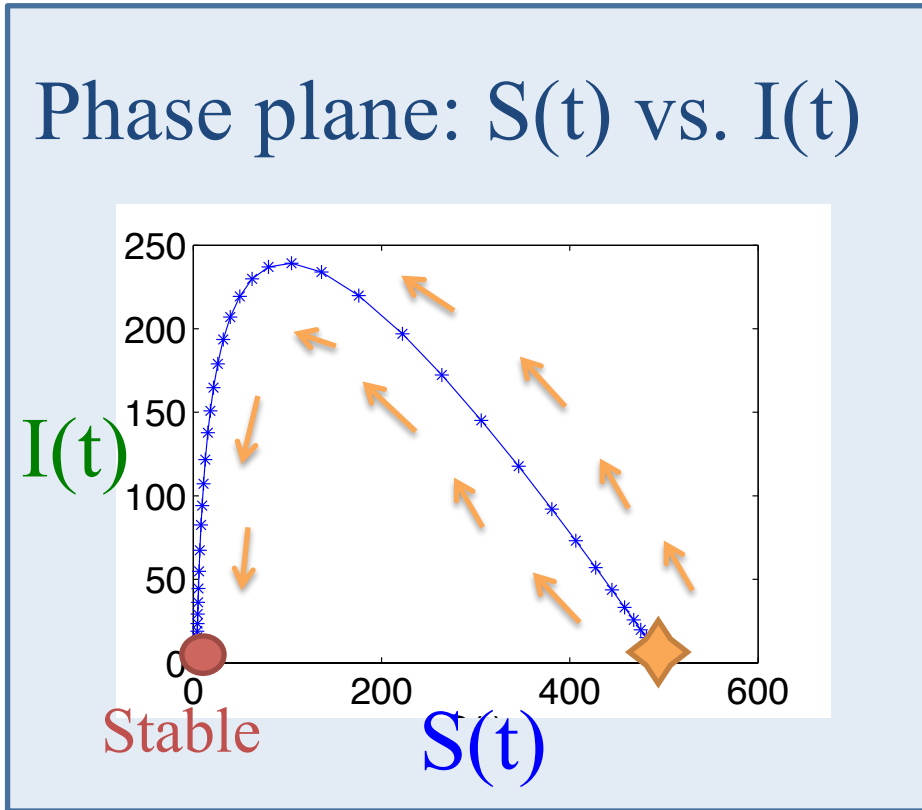


β : Infection rate
 δ : Recovery rate

Susceptible-Infected-recovered (SIR) model

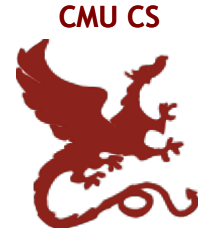


Recovered with immunity

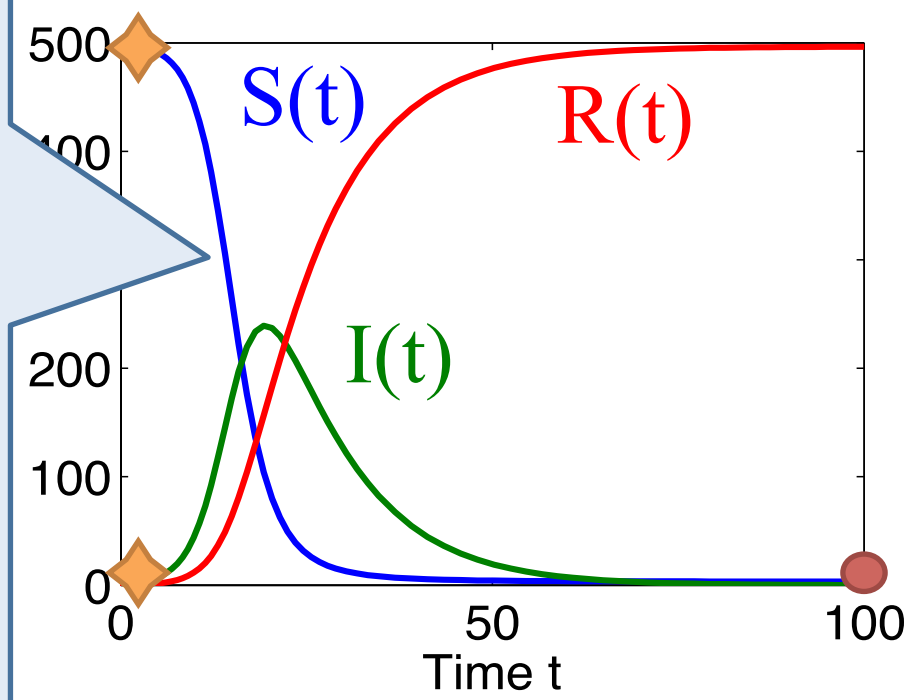
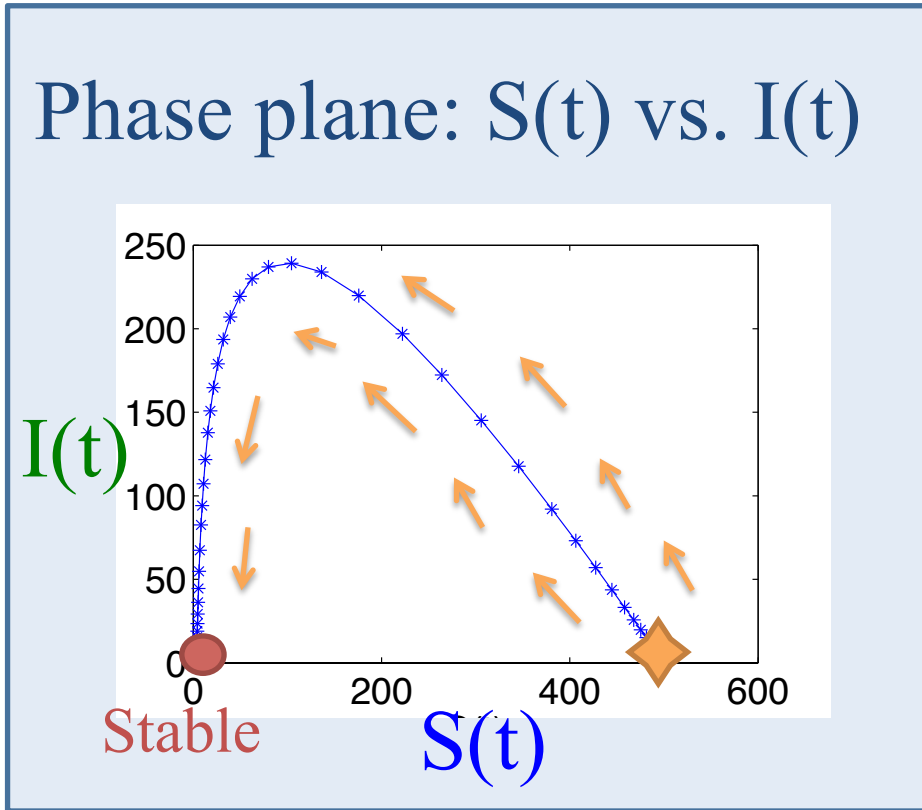


β : Infection rate
 δ : Recovery rate

Susceptible-Infected-recovered (SIR) model



Recovered with immunity



β : Infection rate
 δ : Recovery rate



Other epidemic models

Other virus propagation models (“VPM”)

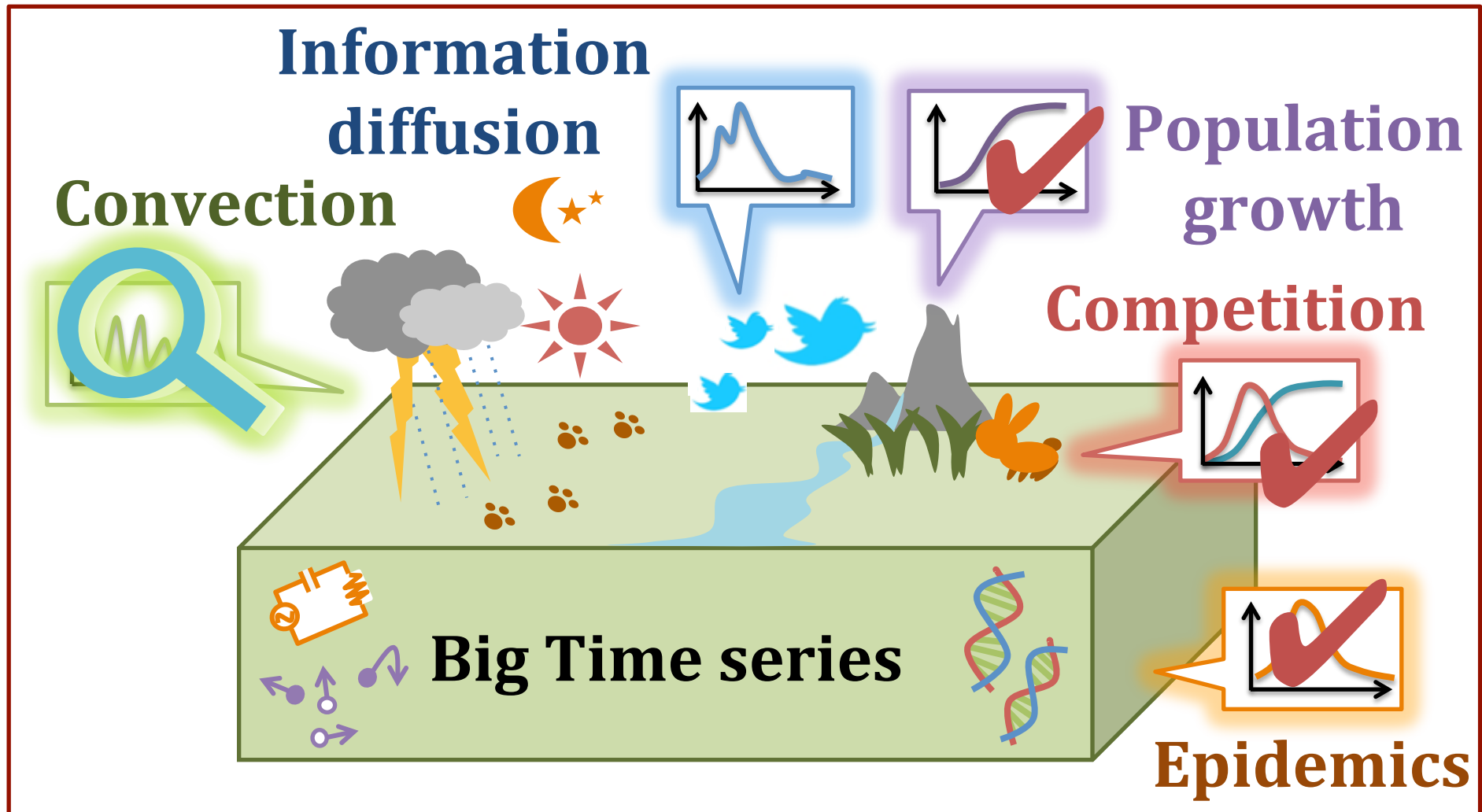
- **SIS** : susceptible-infected-susceptible, flu-like
- **SIRS** : **temporary** immunity, like pertussis
- **SEIR** : mumps-like, with virus **incubation**
(E = Exposed)
- **SEIR-birth/death**: with birth/death rate

Underlying contact-network

- ‘who-can-infect-whom’



Grey-box mining and non-linear equations





Other non-linear models



LORENZ: eqs. for atmospheric convection

$$\frac{dx}{dt} = \sigma(y - x)$$

$$\frac{dy}{dt} = x(\rho - z) - y$$

$$\frac{dz}{dt} = xy - \beta z$$

- x: convective intensity
- y: temperature difference between ascending and descending currents
- z: difference in vertical temperature profile from linearity



Other non-linear models



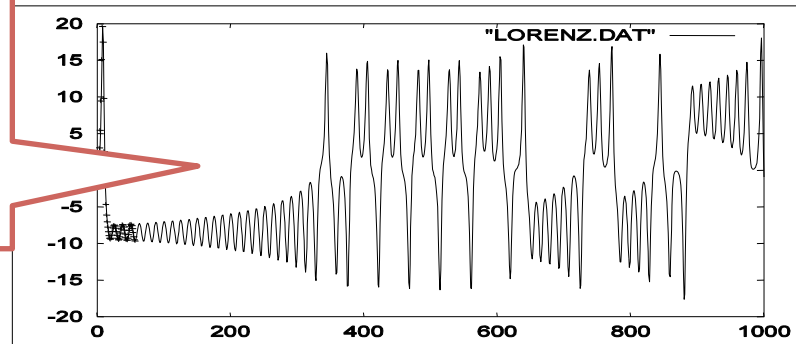
LORENZ: eqs. for atmospheric convection

Butterfly effect
(chaos)

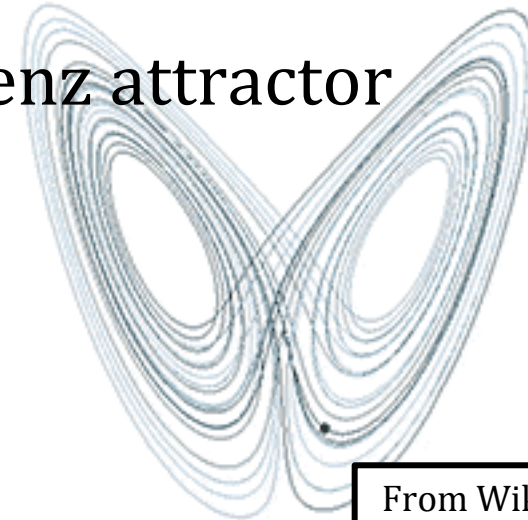
$$\frac{dx}{dt}$$

$$\frac{dy}{dt} = x(\rho - z) - y$$

$$\frac{dz}{dt} = xy - \beta z$$



Lorenz attractor



From Wikipedia

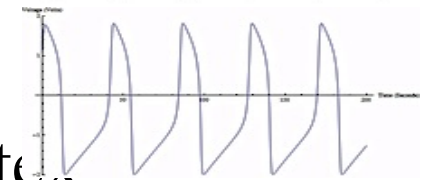
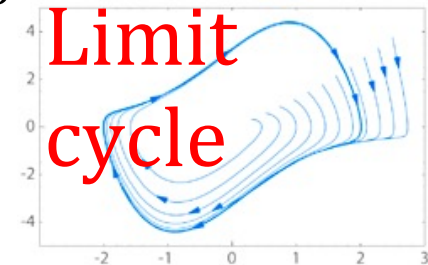
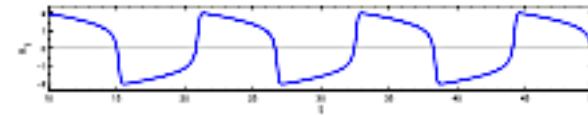


Other non-linear models

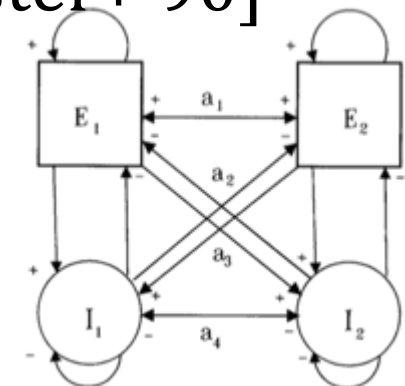


From Wikipedia

- Van del Pol oscillator
 - Electric circuits, heart-beats, neurons
- FitzHugh-Nagumo model
 - An excitable system (e.g., a neuron)
- Excitatory-inhibitory (EI) model
 - Neuronal oscillations in the visual cortex
 - Epilepsy
- ...
- ...



[Schuster+ 90]





Part 2

Roadmap



Problem

- ✓ Why: “non-linear” modeling

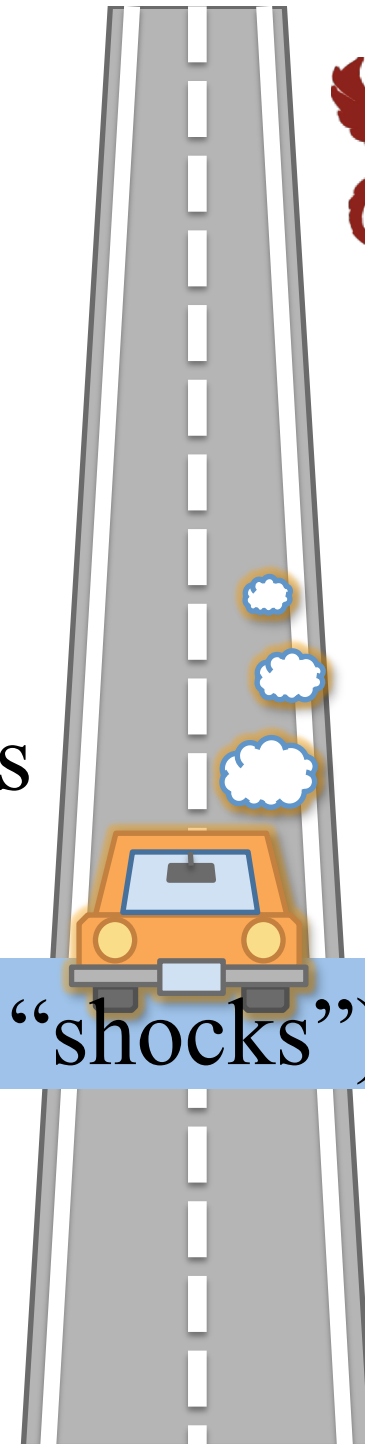
Fundamentals

- ✓ Non-linear (“gray-box”) models

Applications

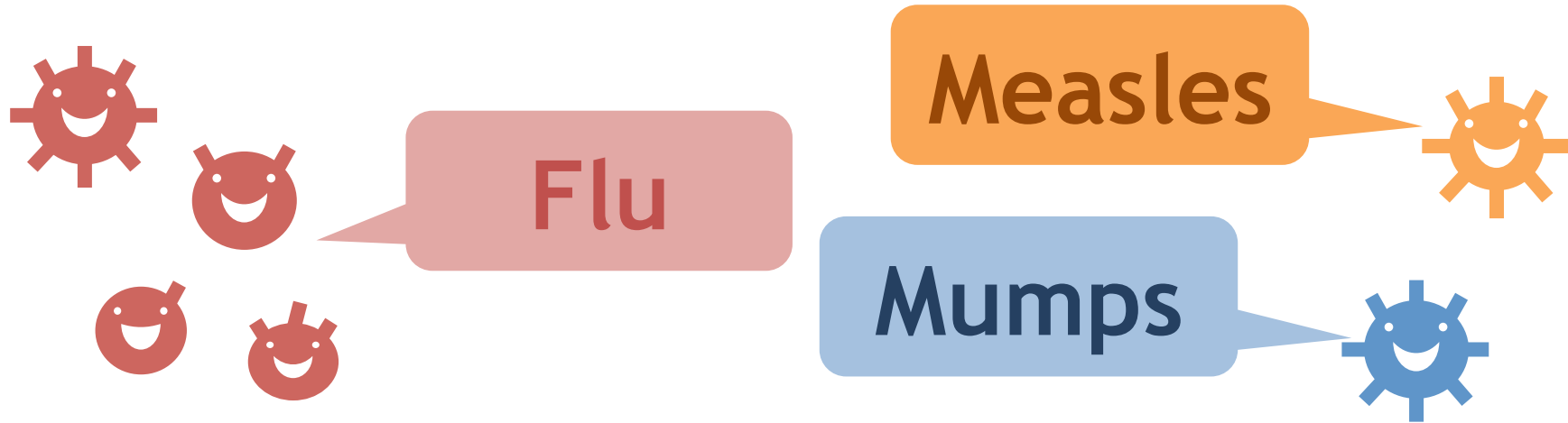


- Epidemics (skips, competition, “shocks”)
- Information diffusion
- Online competition



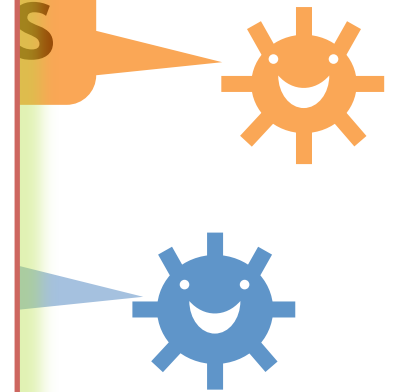
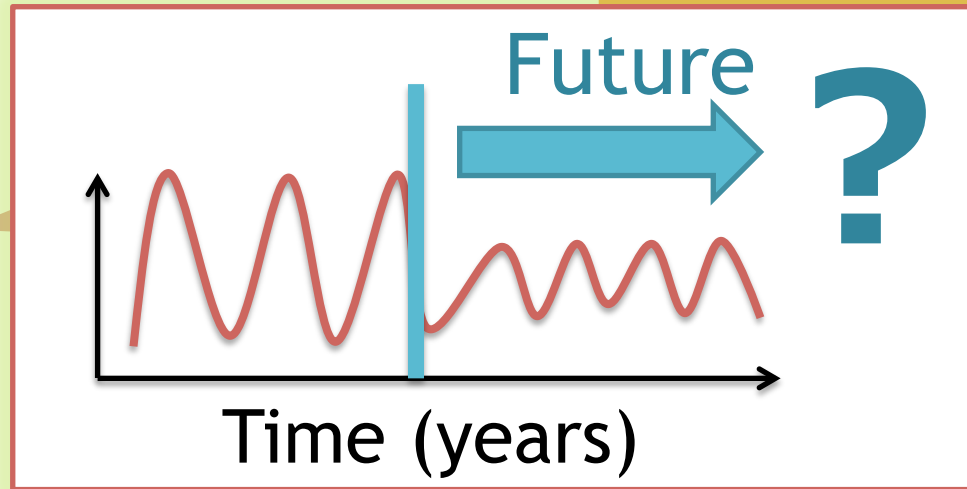
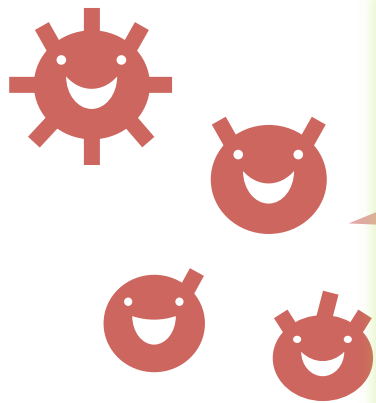


Mining and forecasting of co-evolving epidemics





Mining and forecasting of co-evolving epidemics

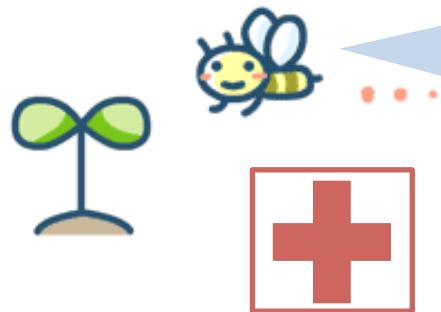


Q. Can we forecast future epidemics? 😊



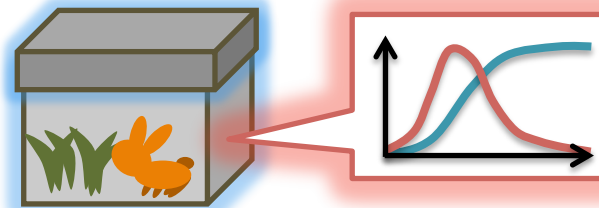


Epidemics - roadmap



A. Non-linear (gray-box) modeling!

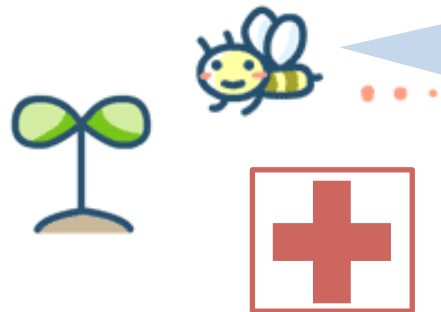
Solutions



- Outbreak vs. Skips [Stone+ Nature'07]
- Interaction between diseases [Rohani+ Nature'03]
- FUNNEL [Matsubara+ KDD'14]

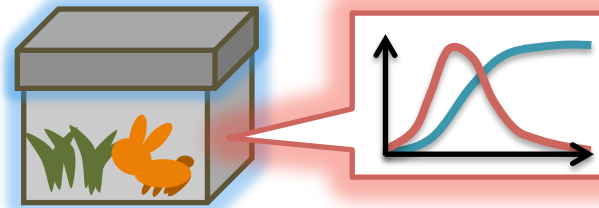


Epidemics - roadmap



A. Non-linear (gray-box) modeling!

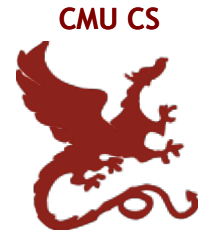
Solutions



- **Outbreak vs. Skips** [Stone+ Nature'07]
- **Interaction between diseases** [Rohani+ Nature'03]
- **FUNNEL** [Matsubara+ KDD'14]



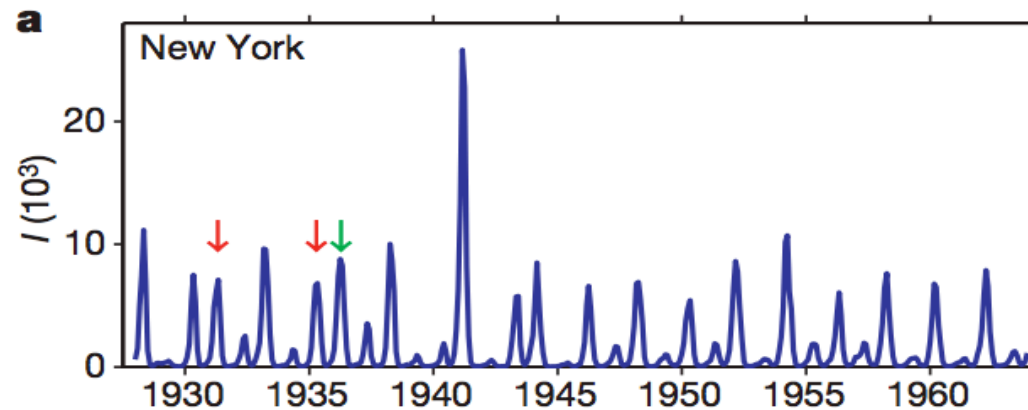
Recurrent epidemics: Outbreak or skip?



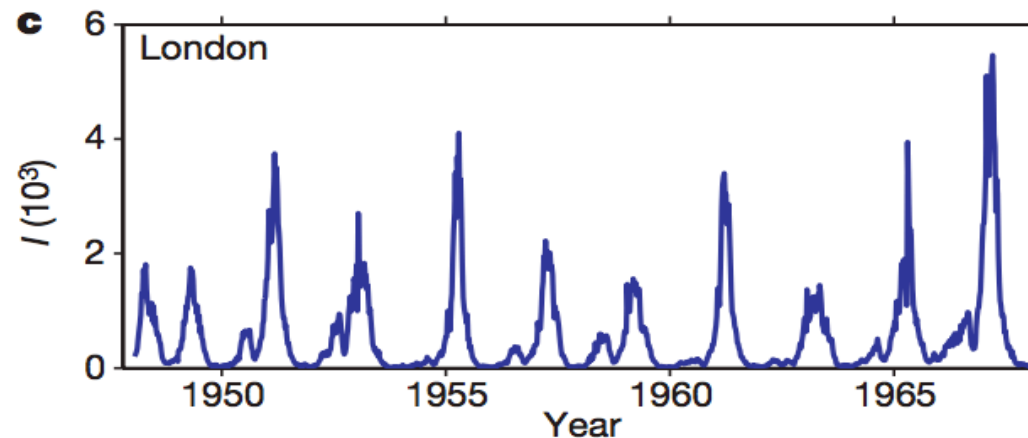
[Stone+ Nature'07]

- Time series of reported measles cases

New York



London





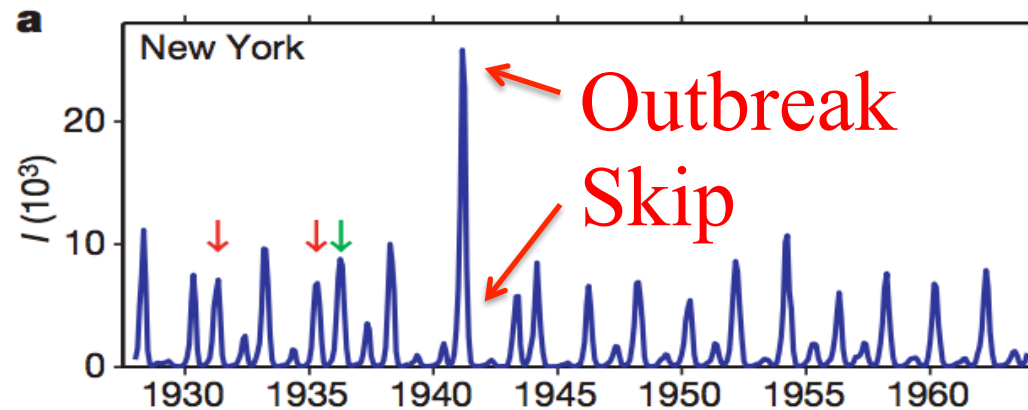
Recurrent epidemics: Outbreak or skip?



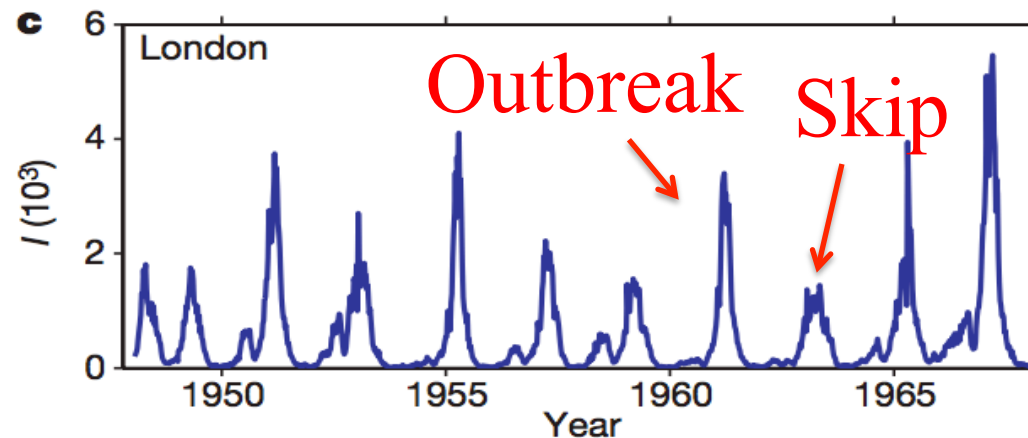
[Stone+ Nature'07]

- Time series of reported measles cases

New York

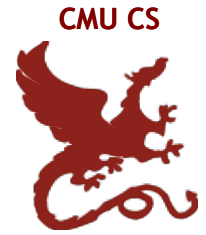


London





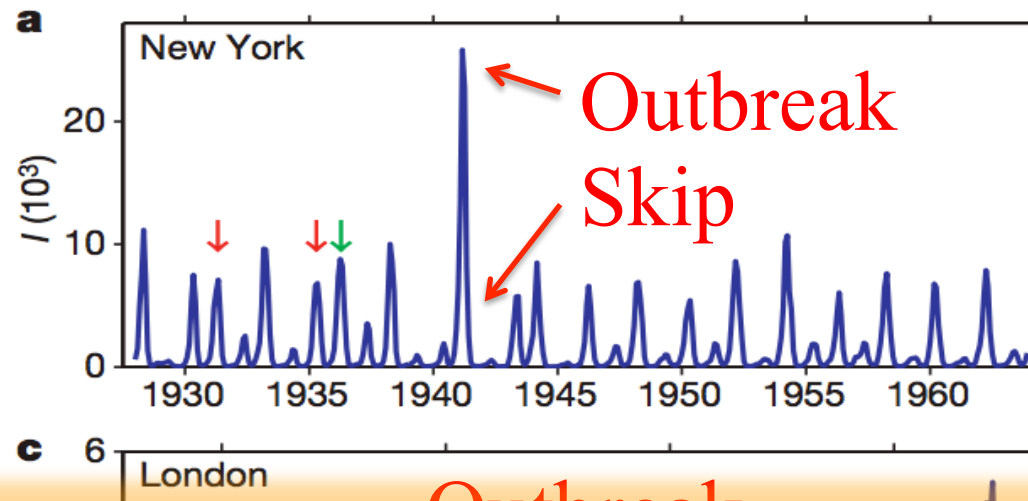
Recurrent epidemics: Outbreak or skip?



[Stone+ Nature'07]

- Time series of reported measles cases

New York



Q. Outbreak vs. skip?

1950 1955 1960 1965

Year

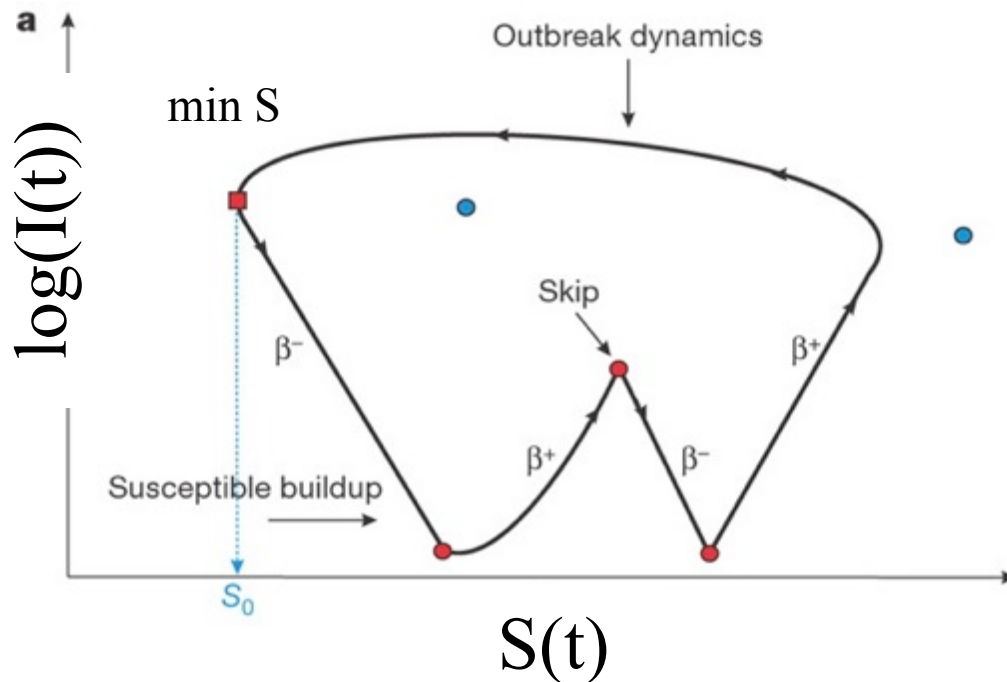
Recurrent epidemics: Outbreak or skip?



[Stone+ Nature'07]

- Conditions for predicting “outbreak vs. skip”
 - SIR model with high/low seasons

Phase plane diagram (S vs. $\log(I)$)



Contact rate
 β^+ : high season
 β^- : low season

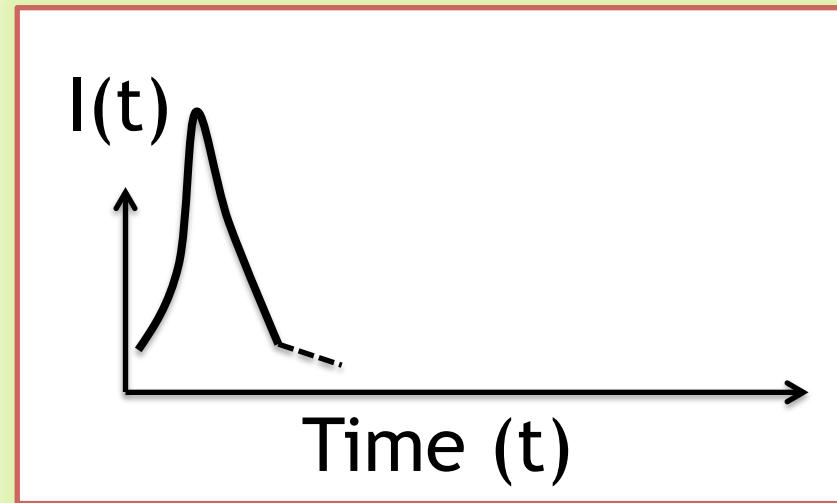
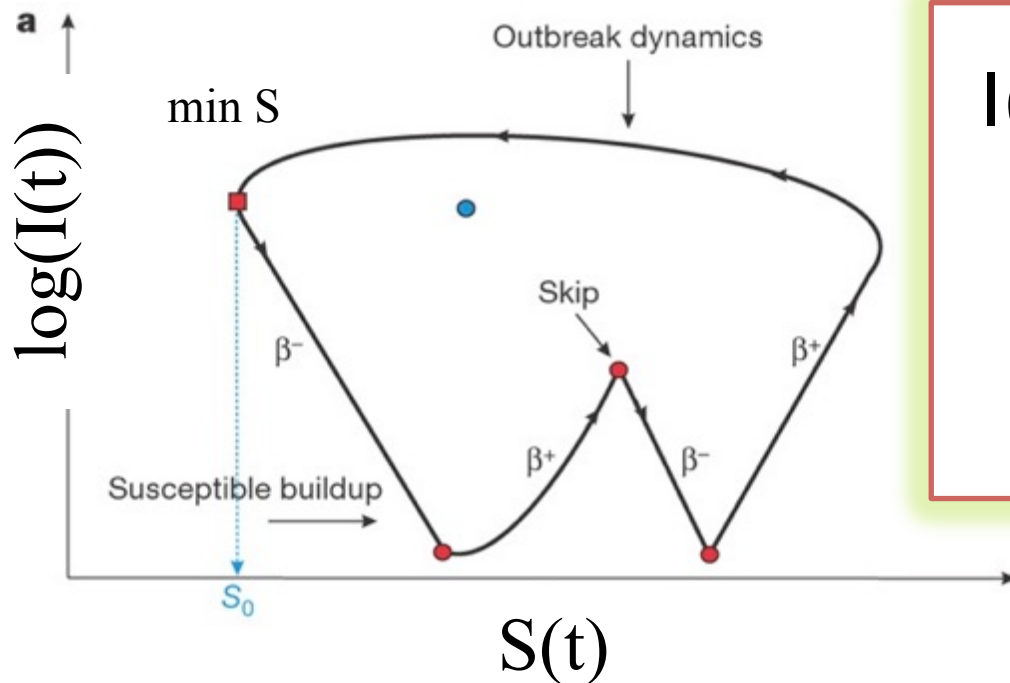
Recurrent epidemics: Outbreak or skip?



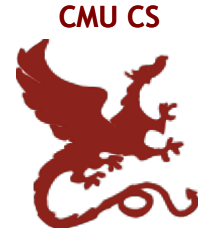
[Stone+ Nature'07]

- Conditions for predicting “outbreak vs. skip”
 - SIR model with high/low seasons

Phase plane diagram (S vs. $\log(I)$)



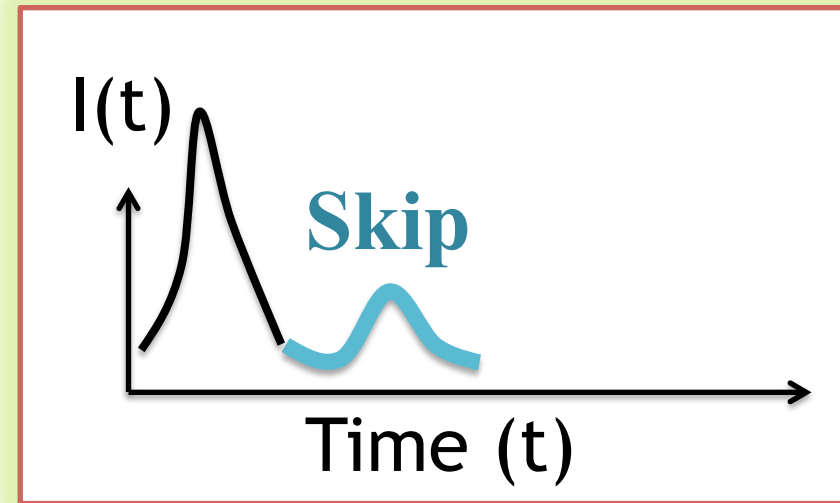
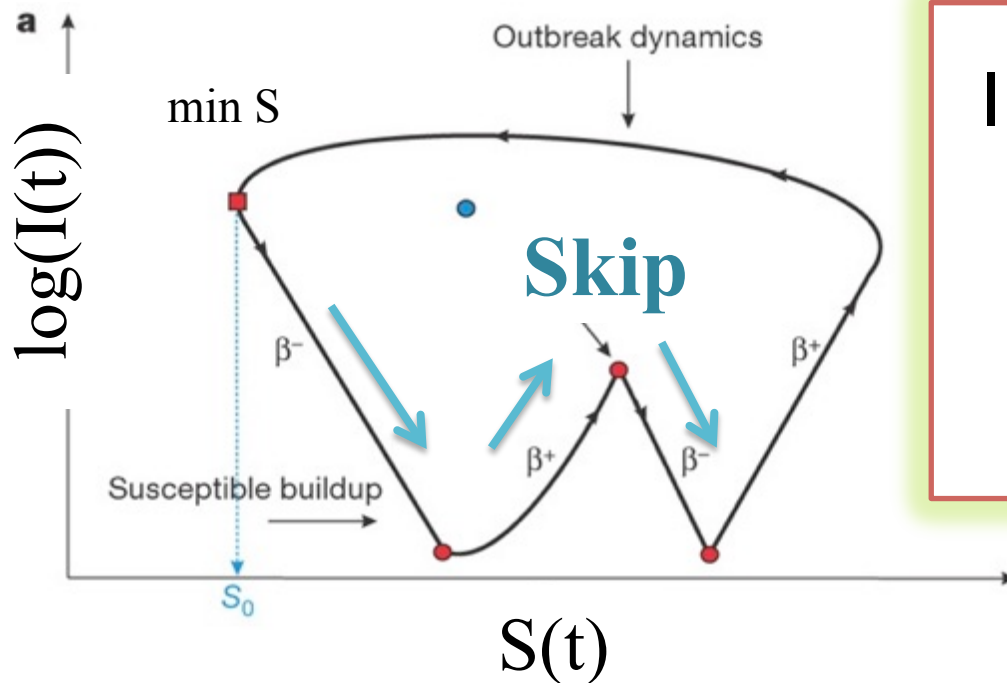
Recurrent epidemics: Outbreak or skip?



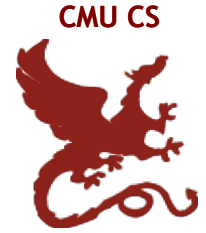
[Stone+ Nature'07]

- Conditions for predicting “outbreak vs. skip”
 - SIR model with high/low seasons

Phase plane diagram (S vs. $\log(I)$)



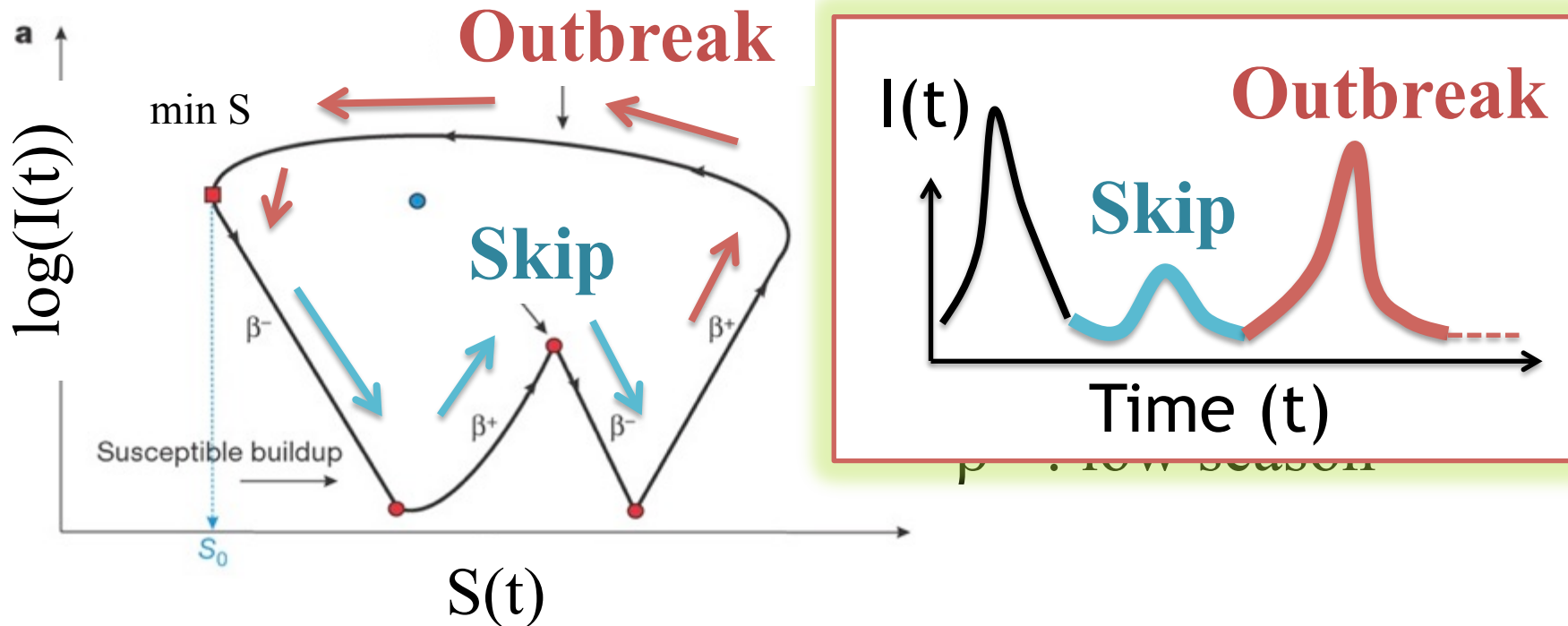
Recurrent epidemics: Outbreak or skip?



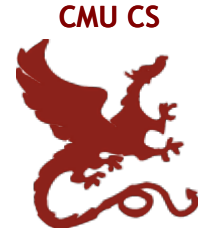
[Stone+ Nature'07]

- Conditions for predicting “outbreak vs. skip”
 - SIR model with high/low seasons

Phase plane diagram (S vs. $\log(I)$)



Recurrent epidemics: Outbreak or skip?

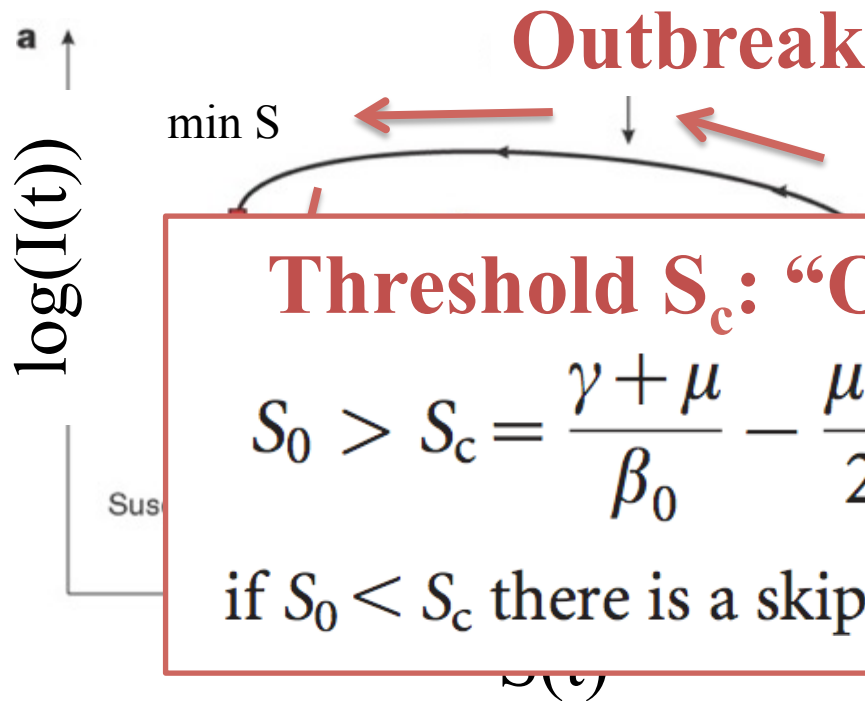


[Stone+ Nature'07]

- Conditions for predicting “outbreak vs. skip”
 - SIR model with high/low seasons

Phase plane diagram (S vs. log(I))

γ : recover rate
 μ : birth/death rate
 β_0 :infection rate
 χ : time period



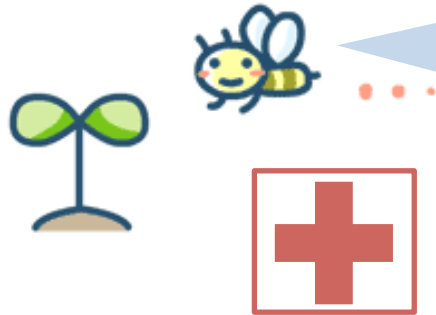
Threshold S_c : “Outbreak vs. Skip”

$$S_0 > S_c = \frac{\gamma + \mu}{\beta_0} - \frac{\mu\chi}{2} \Rightarrow \text{epidemic}$$

if $S_0 < S_c$ there is a skip in the following year.

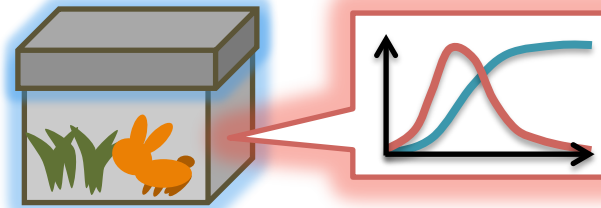


Epidemics - roadmap



A. Non-linear (gray-box) modeling!

Solutions



- Outbreak vs. Skips [Stone+ Nature'07]
- **Interaction between diseases** [Rohani+ Nature'03]
- FUNNEL [Matsubara+ KDD'14]





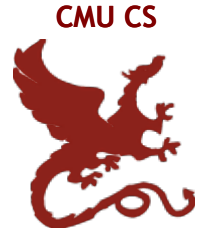
Ecological interference between fatal diseases



Q. Any relationship (i.e., interaction)
between two different diseases
(e.g., measles vs. whooping cough)?

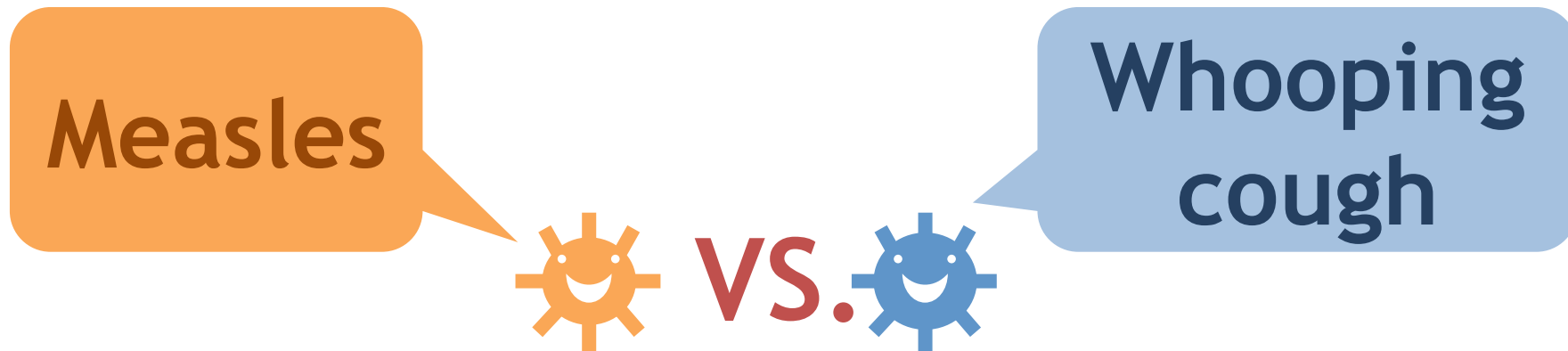


Ecological interference between fatal diseases



Q. Any relationship (i.e., interaction)
between two different diseases
(e.g., measles vs. whooping cough)?

A. Yes. There are “competing” diseases!





Ecological interference between fatal diseases

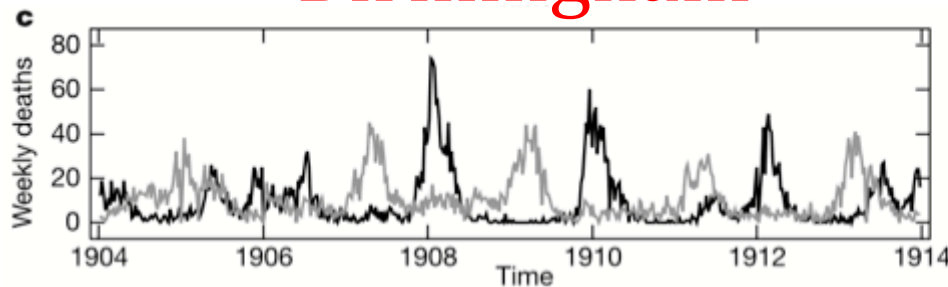


[Rohani+ Nature'03]

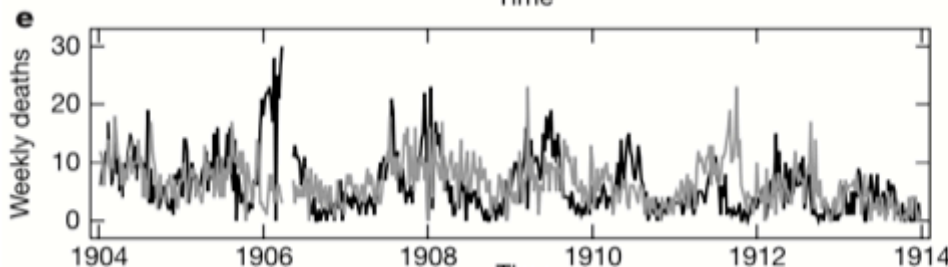
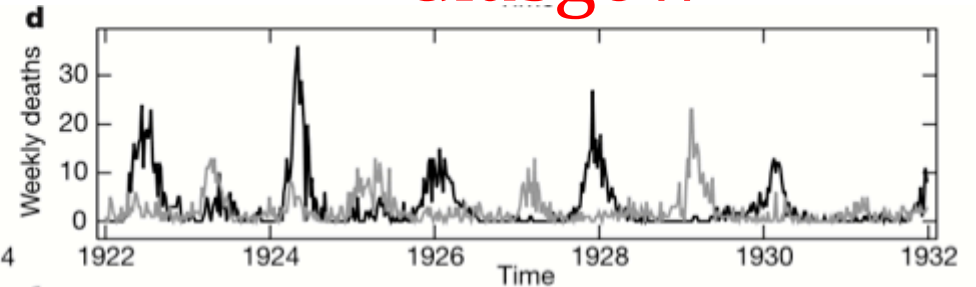
Weekly case fatality reports for two diseases

— measles — Whooping cough

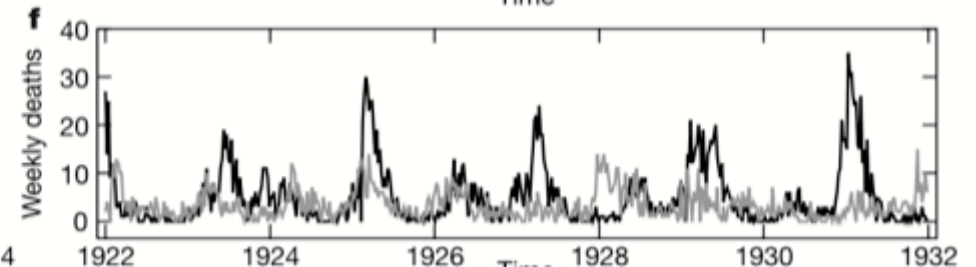
Birmingham



Glasgow



Berlin



Liverpool

Ecological interference between fatal diseases

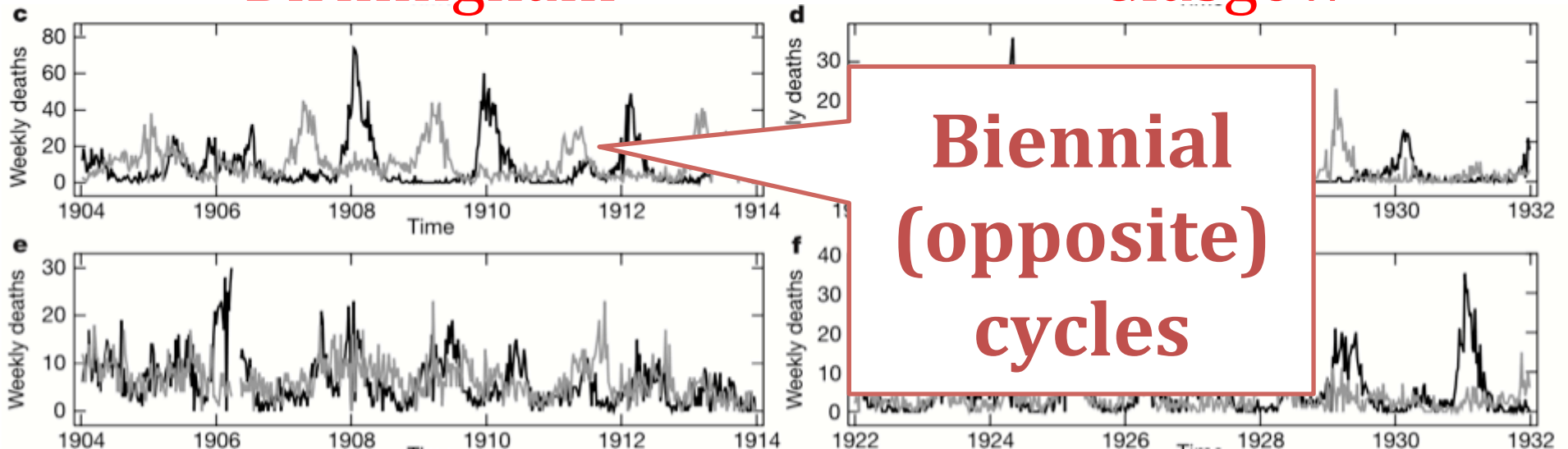
[Rohani+ Nature'03]

Weekly case fatality reports for two diseases

— measles — Whooping cough

Birmingham

Glasgow



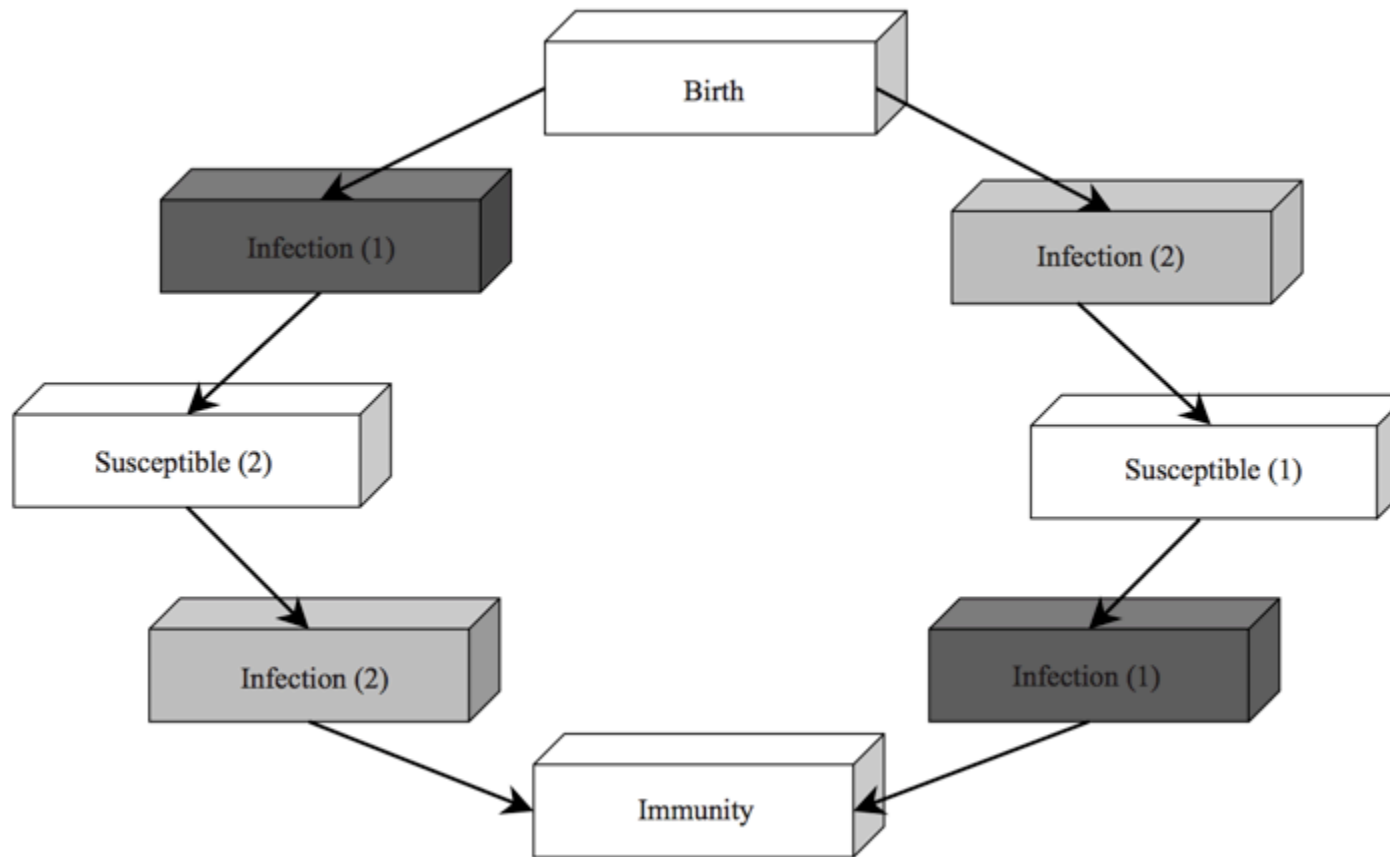
Berlin

Liverpool

Ecological interference between fatal diseases

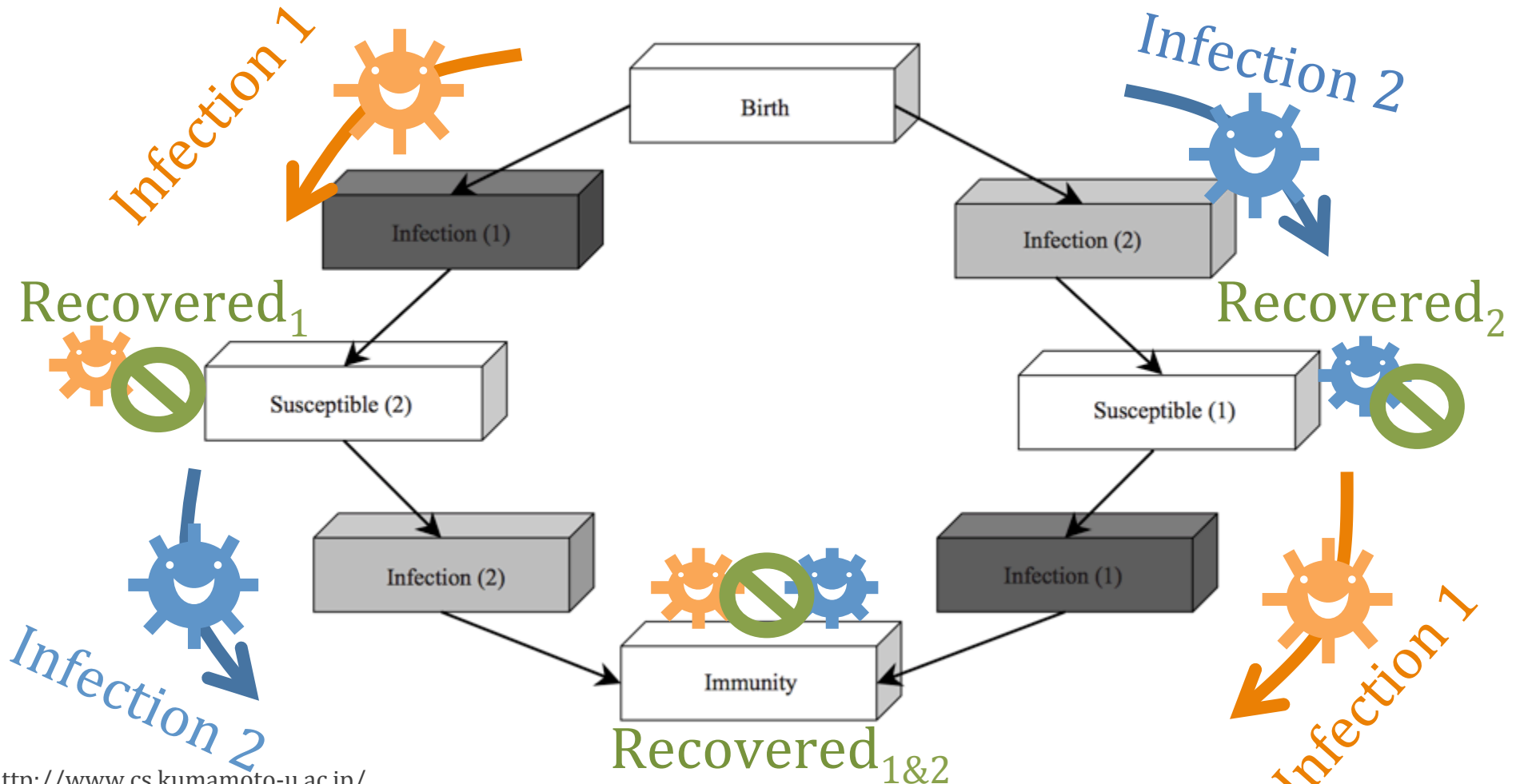


Extension of SIR model [Rohani+'98]



Ecological interference between fatal diseases

Extension of SIR model [Rohani+'98]



Ecological interference between fatal diseases



Equations for 3 disease model

[Rohani+ Nature'03]

$$\frac{dS_{SSS}}{dt} = \nu N(1 - p) - \mu S_{SSS}$$

$$\text{⚙️} - \frac{\beta_1(t) S_{SSS}}{N} (I_{IRR} + I_{IRT} + I_{ITR} + I_{ITT})$$

$$\text{⚙️} - \frac{\beta_2(t) S_{SSS}}{N} (I_{RIR} + I_{RIT} + I_{TIR} + I_{TIT})$$

$$\text{⚙️} - \frac{\beta_3(t) S_{SSS}}{N} (I_{RRI} + I_{RTI} + I_{TRI} + I_{TTI})$$

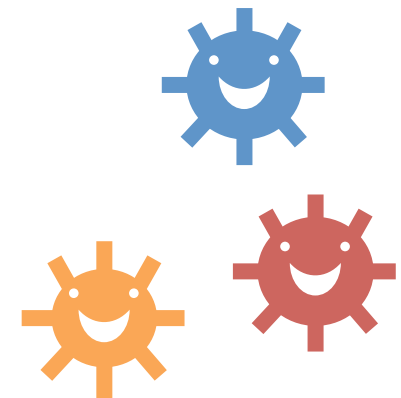
$$\frac{dI_{ITT}}{dt} = \frac{\beta_1(t) S_{SSS}}{N} (I_{IRR} + I_{IRT} + I_{ITR} + I_{ITT})$$

$$- (\mu + \gamma_1) I_{ITT}$$

$$\frac{dI_{IRT}}{dt} = \frac{\beta_1(t) S_{SSS}}{N} (I_{IRR} + I_{IRT} + I_{ITR} + I_{ITT})$$

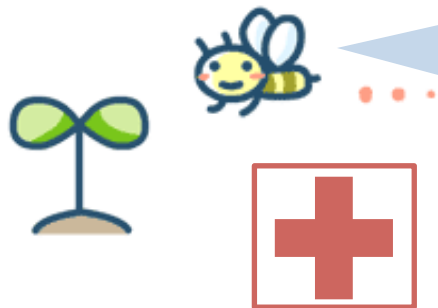
$$- (\mu + \gamma_1) I_{IRT}$$

...



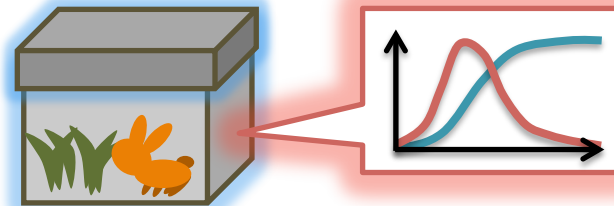


Epidemics - roadmap



Non-linear (gray-box)
modeling!

Solutions

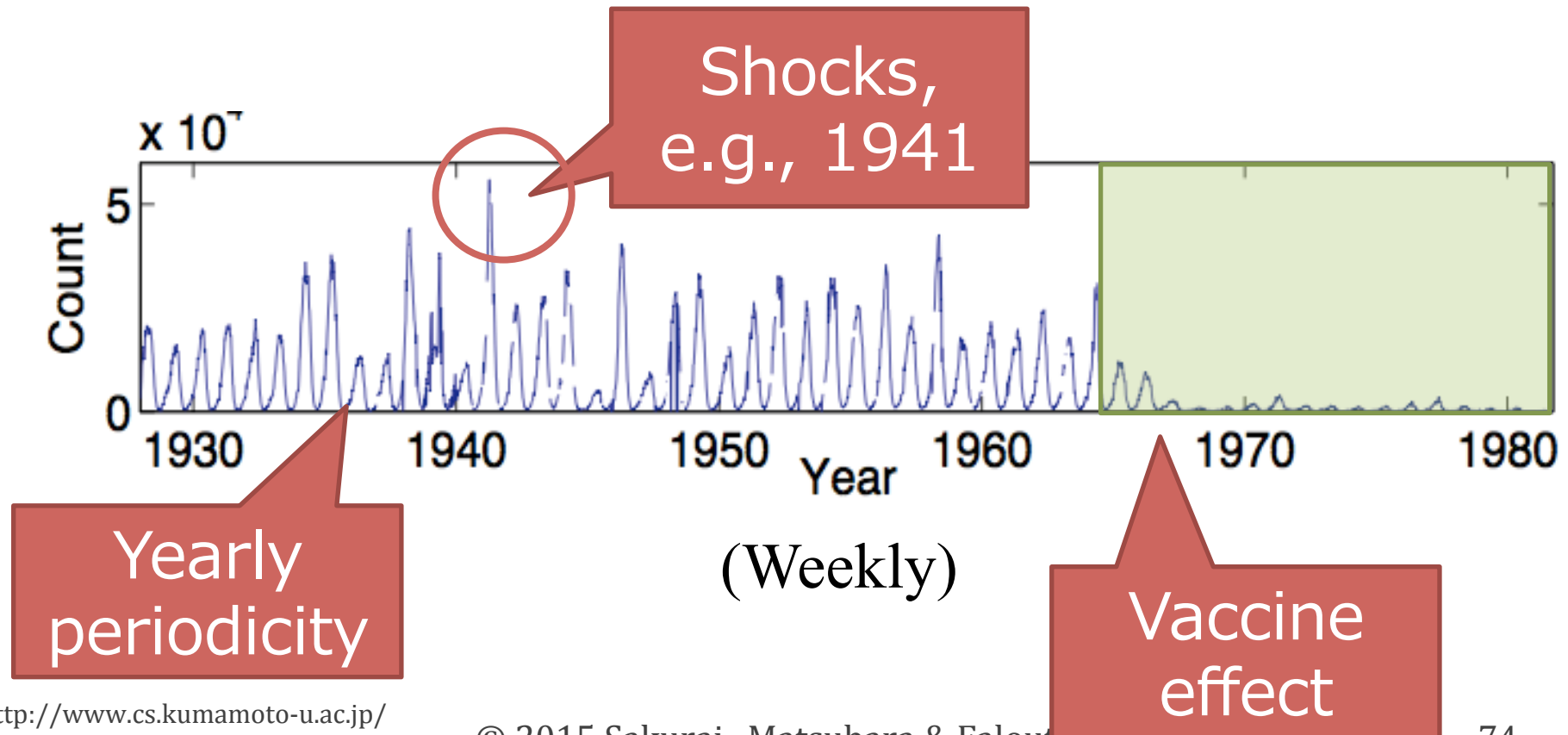


- E1. Outbreak vs. Skips [Stone+ Nature'07]
- E2. Interaction between diseases [Rohani+ Nature'03]
- **E3. FUNNEL** [Matsubara+ KDD'14]



with a single epidemic

e.g., Measles cases in the U.S.





FUNNEL

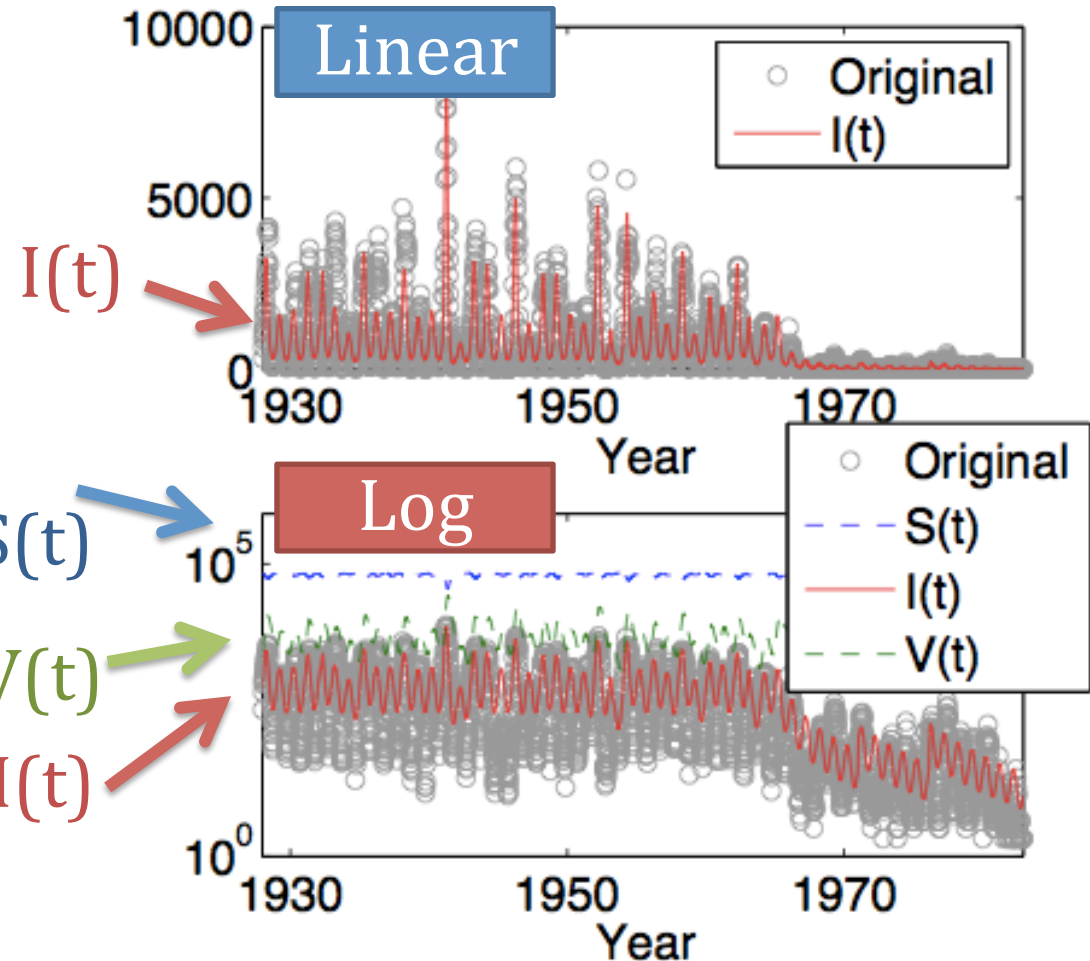
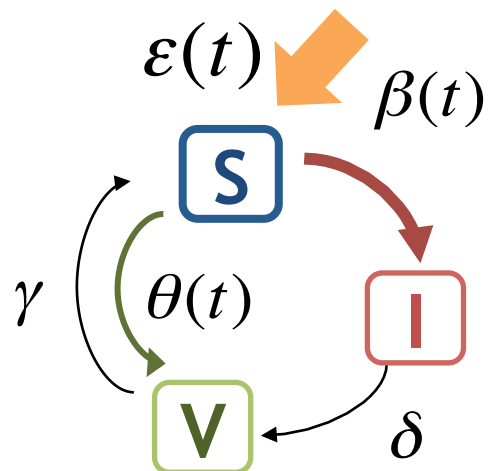
[Matsubara+ KDD'14]



with a single epidemic

With a single epidemic: Funnel-RE

- People of 3 classes
- **S** : Susceptible
 - **I** : Infected
 - **V** : Vigilant/
vaccinated





with a single epidemic

With a single epidemic: Funnel-RE

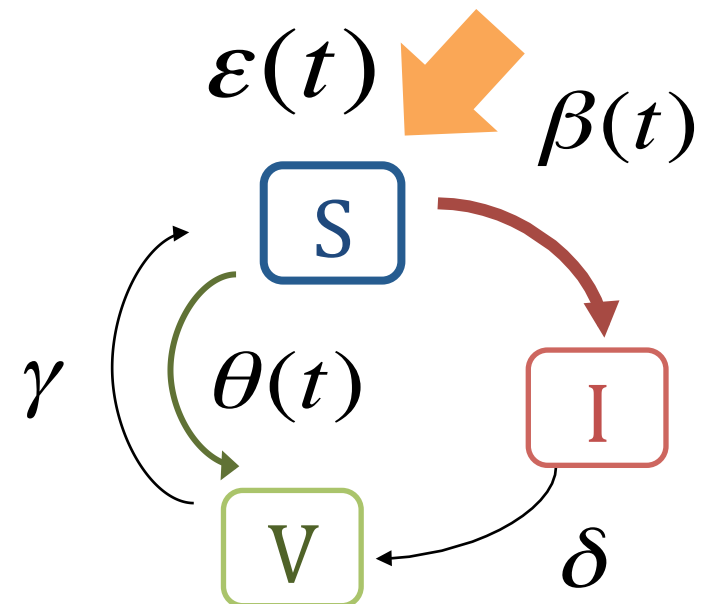
$$\begin{aligned}
 S(t+1) &= S(t) - \beta(t)\epsilon(t)S(t)I(t) + \gamma V(t) - \theta(t)S(t) \\
 I(t+1) &= I(t) + \beta(t)\epsilon(t)S(t)I(t) - \delta I(t) \\
 V(t+1) &= V(t) + \delta I(t) - \gamma V(t) + \theta(t)S(t)
 \end{aligned} \tag{3}$$

S(t) : susceptible

I(t) : Infected

V(t) : Vigilant

/Vaccinated





with a single epidemic

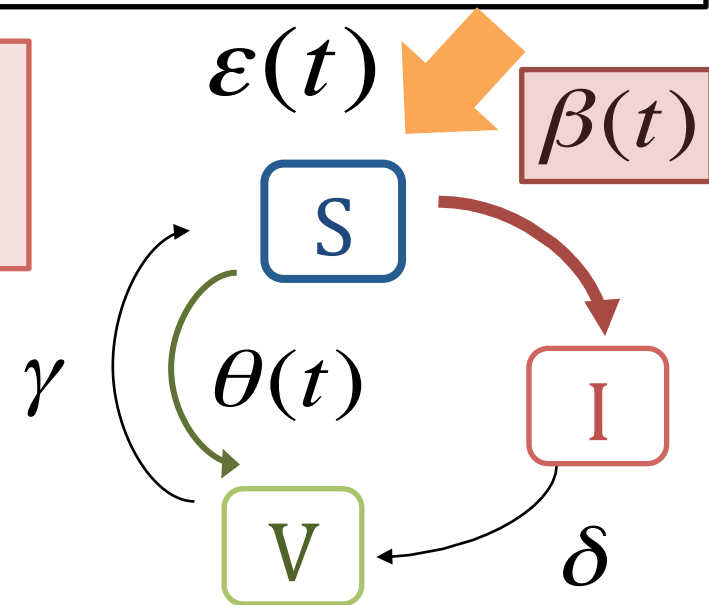
With a single epidemic: Funnel-RE

$$\begin{aligned}
 S(t+1) &= S(t) - \beta(t)\epsilon(t)S(t)I(t) + \gamma V(t) - \theta(t)S(t) \\
 I(t+1) &= I(t) + \beta(t)\epsilon(t)S(t)I(t) - \delta I(t) \\
 V(t+1) &= V(t) + \delta I(t) - \gamma V(t) + \theta(t)S(t)
 \end{aligned} \tag{3}$$

$\beta(t)$: strength of infection
(yearly periodic func)

$$\beta(t) = \beta_0 \cdot \left(1 + P_a \cdot \cos\left(\frac{2\pi}{P_p}(t + P_s)\right) \right)$$

$P_p = 52$





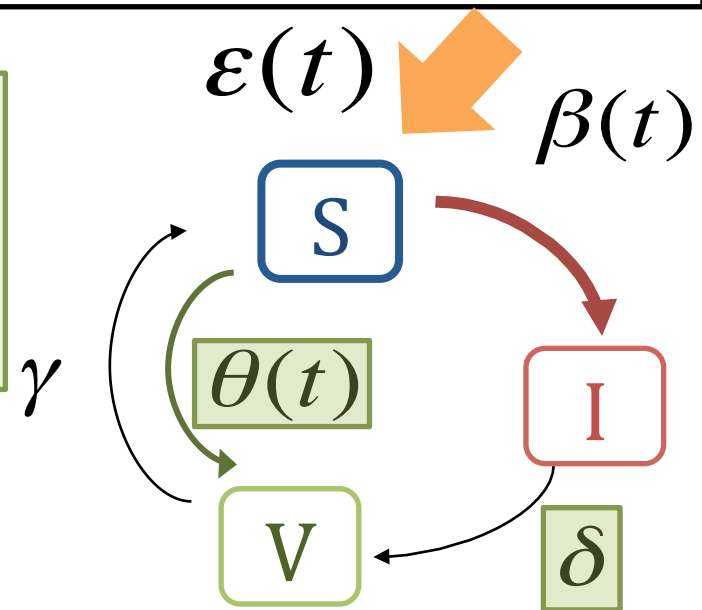
with a single epidemic

With a single epidemic: Funnel-RE

$$\begin{aligned}
 S(t+1) &= S(t) - \beta(t)\epsilon(t)S(t)I(t) + \gamma V(t) - \theta(t)S(t) \\
 I(t+1) &= I(t) + \beta(t)\epsilon(t)S(t)I(t) - \delta I(t) \\
 V(t+1) &= V(t) + \delta I(t) - \gamma V(t) + \theta(t)S(t)
 \end{aligned} \tag{3}$$

δ : healing rate
 $\theta(t)$: disease reduction effect

$$\theta(t) = \begin{cases} 0 & (t < t_\theta) \\ \theta_0 & (t \geq t_\theta) \end{cases}$$



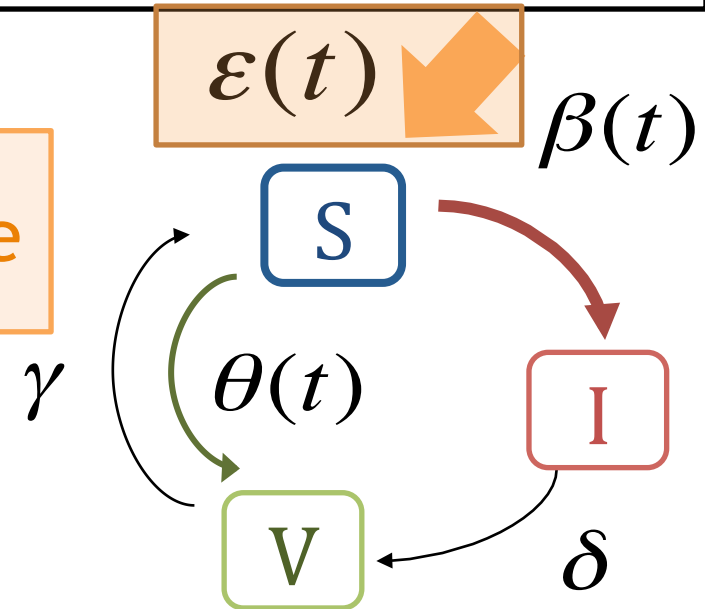


with a single epidemic

With a single epidemic: Funnel-RE

$$\begin{aligned}
 S(t+1) &= S(t) - \beta(t)\epsilon(t)S(t)I(t) + \gamma V(t) - \theta(t)S(t) \\
 I(t+1) &= I(t) + \beta(t)\epsilon(t)S(t)I(t) - \delta I(t) \\
 V(t+1) &= V(t) + \delta I(t) - \gamma V(t) + \theta(t)S(t)
 \end{aligned} \tag{3}$$

$\epsilon(t)$: temporal susceptible rate





FUNNEL



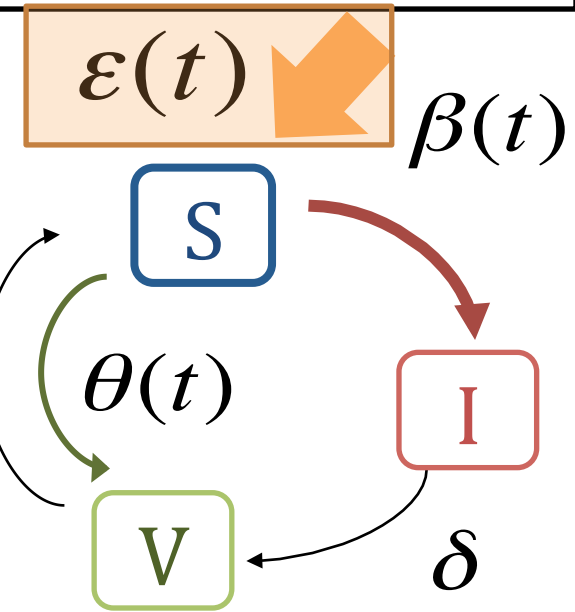
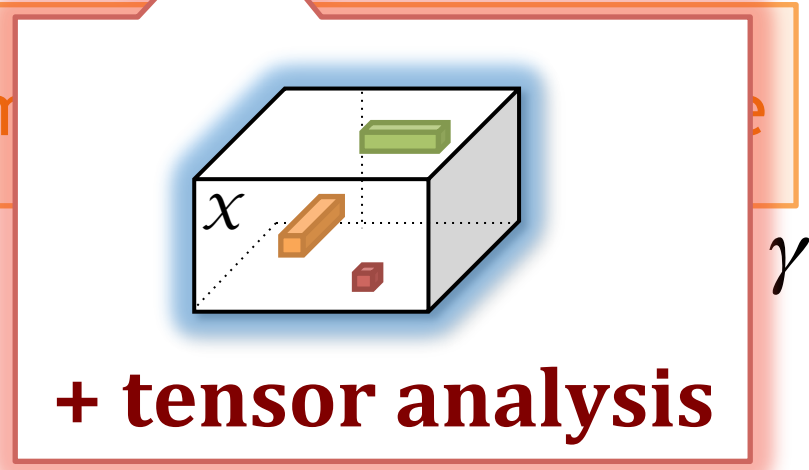
with a single epidemic

With a single epidemic: Funnel-RF

$$\begin{aligned}
 S(t+1) &= S(t) - \beta(t)\epsilon(t)S(t) - \theta(t)S(t) \\
 I(t+1) &= I(t) + \beta(t)\epsilon(t)S(t) - \delta I(t) \\
 V(t+1) &= V(t) + \delta I(t) - \gamma V(t) + \theta(t)S(t)
 \end{aligned} \tag{3}$$

FUNNEL: Details @ part3

$\epsilon(t)$: tem





Part 2

Roadmap



Problem

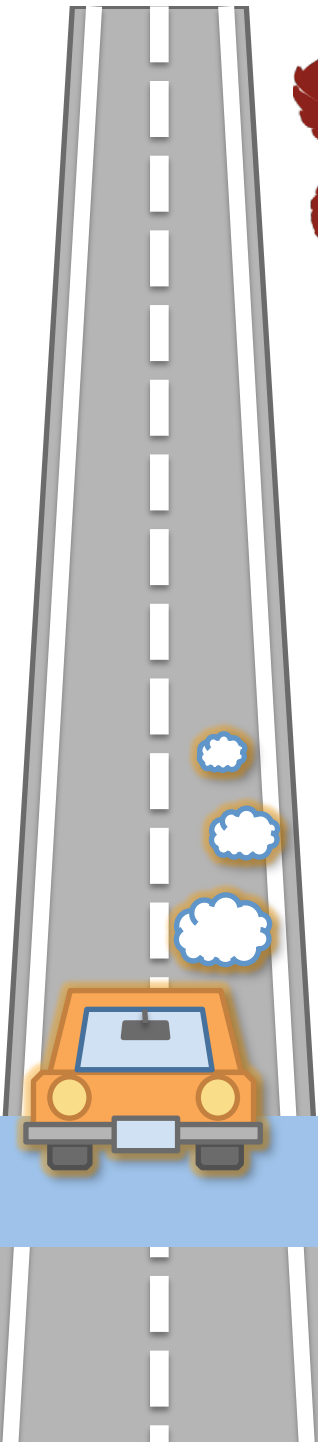
- ✓ Why: “non-linear” modeling

Fundamentals

- ✓ Non-linear (grey-box) models

Applications

- ✓ Epidemics
- Information diffusion
- Online competition



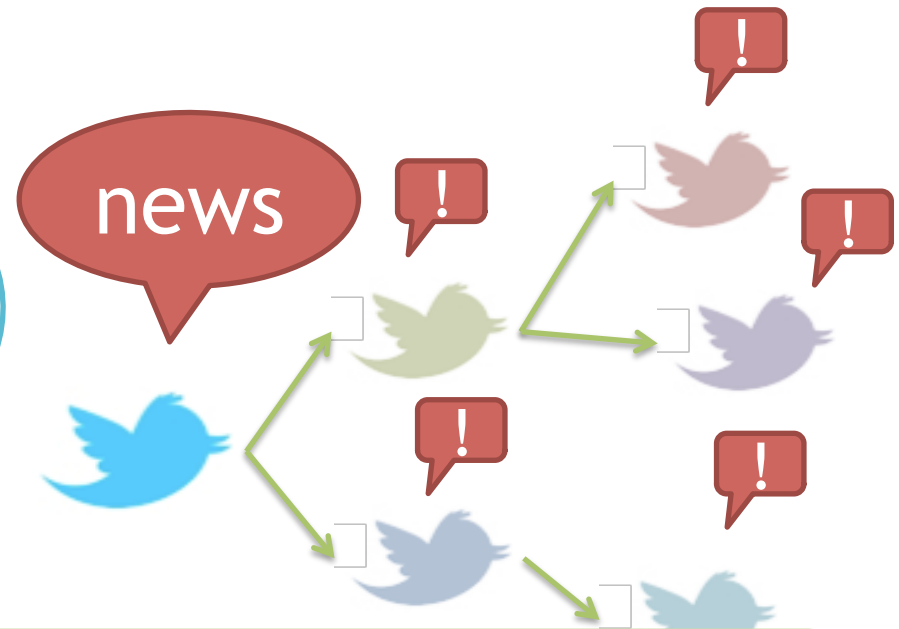


Information diffusion in social networks





Information diffusion in social networks



Q. How news/rumors spread in social media?



News spread in social media

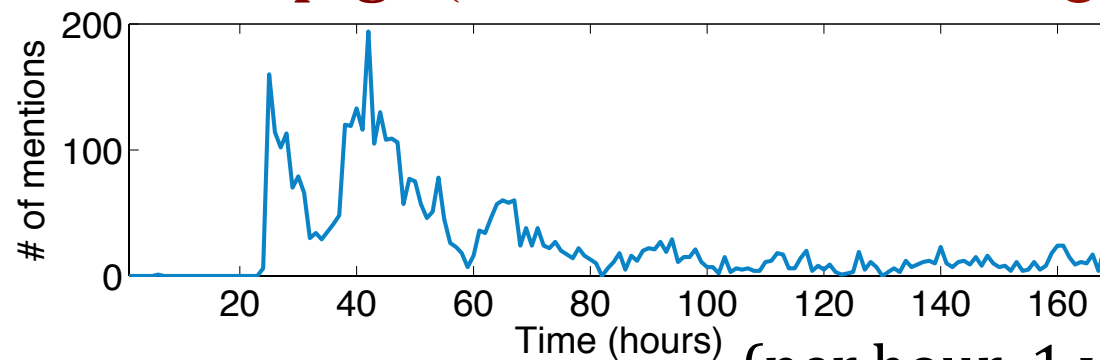


MemeTracker [Leskovec+ KDD'09]



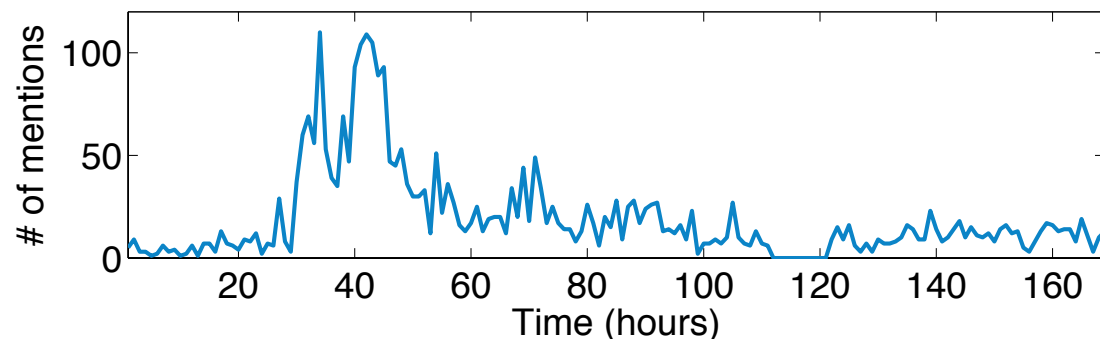
- Short phrases sourced from U.S. politics in 2008

“you can put lipstick on a pig” (# of mentions in blogs)



(per hour, 1 week)

“yes we can”



News spread in social media

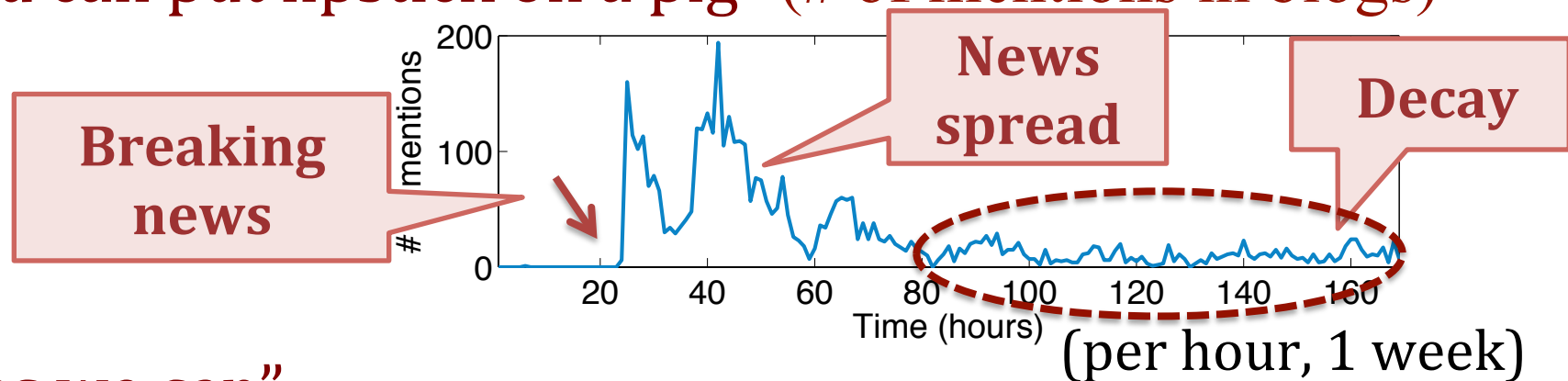


MemeTracker [Leskovec+ KDD'09]

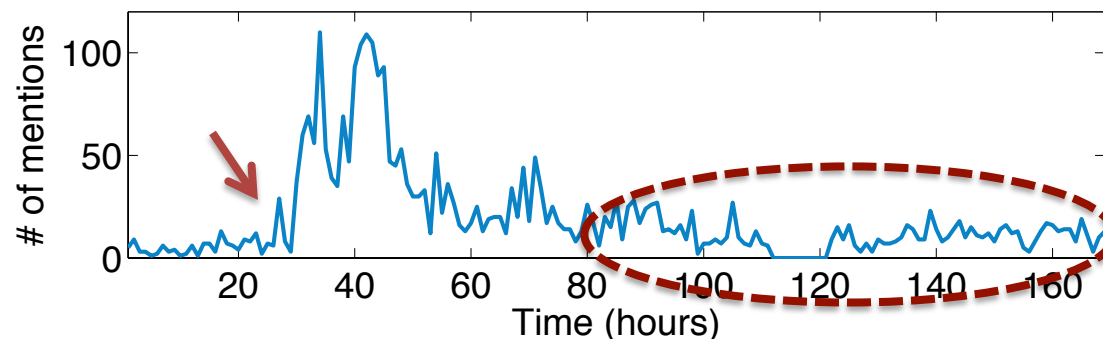


- Short phrases sourced from U.S. politics in 2008

“you can put lipstick on a pig” (# of mentions in blogs)



“yes we can”

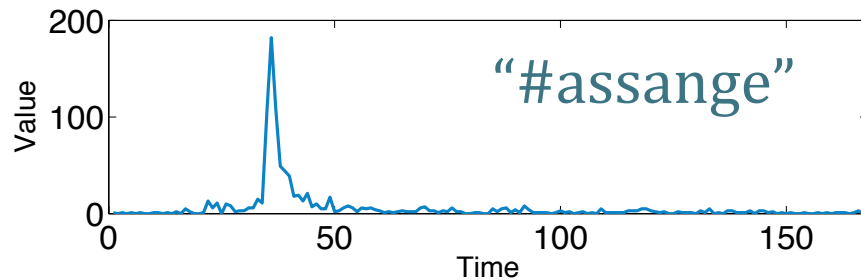




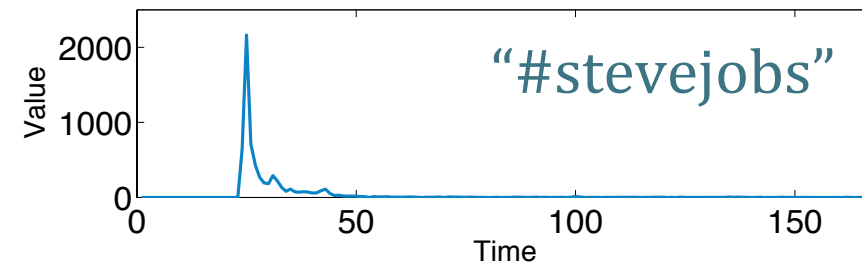
News spread in social media



- Twitter (# of hashtags per hour)



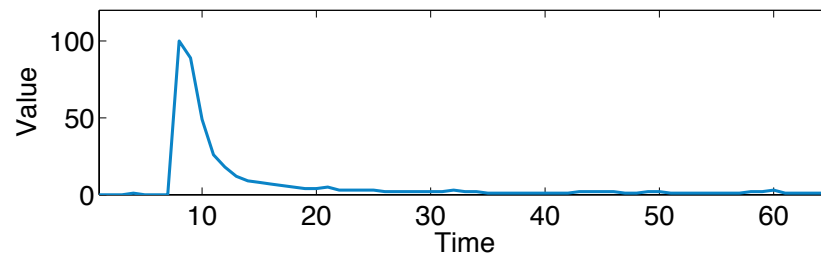
(per hour, 1week)



(per hour, 1 week)

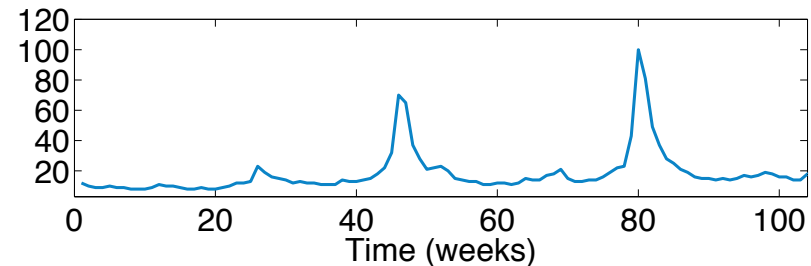
- Google trend (# of queries per week)

“tsunami” (in 2005)



(per week, 1 year)

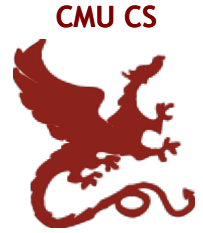
“harry potter” (2010 - 2011)



(per week, 2 years)



News spread in social media



Q. How many patterns are there?

– Four classes on YouTube, etc.

[Crane et al. PNAS'08]

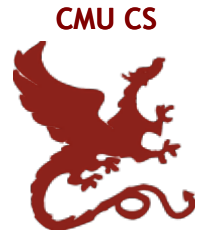
– Six classes on Social media

[Yang et al. WSDM'11]



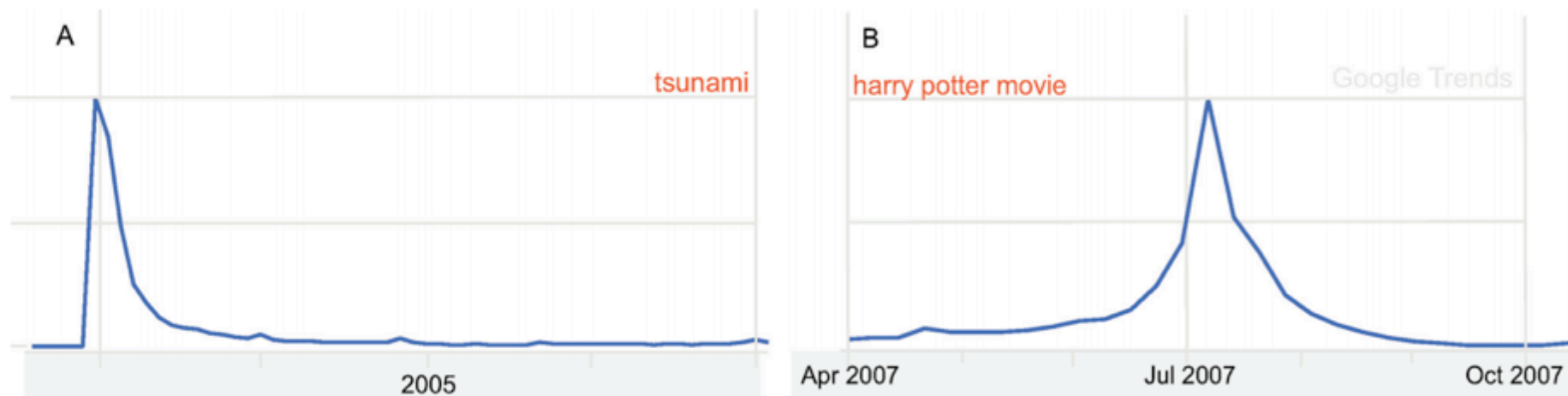


News spread in social media



[Crane et al. PNAS'08]

- The volume of Google searches



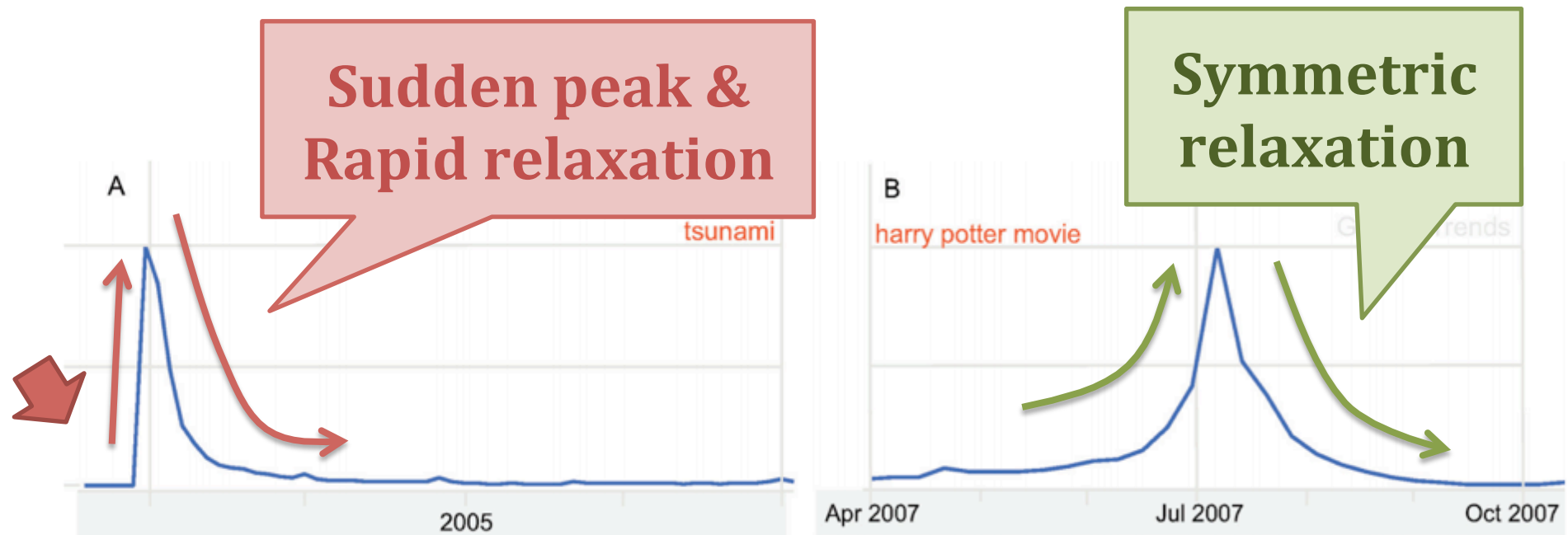
“Tsunami”

“Harry Potter movie”

News spread in social media

[Crane et al. PNAS'08]

- The volume of Google searches



“Tsunami”
(Exogenous)

“Harry Potter movie”
(Endogenous)



News spread in social media



[Crane et al. PNAS'08]

- Based on self-excited Hawkes Poisson process*

$$\frac{dB(t)}{dt} = S(t) + \sum_{i, t_i \leq t} \mu_i \cdot \phi(t - t_i)$$

*[Hawkes+ 1974]



News spread in social media



[Crane et al. PNAS'08]

- Based on self-excited Hawkes Poisson process*

$$\frac{dB(t)}{dt} = S(t) + \sum_{i, t_i \leq t} \mu_i \cdot \phi(t - t_i)$$

Rate of
spread of
infection/
propagation

Exogenous
/External
source

of
Potential
viewers

Decaying
virus/news
strength

*[Hawkes+ 1974]



News spread in social media



[Crane et al. PNAS'08]

- Based on self-excited Hawkes Poisson process*

$$\frac{dB(t)}{dt} = S(t) + \sum_{i, t_i \leq t} \mu_i \cdot \phi(t - t_i)$$

Rate of

Excitement

of

initial

users

Decaying
virus/news
strength
(Power law)

$$\phi(t) \sim \frac{1}{t^{1+\theta}} \quad (0 < \theta < 1)$$

propagation

*[Hawkes+ 1974]



News spread in social media

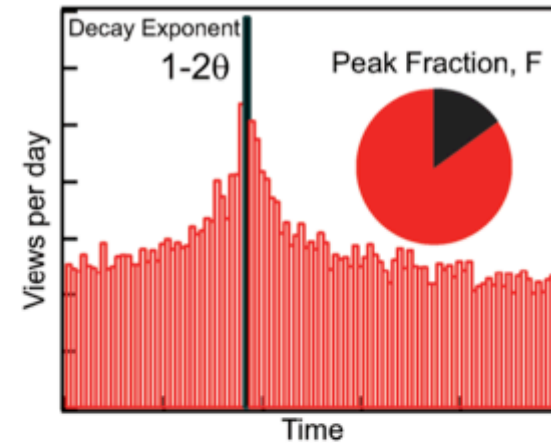
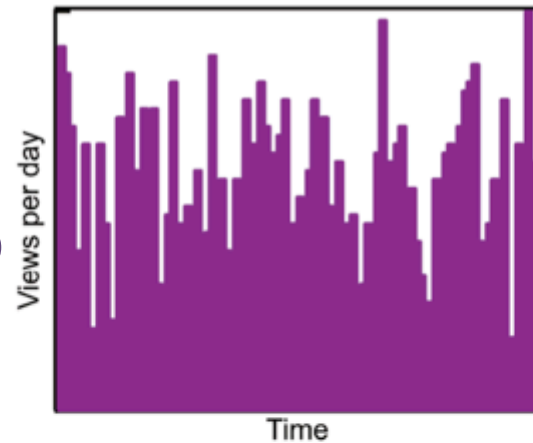


[Crane et al. PNAS'08]

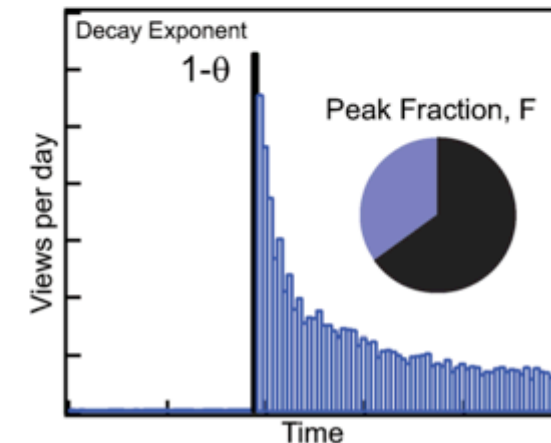
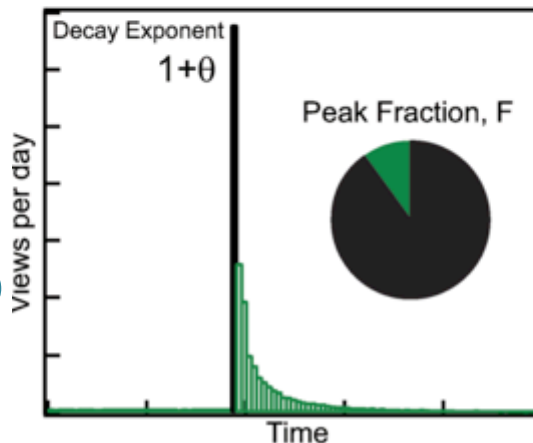
- Four classes on YouTube
Sub-Critical

Critical

Endogenous



Exogenous



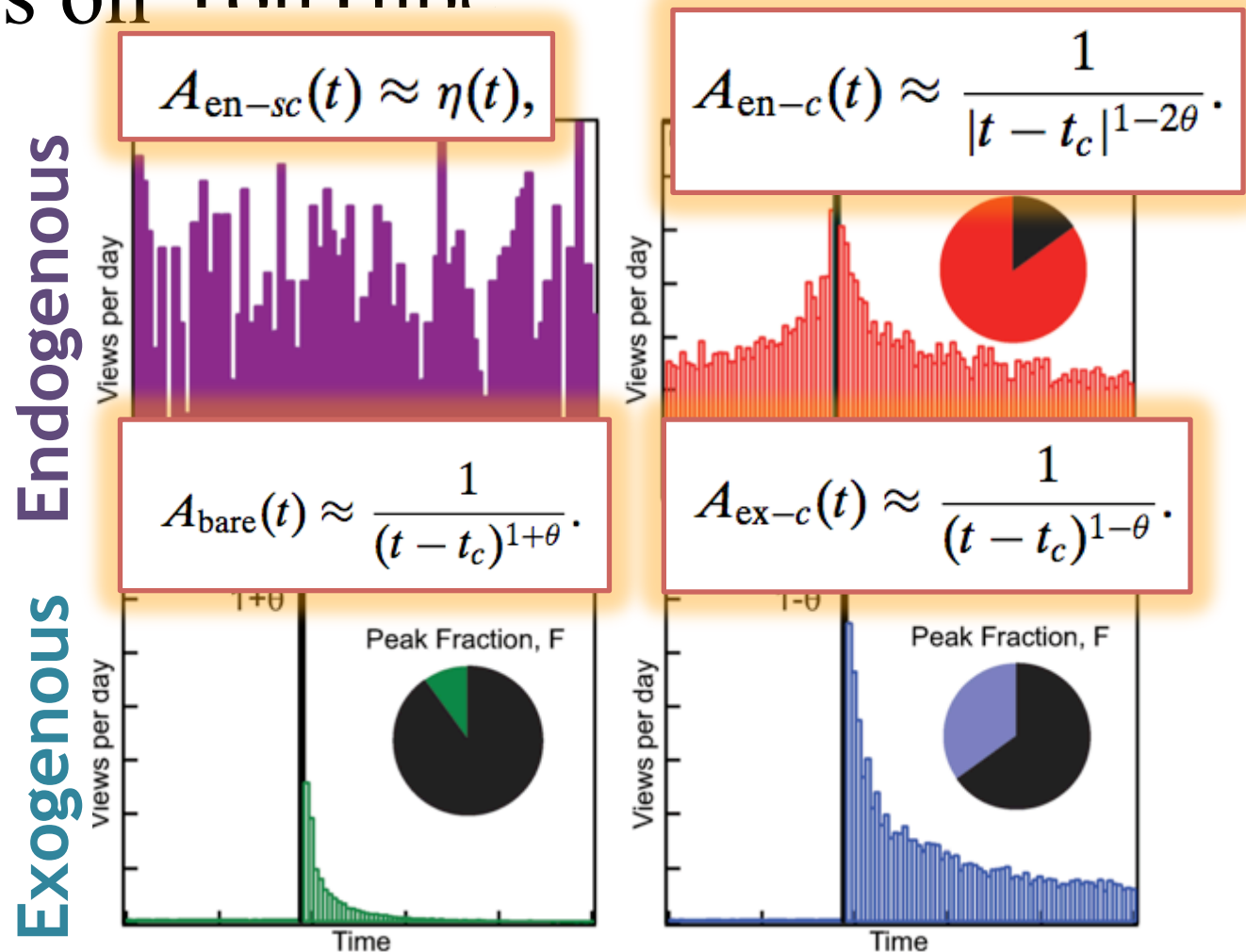


News spread in social media



- Four classes on YouTube

[Crane et al. PNAS'08]



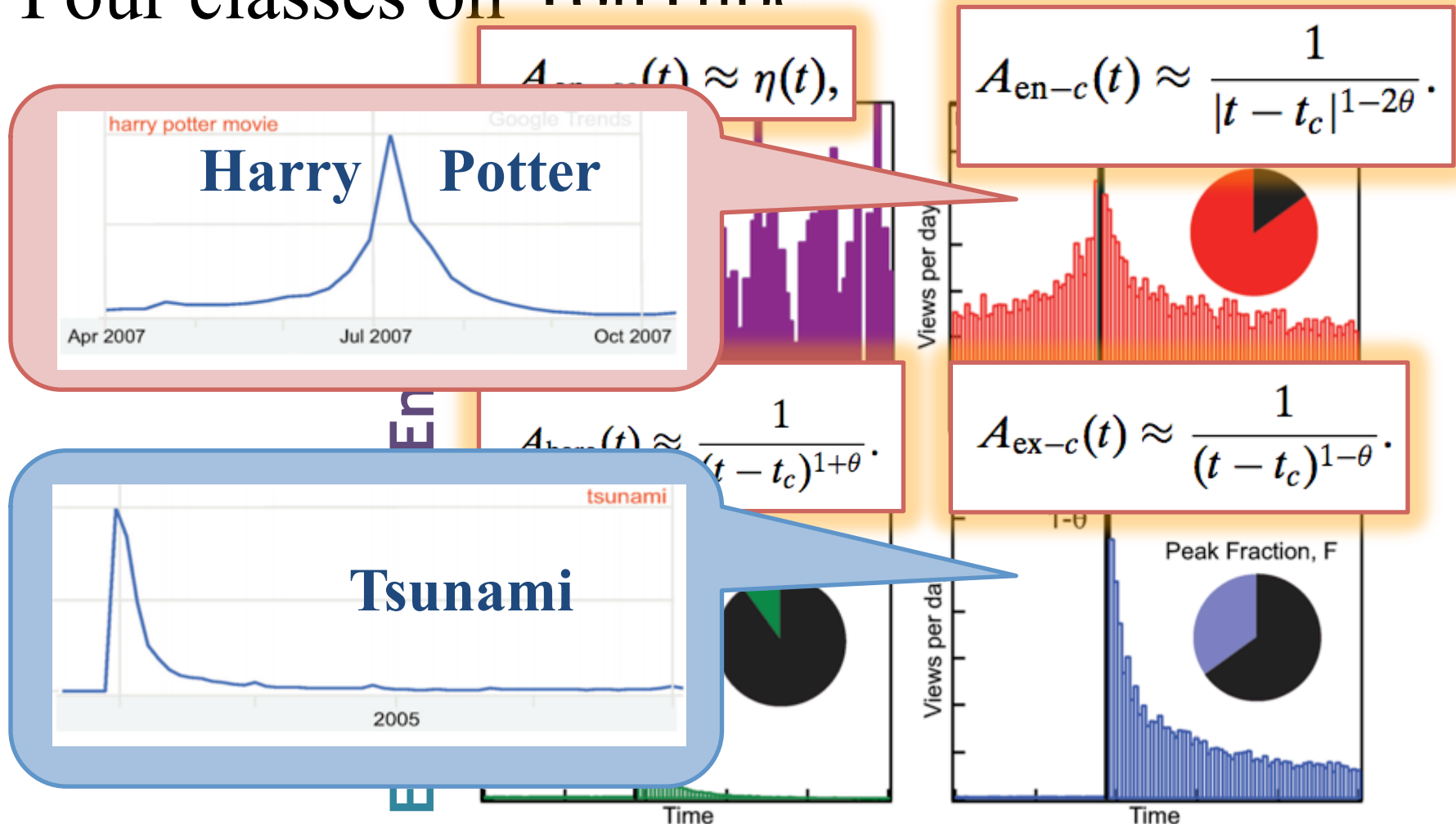


News spread in social media



- Four classes on YouTube

[Crane et al. PNAS'08]

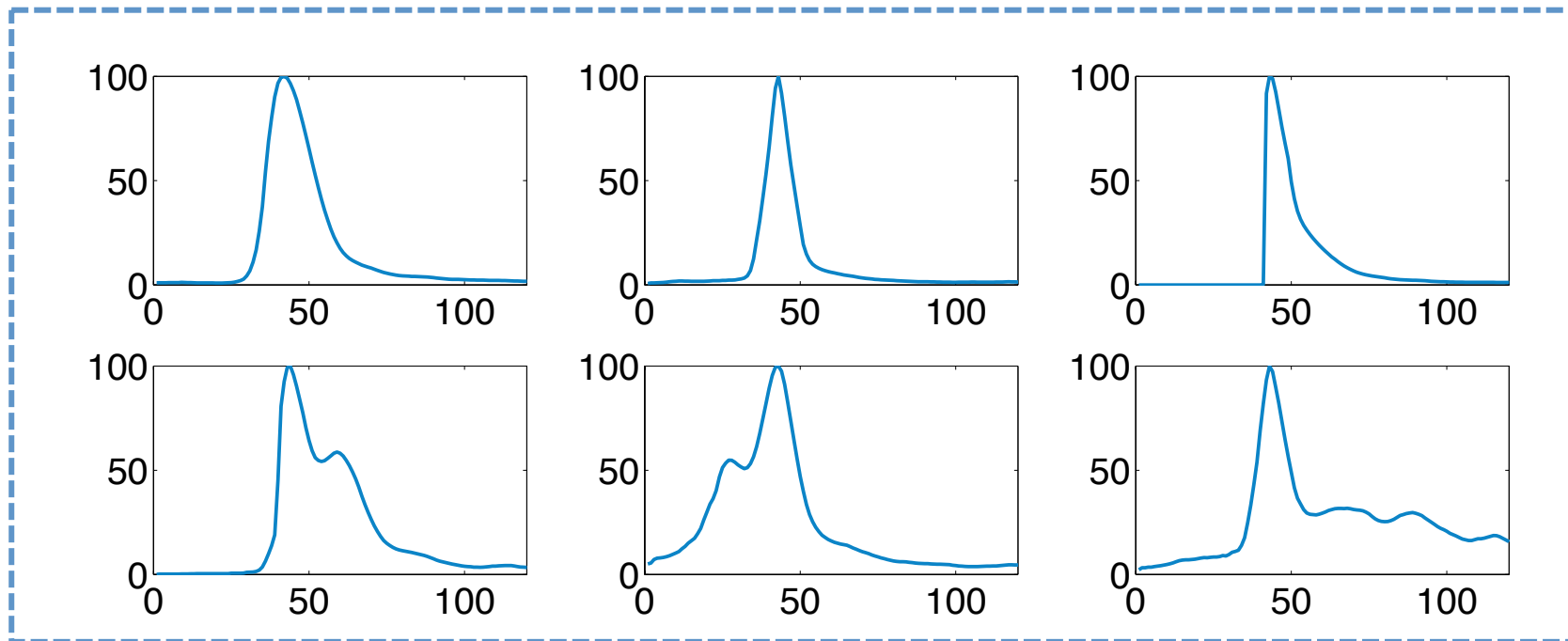




News spread in social media



- Six classes of information diffusion patterns on social media [Yang et al. WSDM'11]

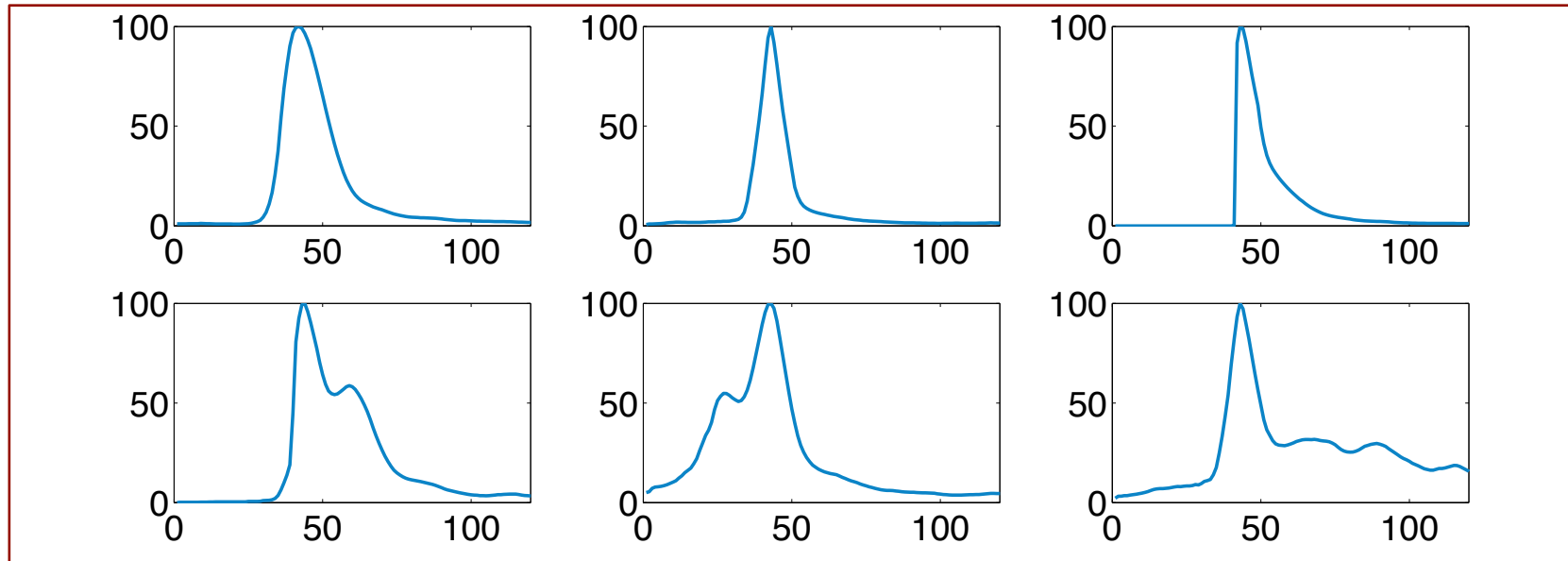
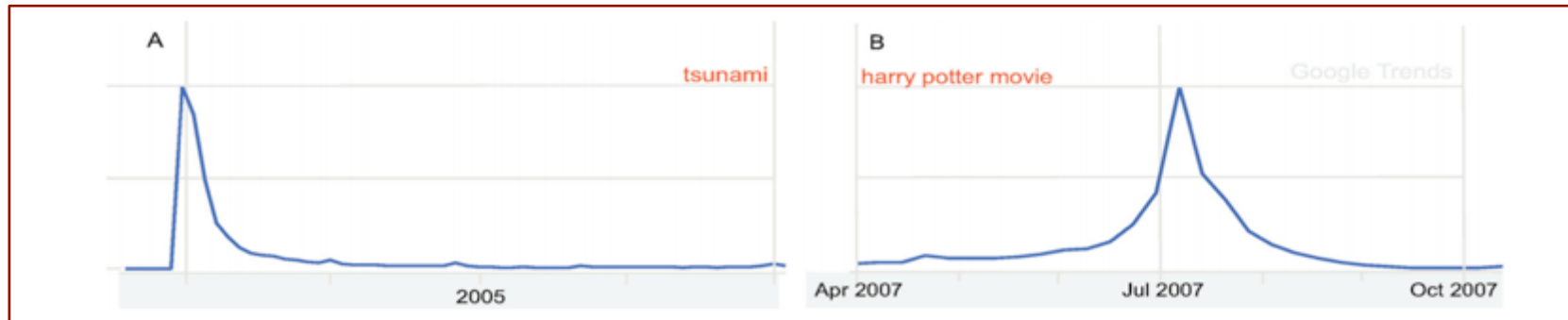




News spread in social media



Q. How many patterns are there, after all?





News spread in social media



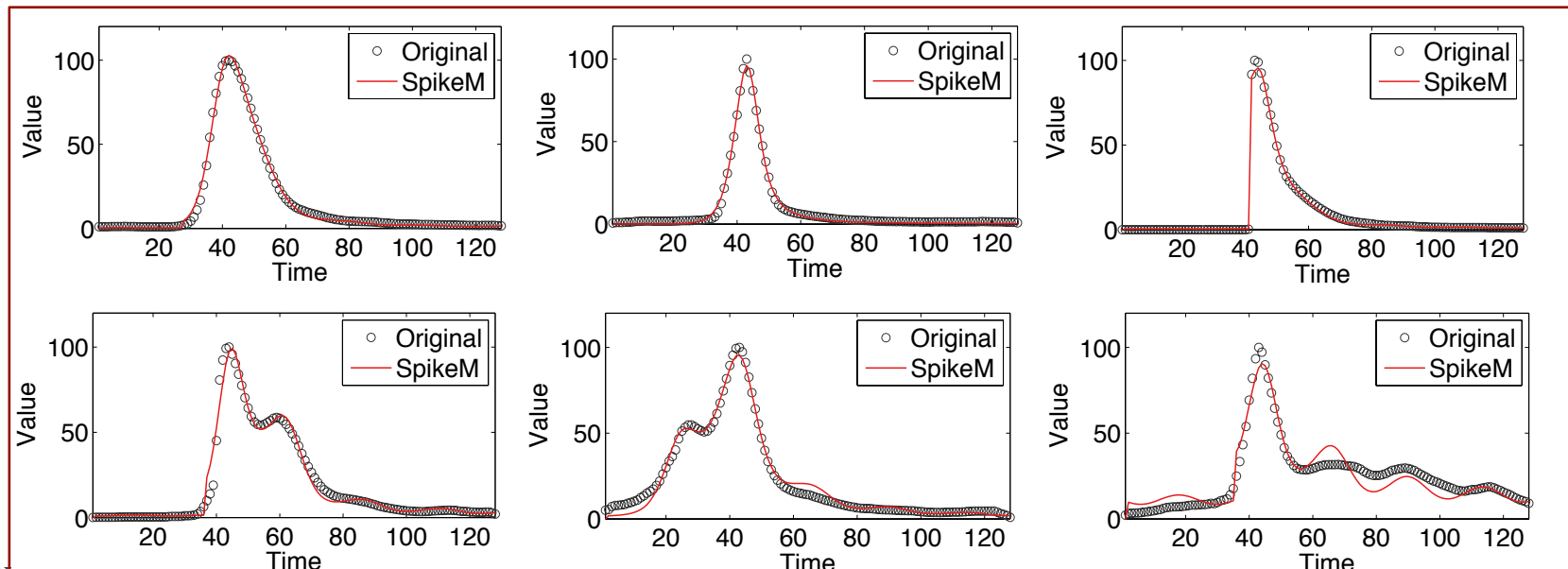
A. Our answer is “ONE”!



A single non-linear model !



“SpikeM”





[Matsubara+ KDD'12]

Rise and Fall Patterns of Information Diffusion: Model and Implications

Yasuko Matsubara (Kyoto University),



Yasushi Sakurai (NTT),



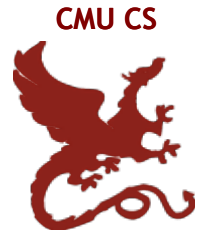
B. Aditya Prakash (CMU),

Lei Li (UCB), Christos Faloutsos (CMU)



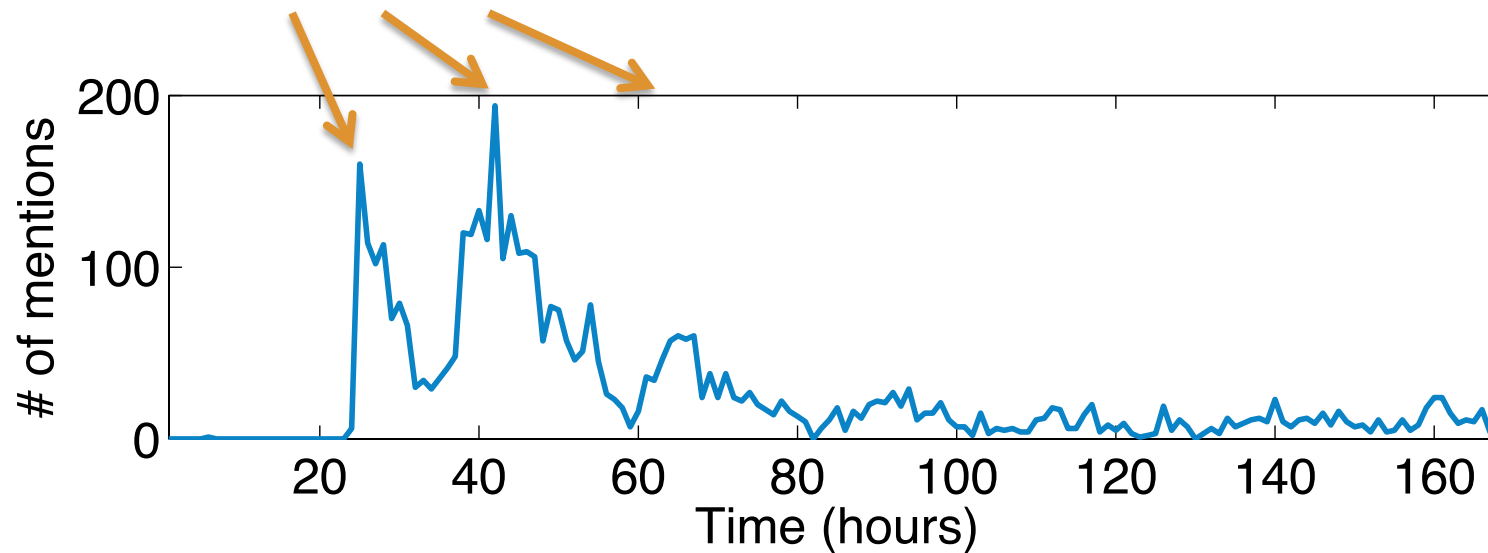


Rise and fall patterns in social media



SpikeM captures 3 properties of real spike

1. periodicities



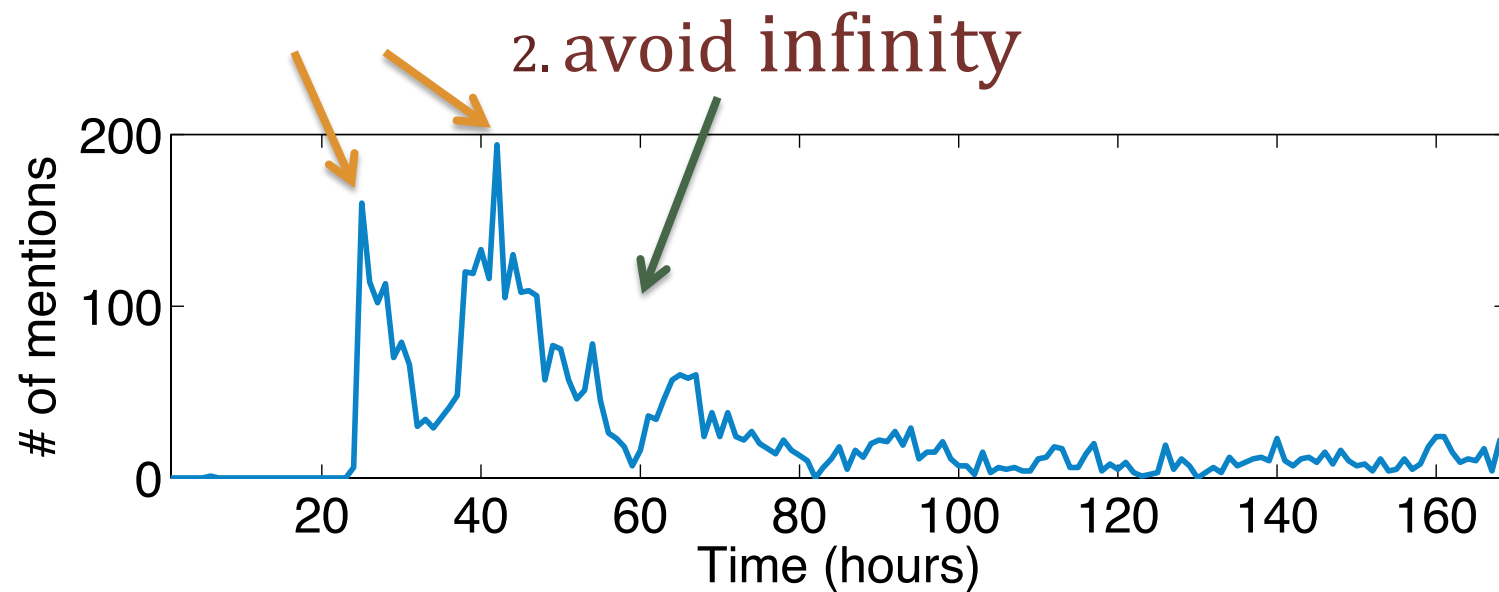


Rise and fall patterns in social media



SpikeM captures 3 properties of real spike

1. periodicities





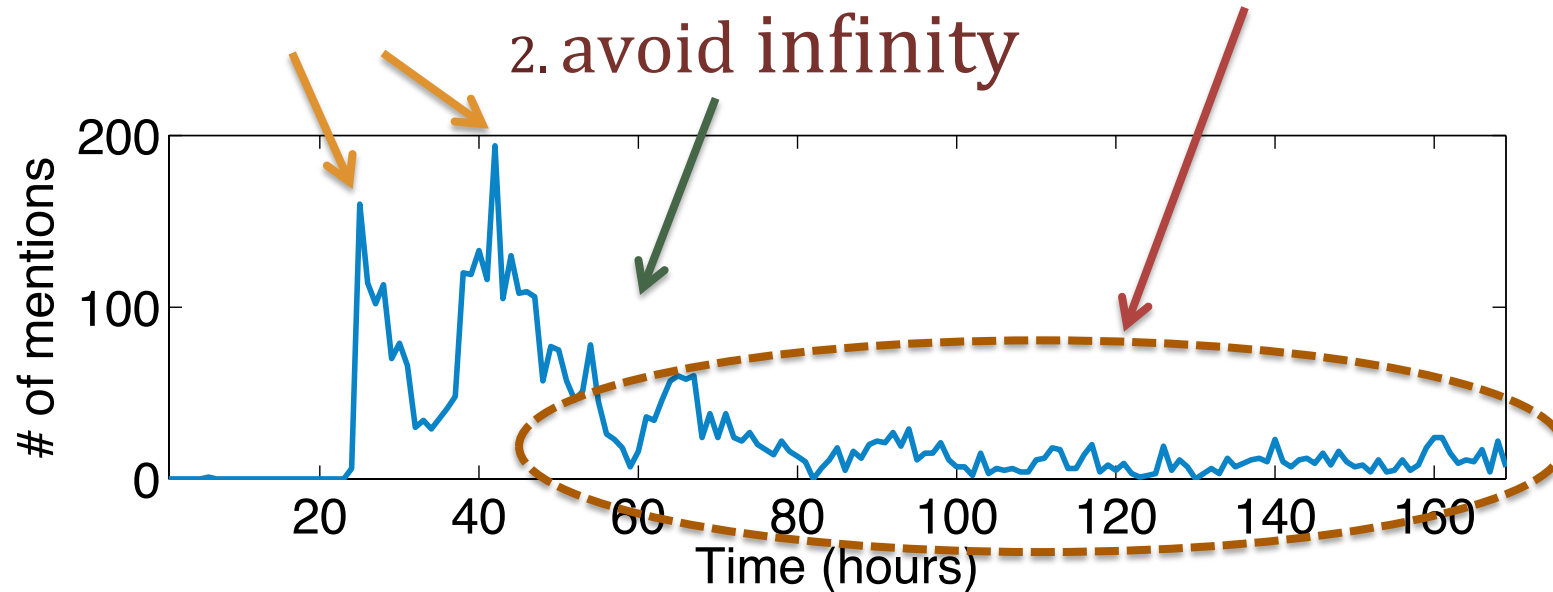
Rise and fall patterns in social media



SpikeM captures 3 properties of real spike

1. periodicities

3. power-law fall





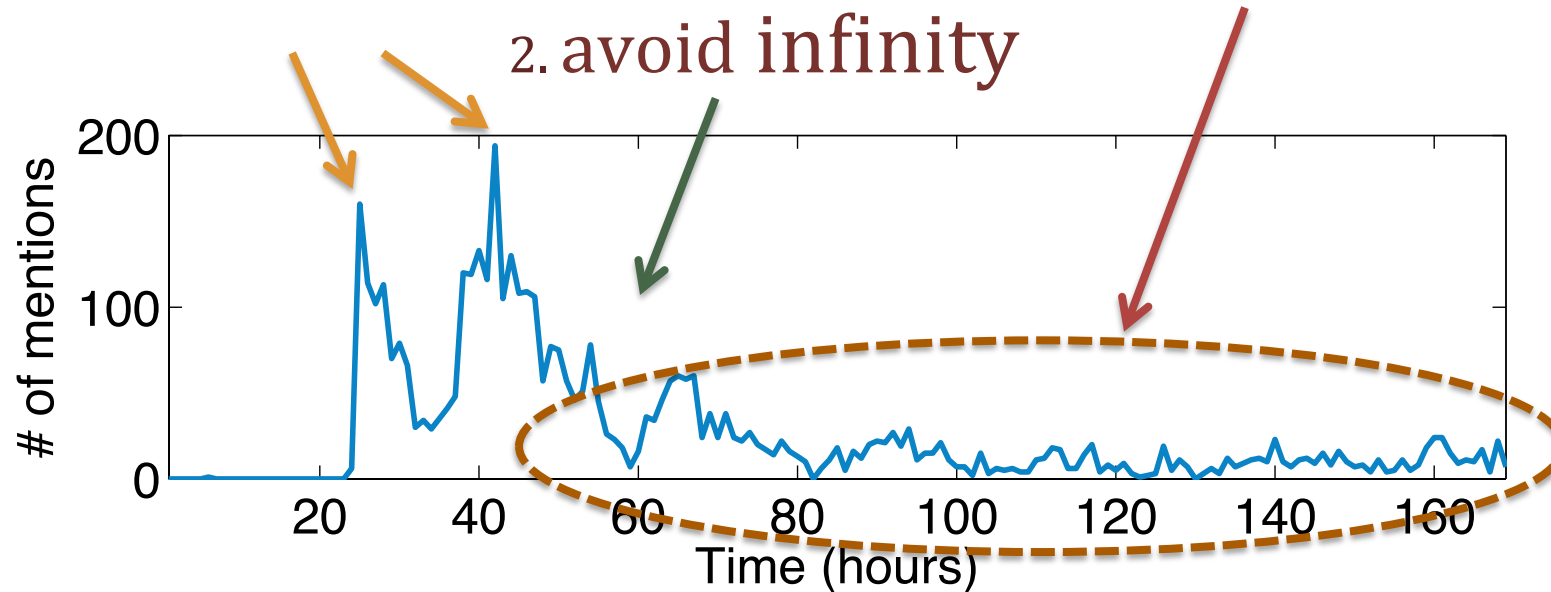
Rise and fall patterns in social media



SpikeM captures 3 properties of real spike

1. periodicities

3. power-law fall

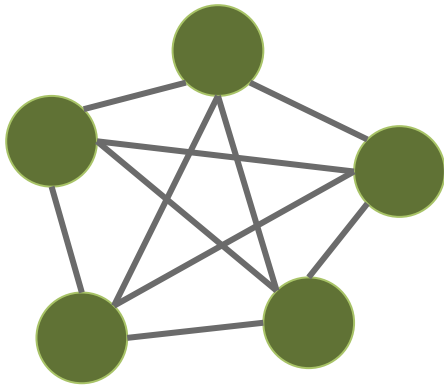


SpikeM can capture behavior of real spikes
using few parameters



Main idea (details)

- 1. **Un-informed bloggers** (clique of N bloggers/nodes)



Time n=0

Nodes (bloggers) consist of two states



– **U**n-informed of rumor

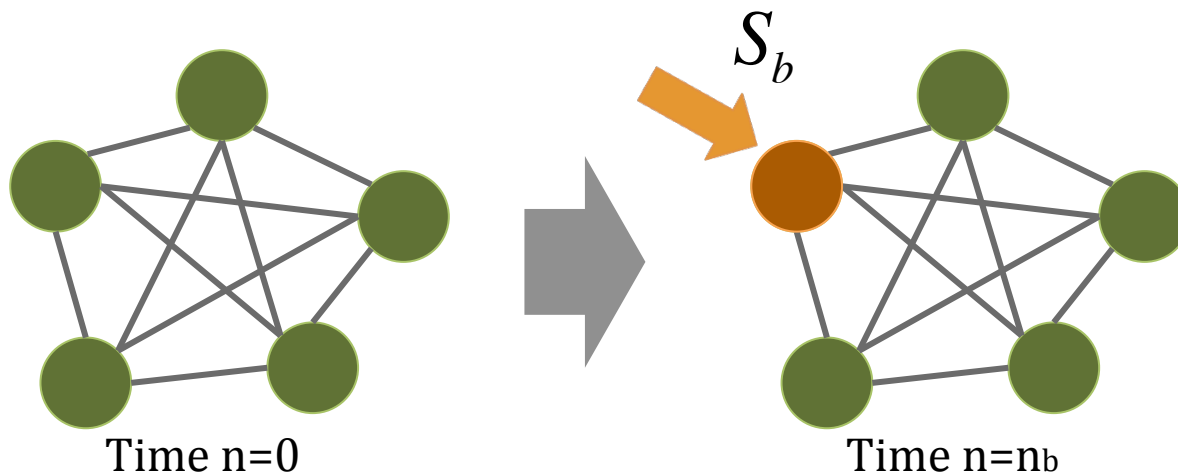


– informed, and **B**logged about rumor



Main idea (details)

- 1. **Un-informed bloggers** (clique of N bloggers/nodes)
- 2. **External shock** at time n_b (e.g, breaking news)



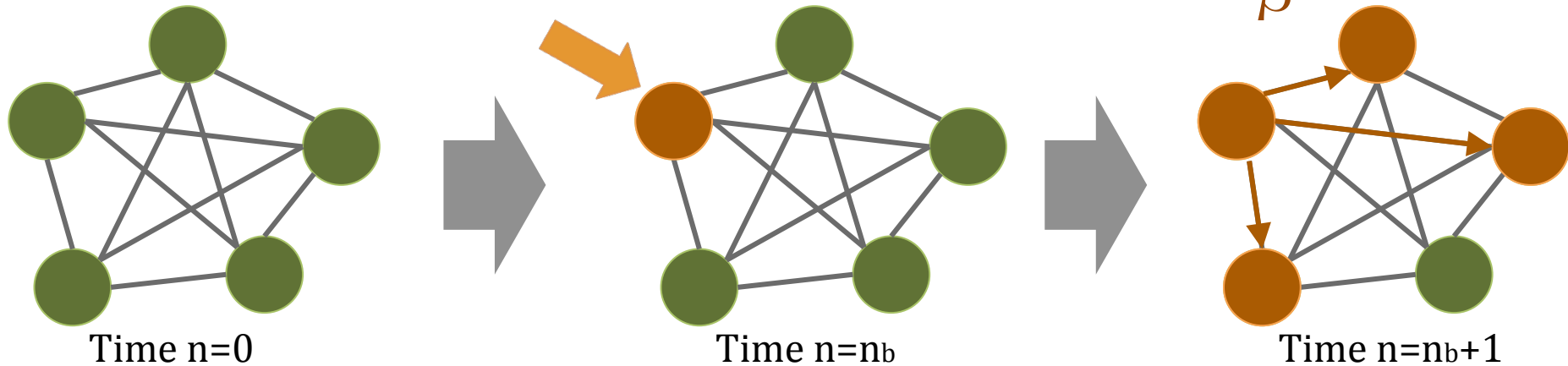
External shock

- Event happened at time n_b
- S_b bloggers are informed, blog about news



Main idea (details)

- 1. **Un-informed bloggers** (clique of N bloggers/nodes)
- 2. **External shock** at time n_b (e.g, breaking news)
- 3. **Infection** (word-of-mouth effects)



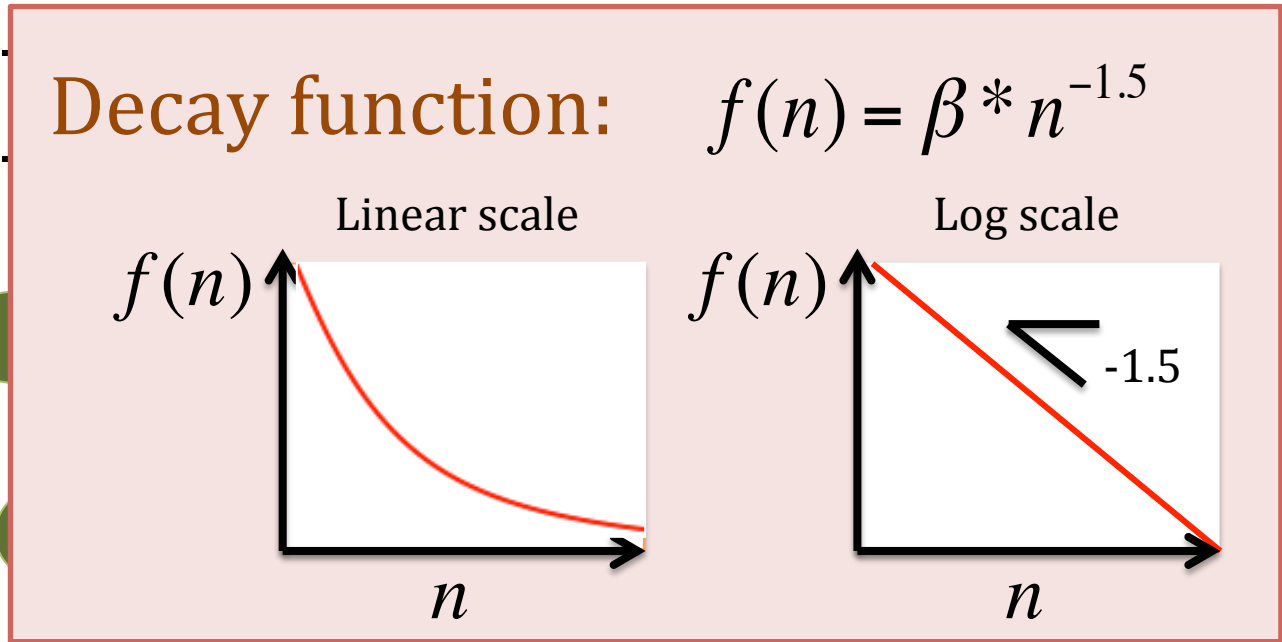
Infectiveness of a blog-post

β – Strength of infection (quality of news)

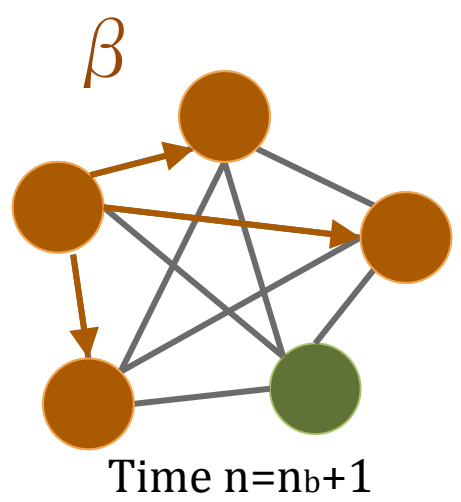
$f(n)$ – Decay function (how infective a blog posting is)

Main idea (details)

- 1. **Un-informed bloggers** (clique of N bloggers/nodes)



news)



Infectiveness of a blog-post

β – Strength of infection (quality of news)

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



SpikeM-base (details)

Equations of SpikeM (base)

$$\underline{\Delta B(n+1)} = U(n) \cdot \sum_{t=n_b}^n (\Delta B(t) + S(t)) \cdot f(n+1-t) + \varepsilon$$

Blogged

$$\underline{U(n+1)} = U(n) - \Delta B(n+1)$$

Un-informed

- N – Total population of available bloggers
- β – Strength of infection/news
- n_b, S_b – External shock S_b at birth (time n_b)
- ε – Background noise

SpikeM - periodicity

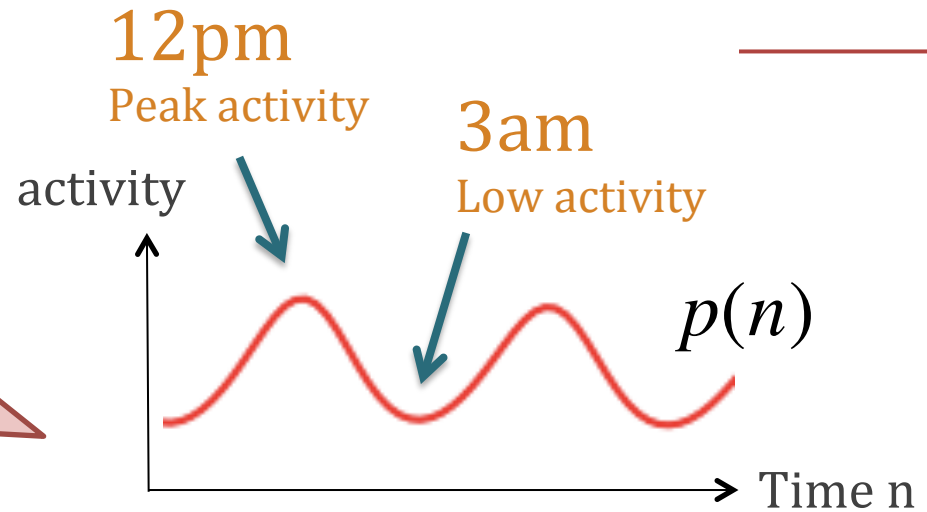
Full equation of SpikeM

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

$$U(n+1) = U(n) - \Delta B(n+1)$$

Un-informed

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





Model fitting (Details)

- SpikeM consists of 7 parameters

$$\theta = \{N, \beta, n_b, S_b, \varepsilon, P_a, P_s\}$$

Learning parameters

- Given a real time sequence

$$X = \{X(1), \dots, X(n), \dots, X(n_d)\}$$

- Minimize the error

(Levenberg-Marquardt (LM) fitting)

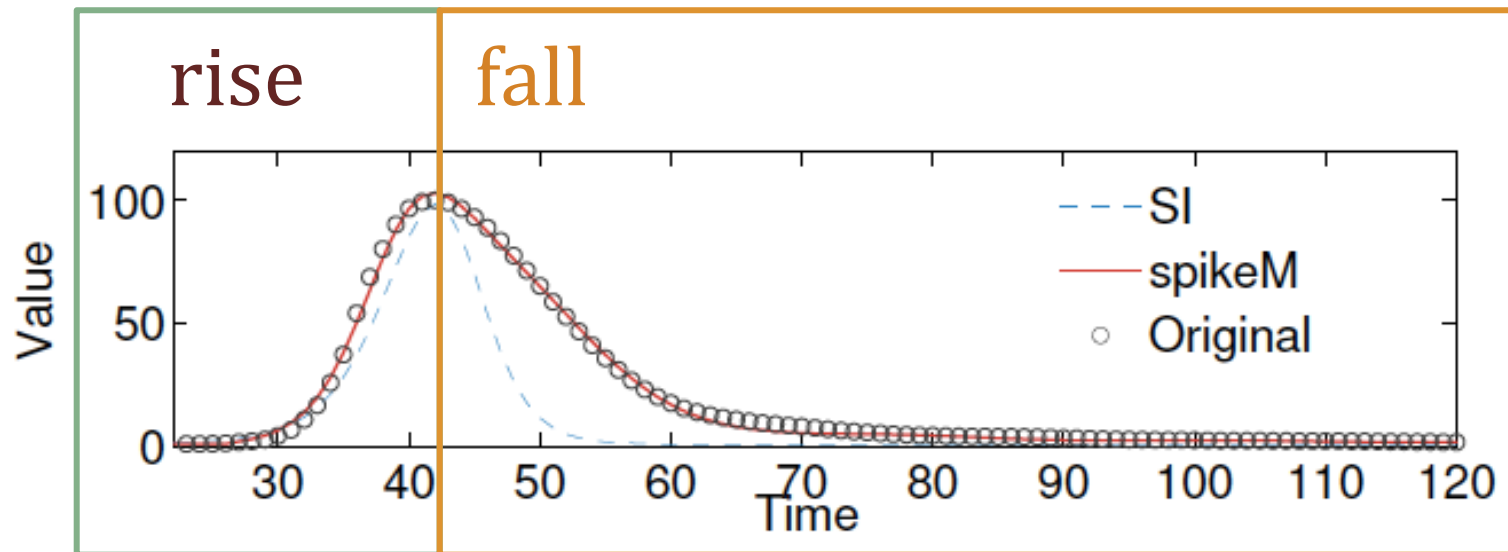
$$D(X, \theta) = \sum_{n=1}^{n_d} (X(n) - \Delta B(n))^2$$



Analysis

SpikeM matches reality

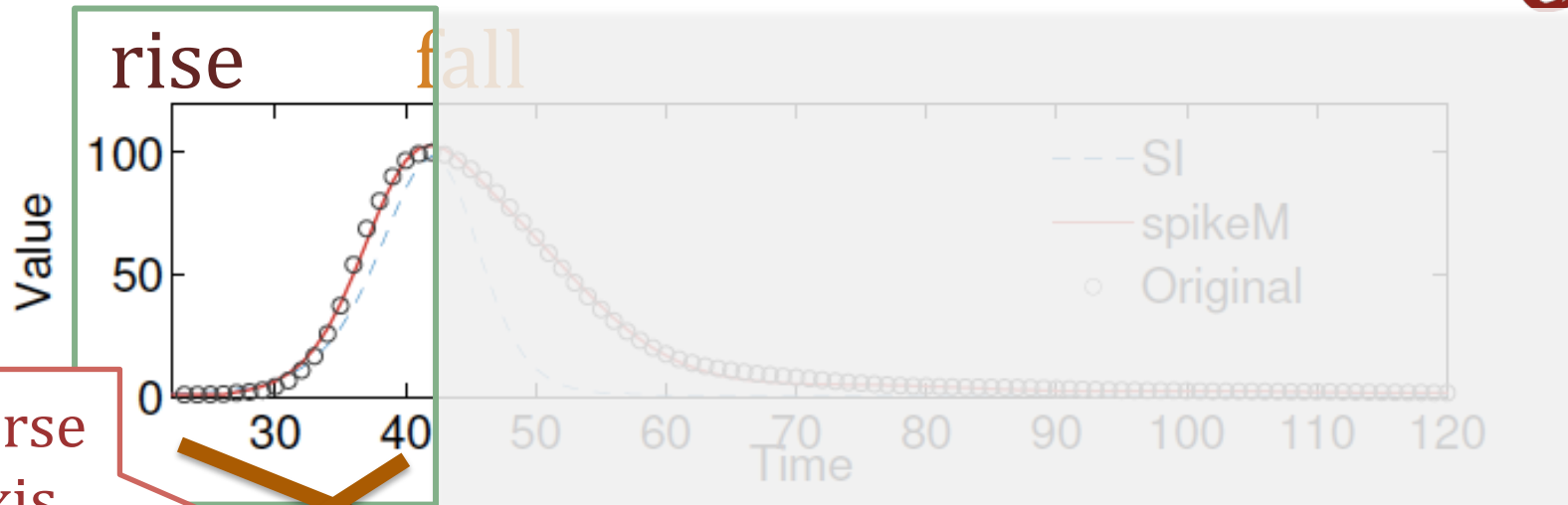
exponential rise and power-law fall



SpikeM vs. **SI model** (susceptible infected model)

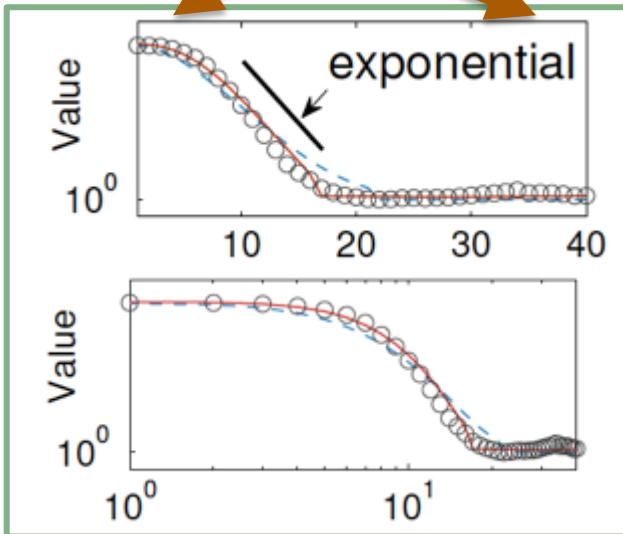


Analysis



Reverse x-axis

Linear-log



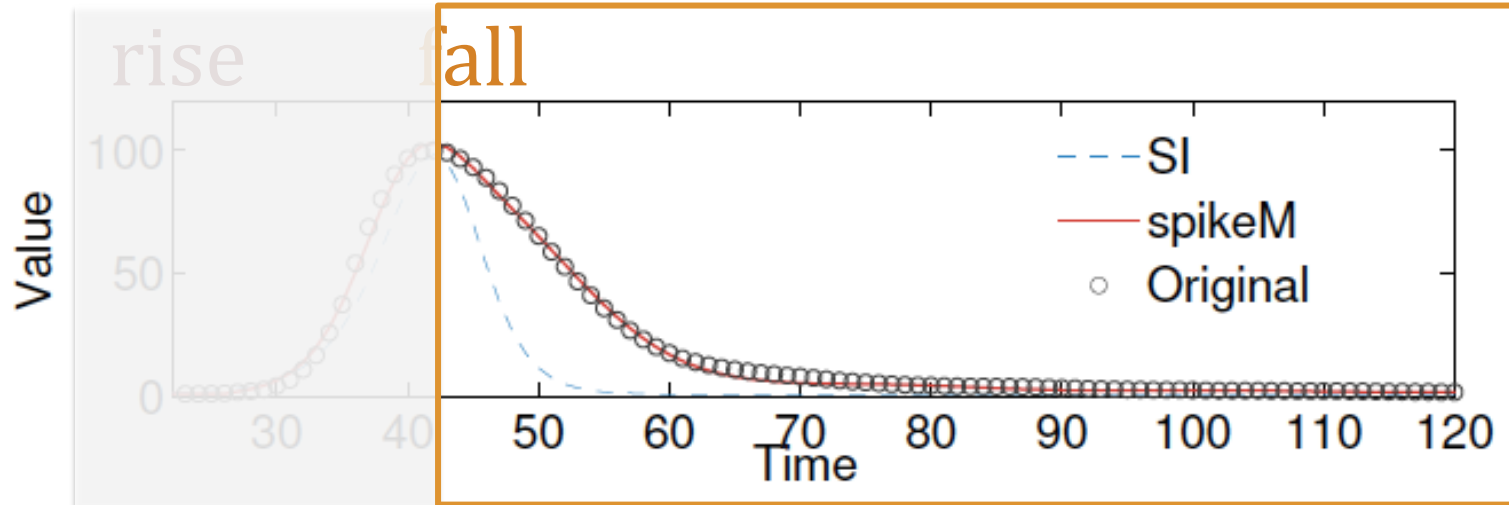
Log-log

Rise-part

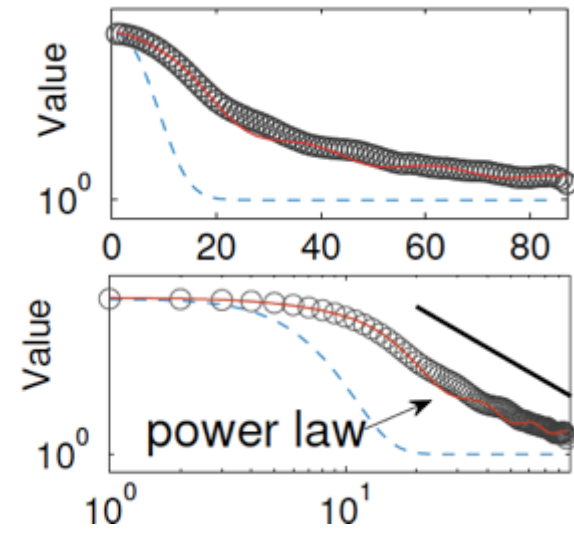
SpikeM: exponential
SI model: exponential



Analysis



Fall-part
SpikeM: power law
SI model: exponential
SpikeM matches reality



Linear-log

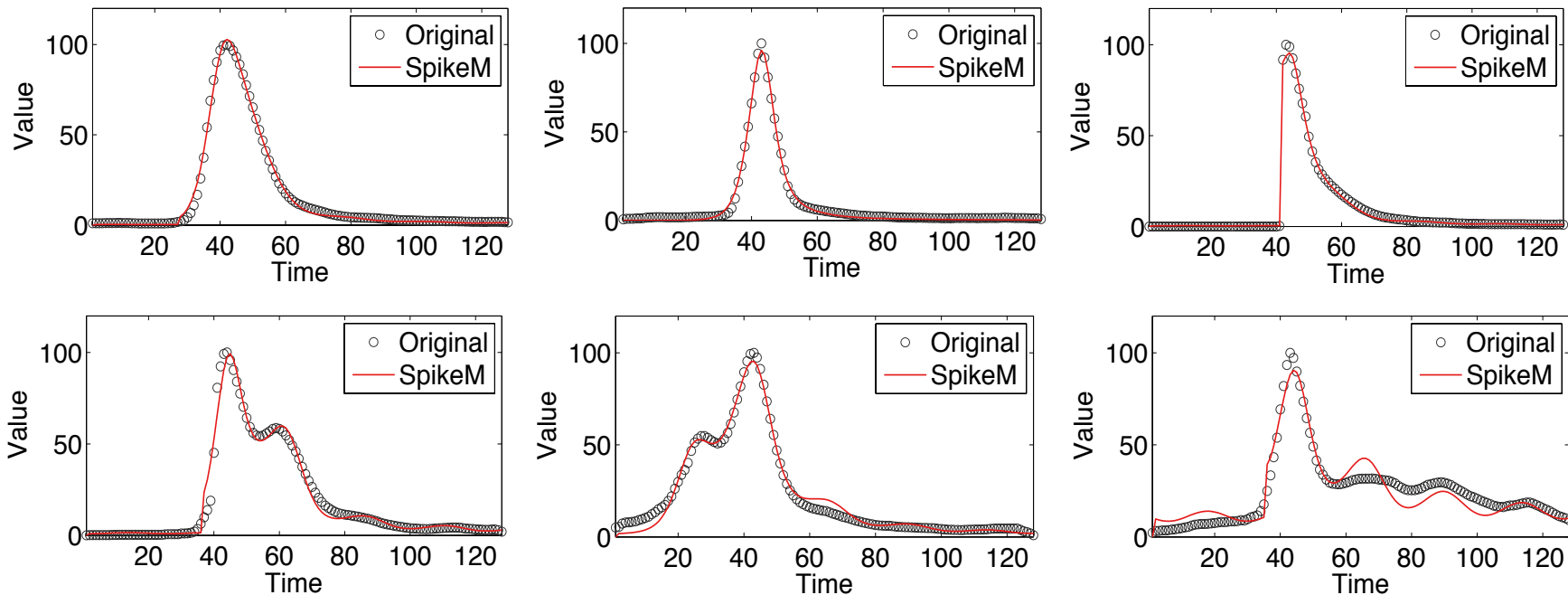
Log-log



Q1-1 Explaining K-SC clusters



–Six patterns of K-SC [Yang et al. WSDM'11]



- **SpikeM** can generate all patterns in K-SC



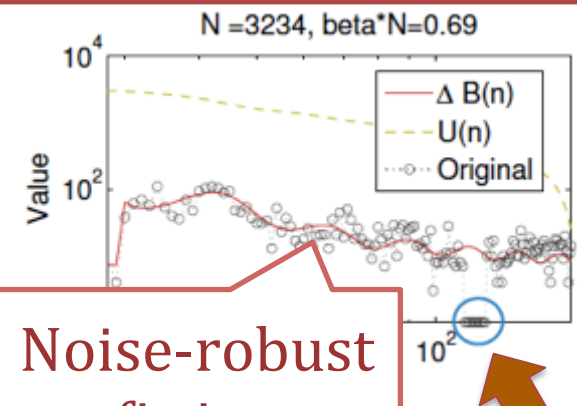
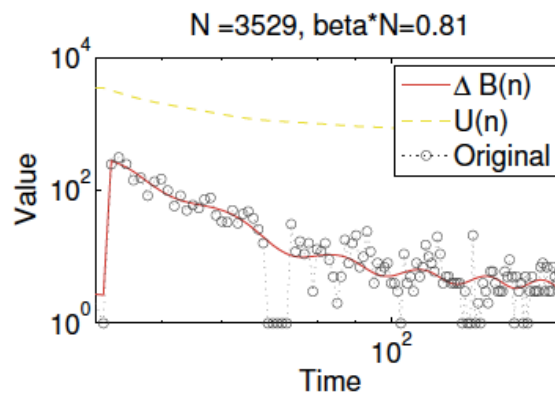
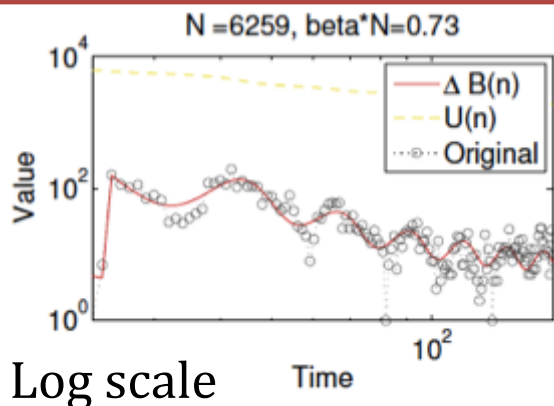
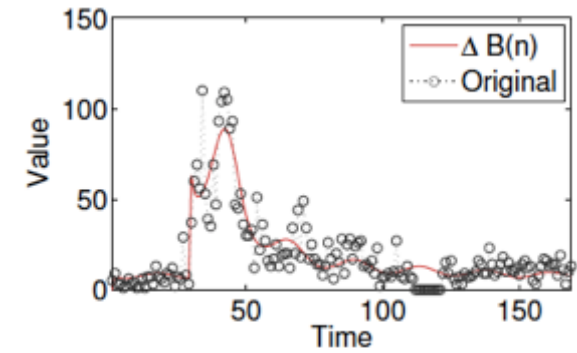
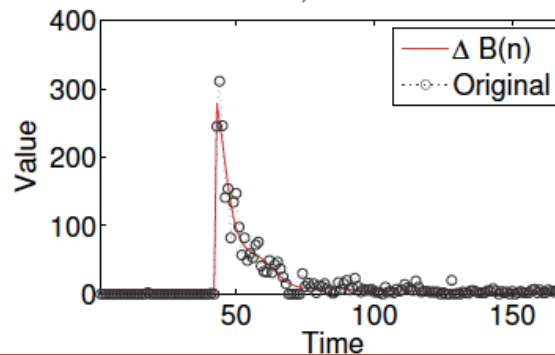
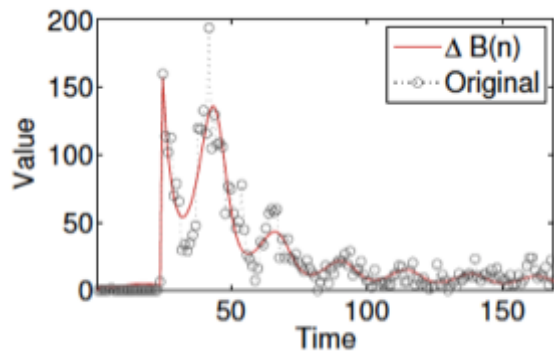
Q1-2 Matching

MemeTracker patterns



MemeTracker (memes in blogs) [Leskovec et al. KDD'09]

Linear scale



Noise-robust fitting

Outliers

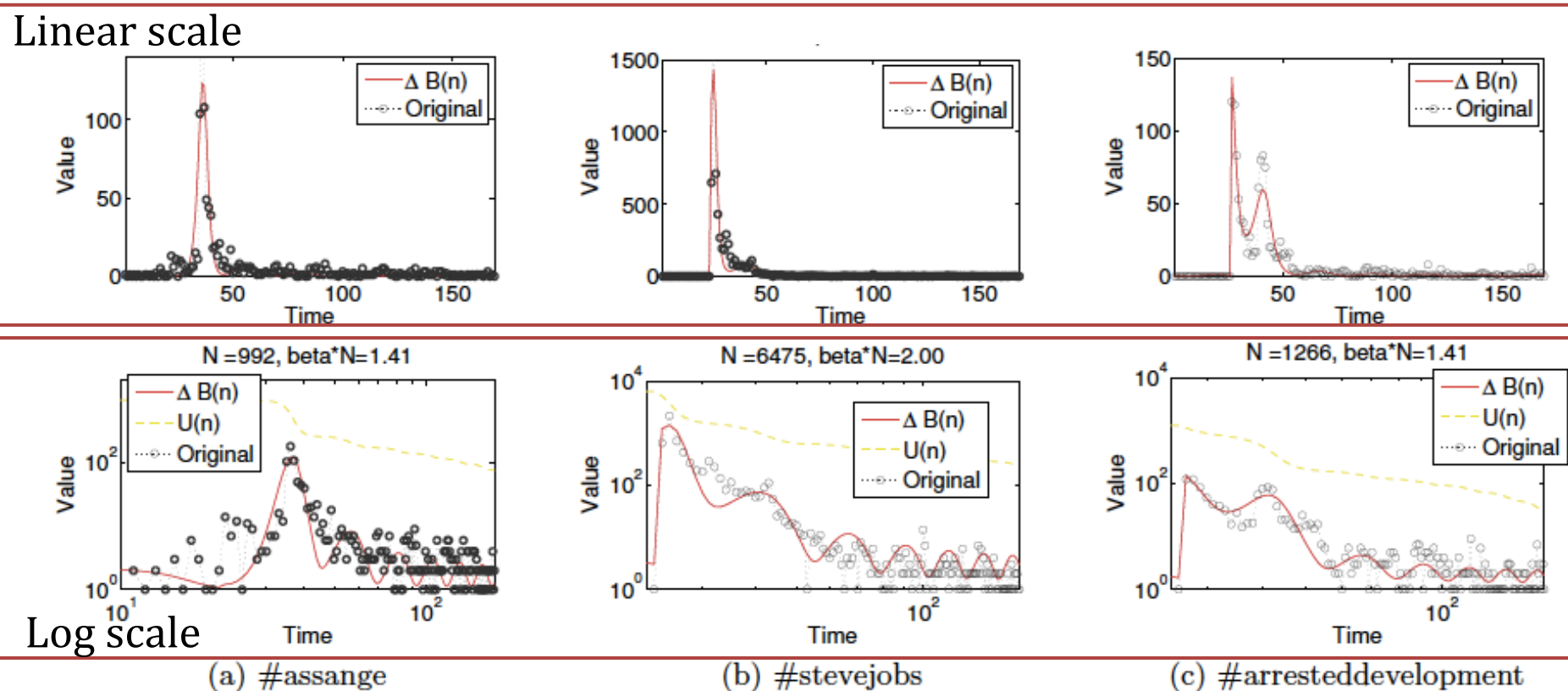
SpikeM can fit various patterns in blog



Q1-3 Matching Twitter data



Twitter data (hashtags)



It can generate various patterns in social media

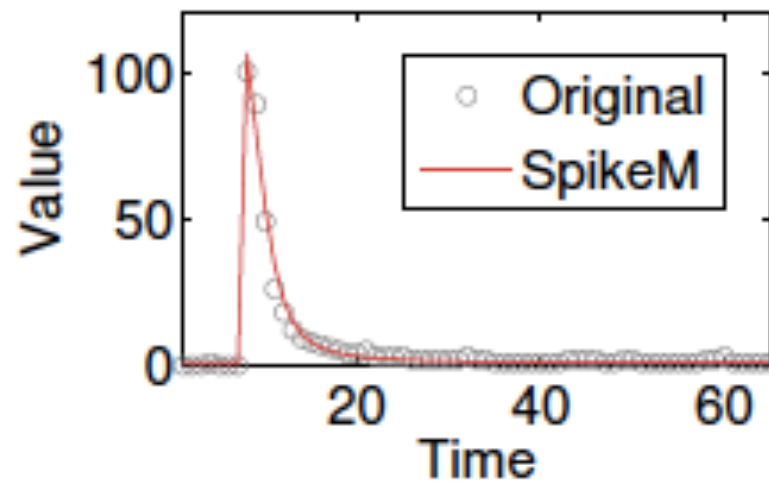


Q1-4 Matching

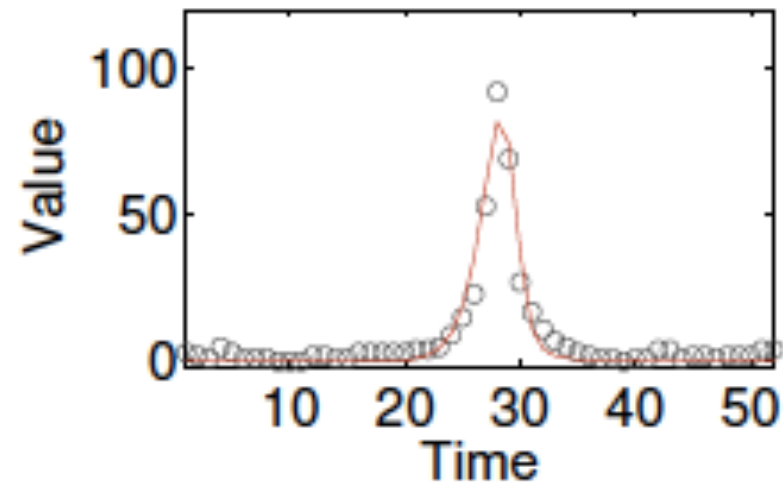
Google trend data



Volume of searches for queries on Google



(a) “tsunami” (2005)



(b) “Harry Potter” (2007)

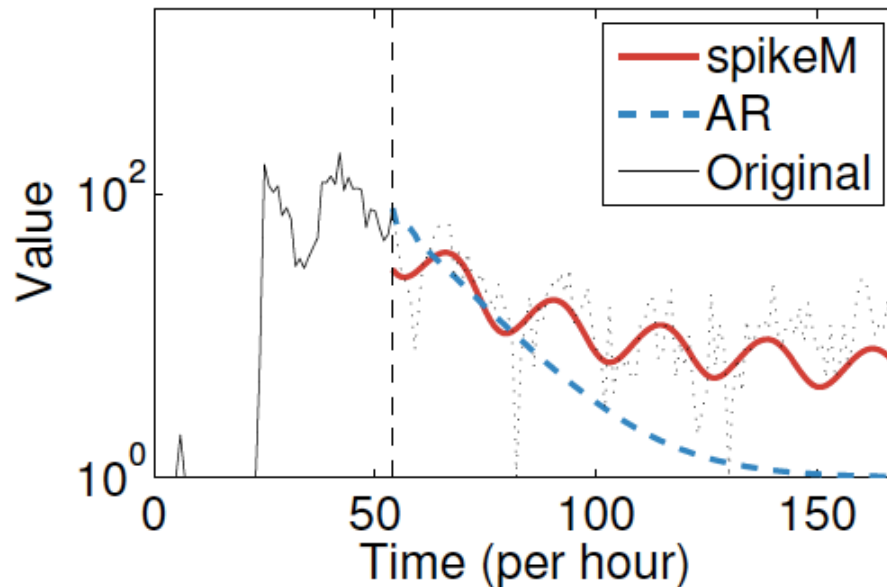
SpikeM can capture various patterns



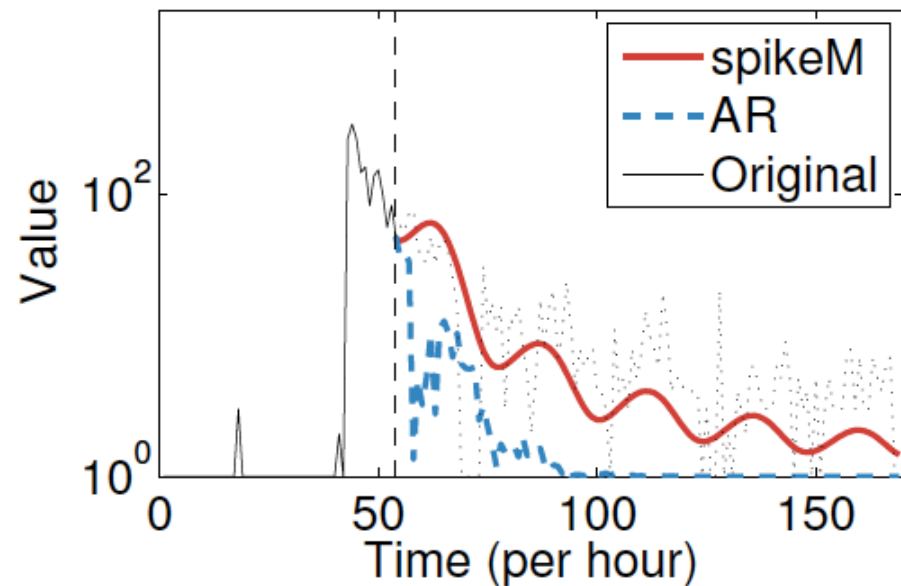
Q2 Tail-part forecasts

- Given a first part of the spike
- forecast the tail part

$N = 5960$, $\beta * N = 0.7$



$N = 3481$, $\beta * N = 1.2$

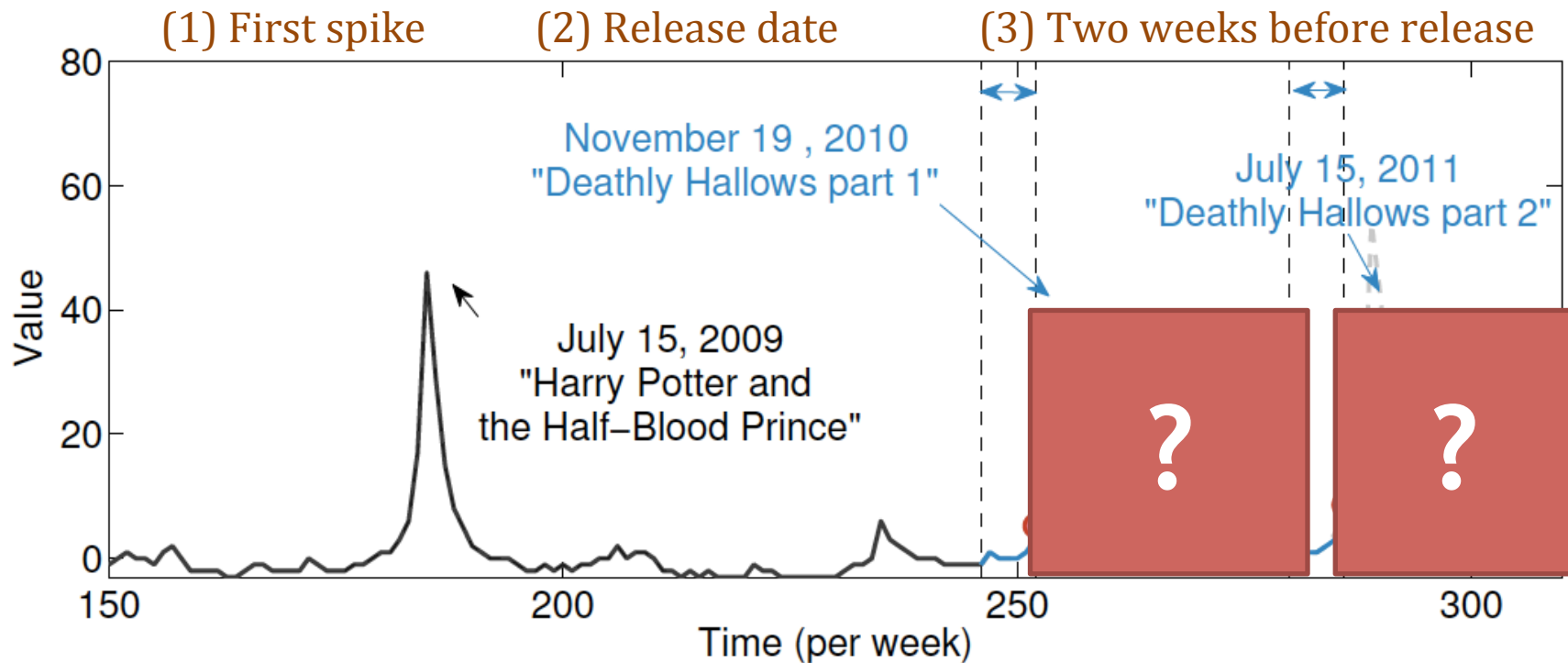


SpikeM can capture tail part (AR: fail)



A1. “What-if” forecasting

Forecast not only tail-part, but also **rise-part!**



e.g., given (1) first spike,

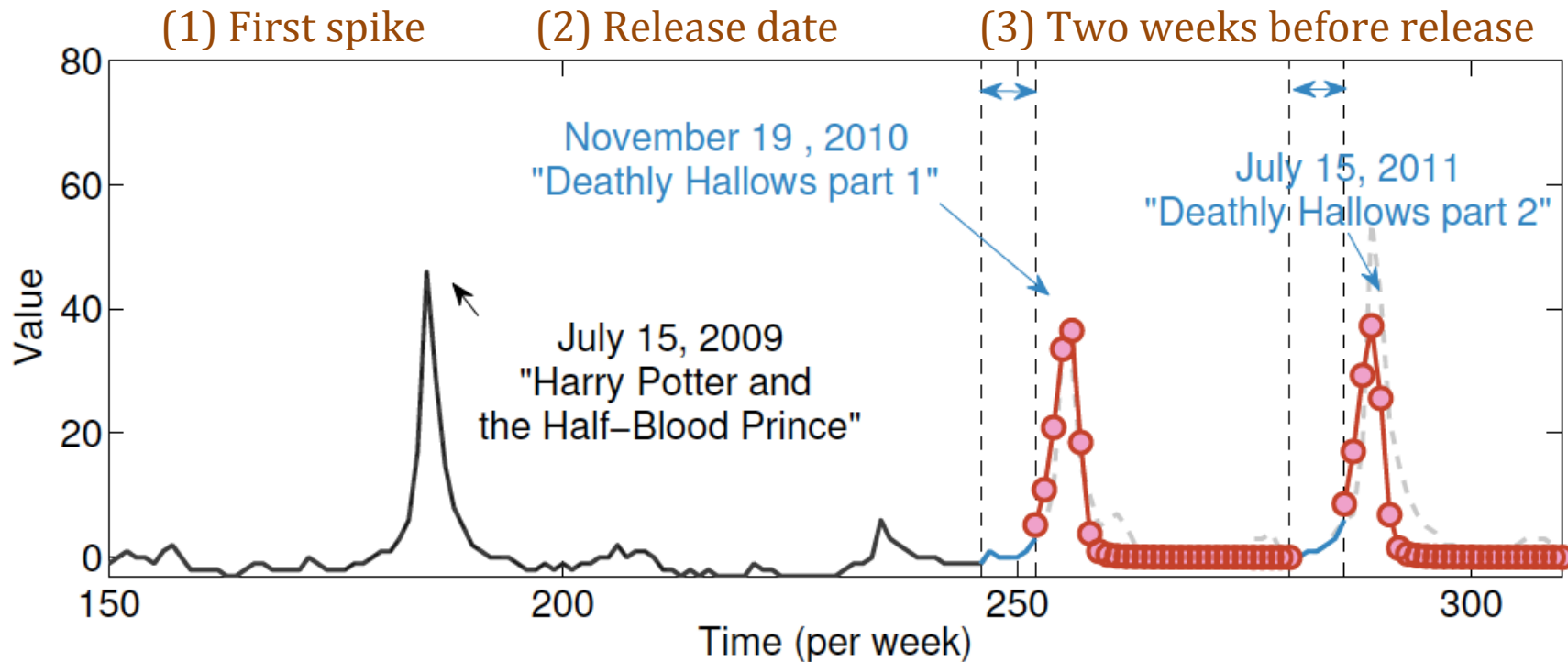
(2) release date of two sequel movies

(3) access volume before the release date



A1. “What-if” forecasting

Forecast not only tail-part, but also **rise-part!**

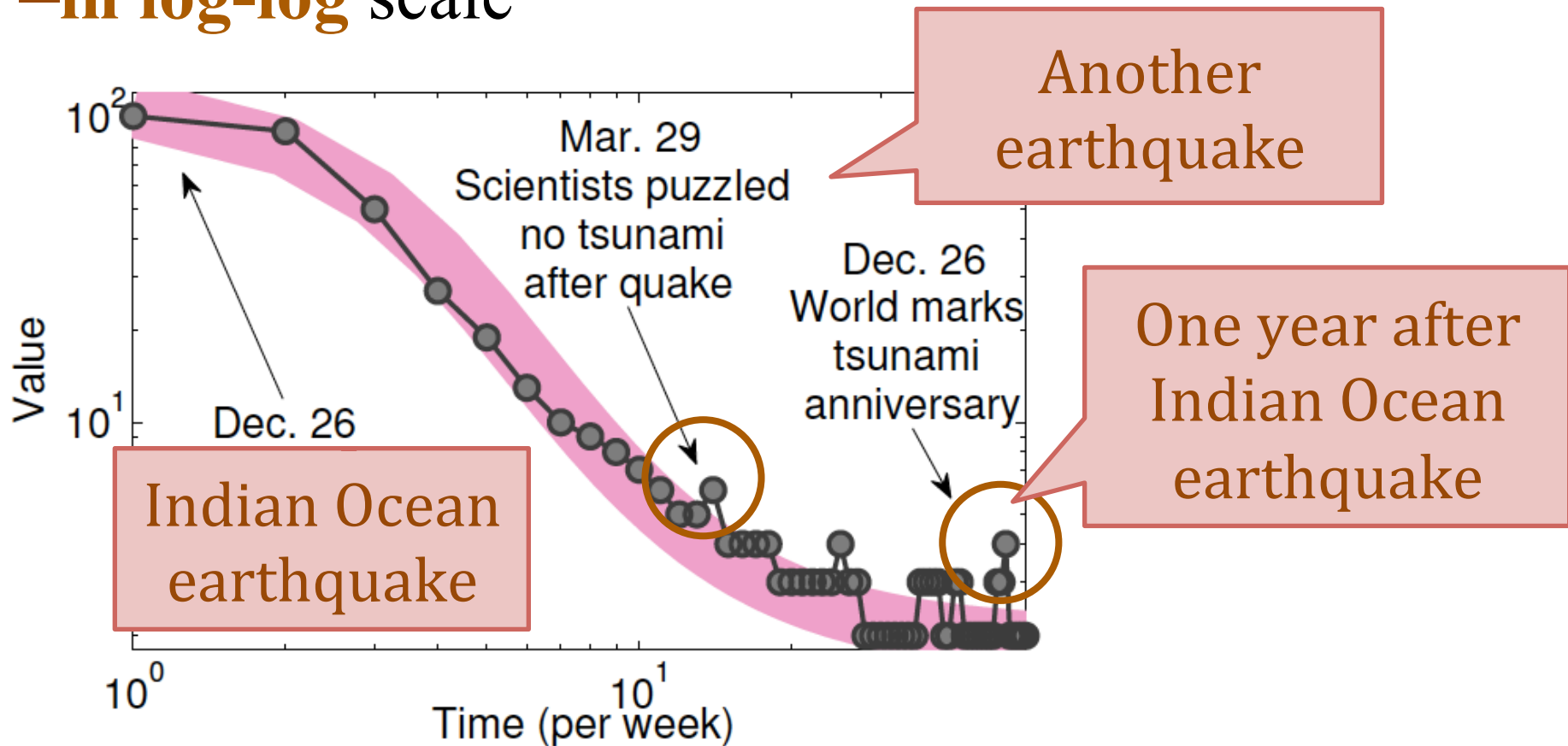


SpikeM can forecast **upcoming spikes!**



A2. Outlier detection

- Fitting result of “tsunami (Google trend)”
- in log-log scale

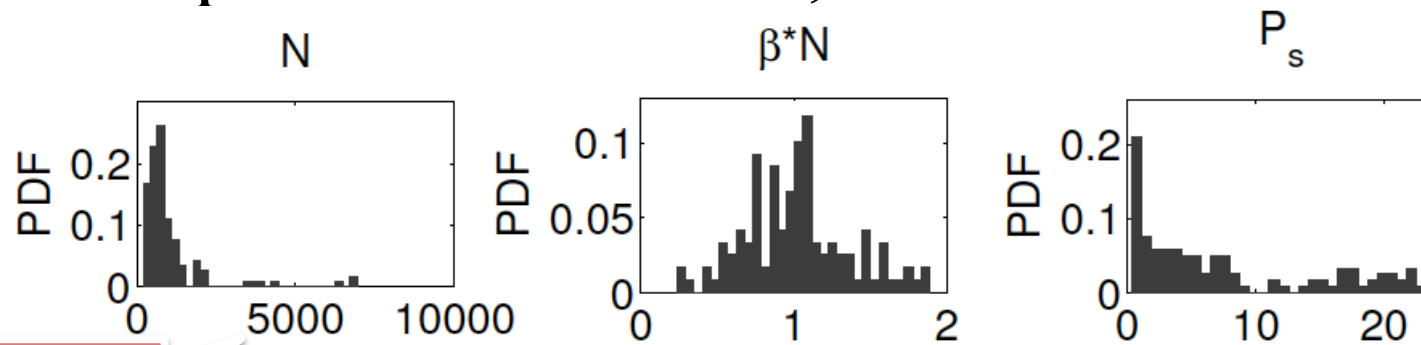




A3. Reverse engineering

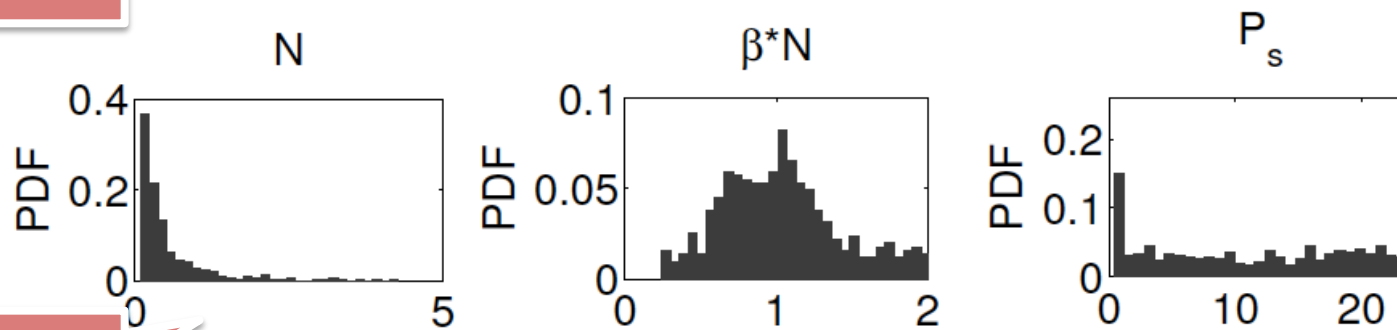
SpikeM provide an intuitive explanation

PDF of parameters over 1,000 memes/hashtags



Meme

(a) *MemeTracker*



Twitter

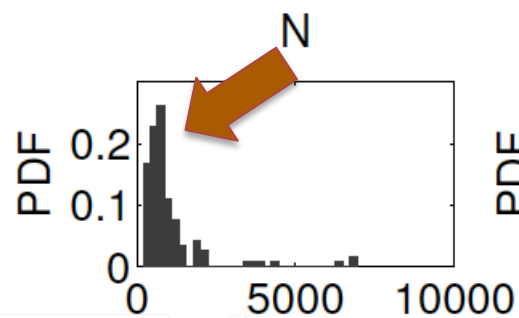
(b) *Twitter*



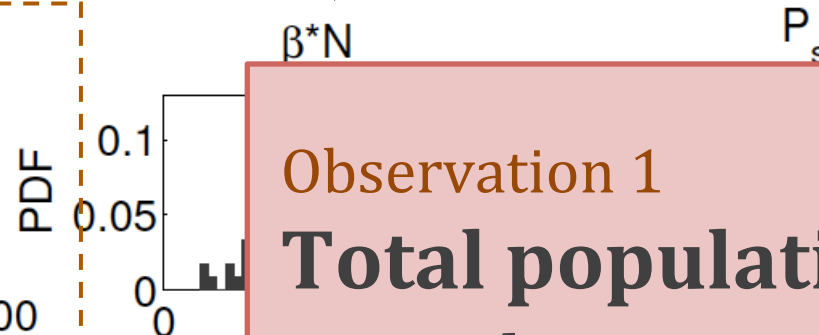
A3. Reverse engineering

SpikeM provide an intuitive explanation

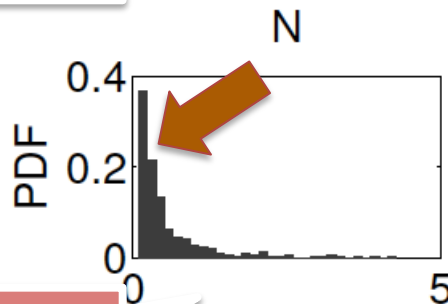
PDF of parameters over 1,000 memes/hashtags



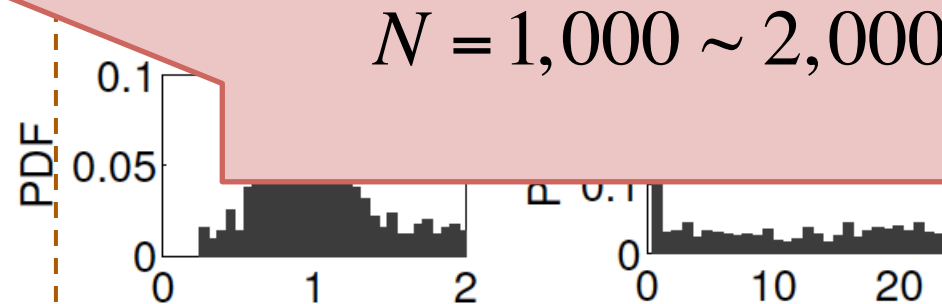
Meme



(a) Memes



Twitter



(b) Twitter

Observation 1

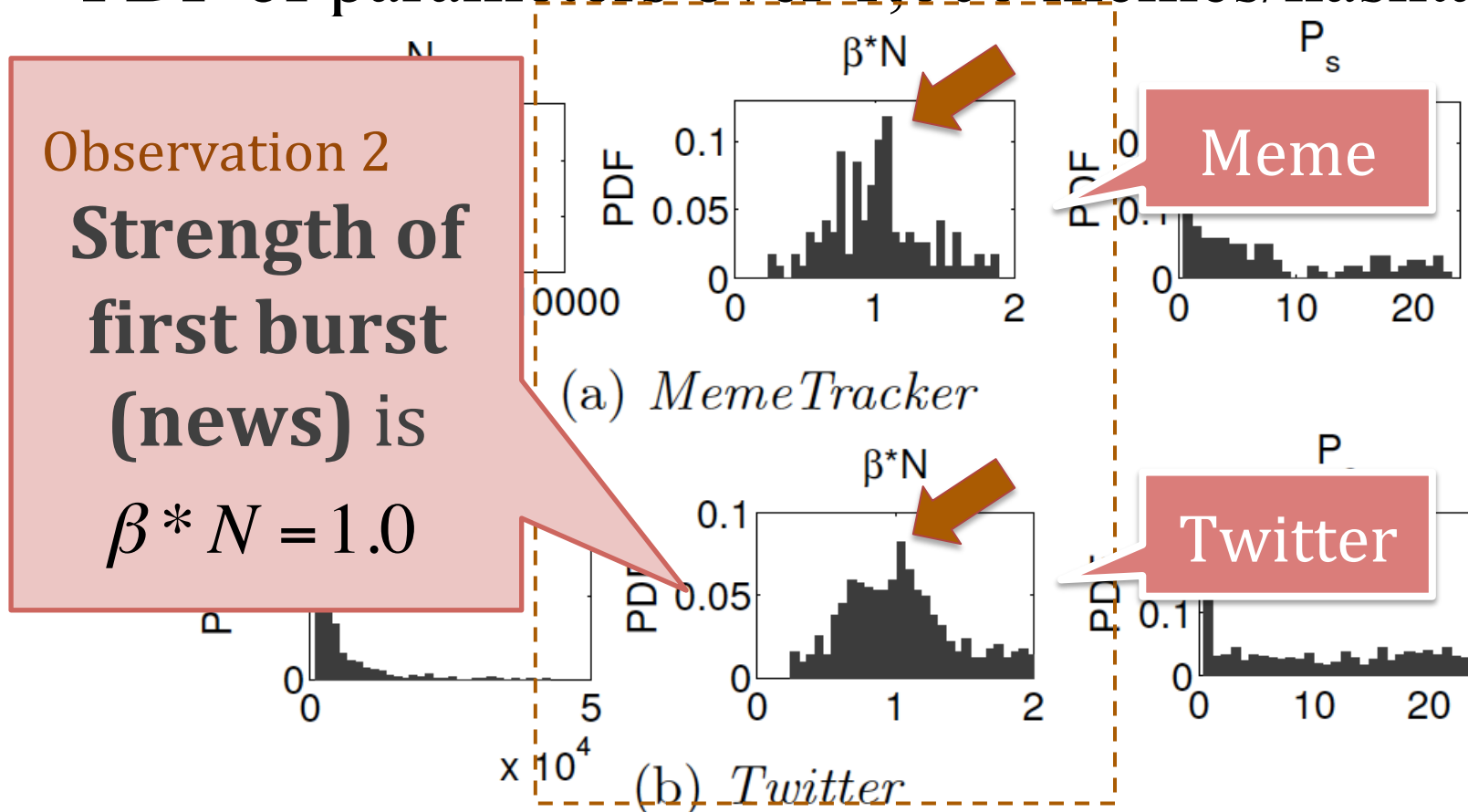
Total population N is almost same

$$N = 1,000 \sim 2,000$$

A3. Reverse engineering

SpikeM provide an intuitive explanation

PDF of parameters over 1,000 memes/hashtags





A3. Reverse engineering

SpikeM provide an intuitive explanation

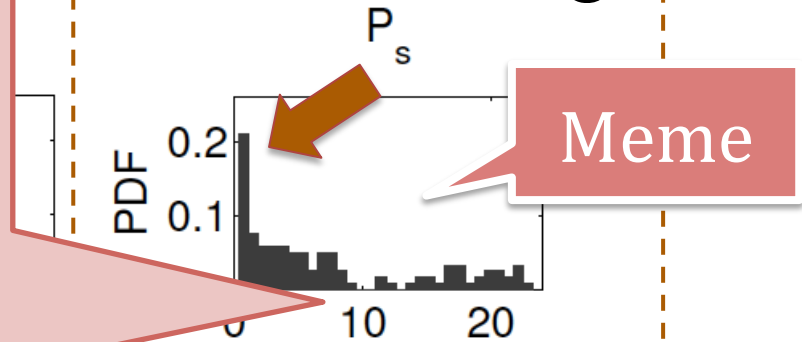
Observation 3

Daily periodicity

with phase shift $P_s = 0$

Every meme has the same periodicity without lag

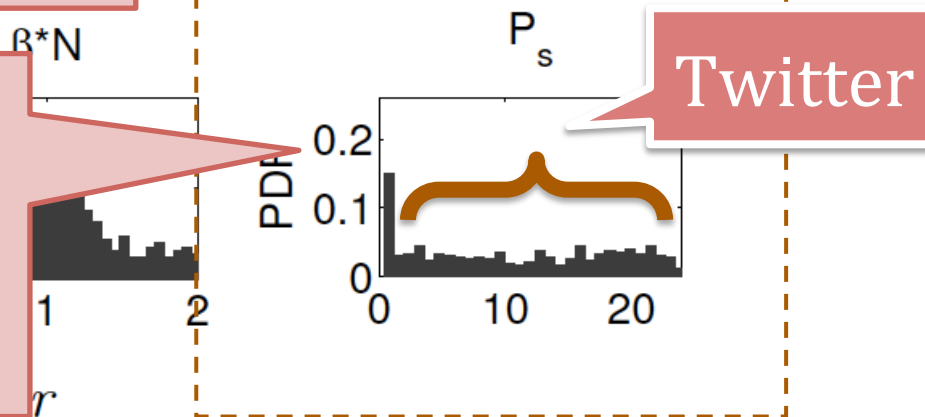
0 memes/hashtags



(Twitter)

Daily periodicity with

more spread in P_s
(i.e., Multiple time zone)





Part 2

Roadmap



Problem

- ✓ Why: “non-linear” modeling

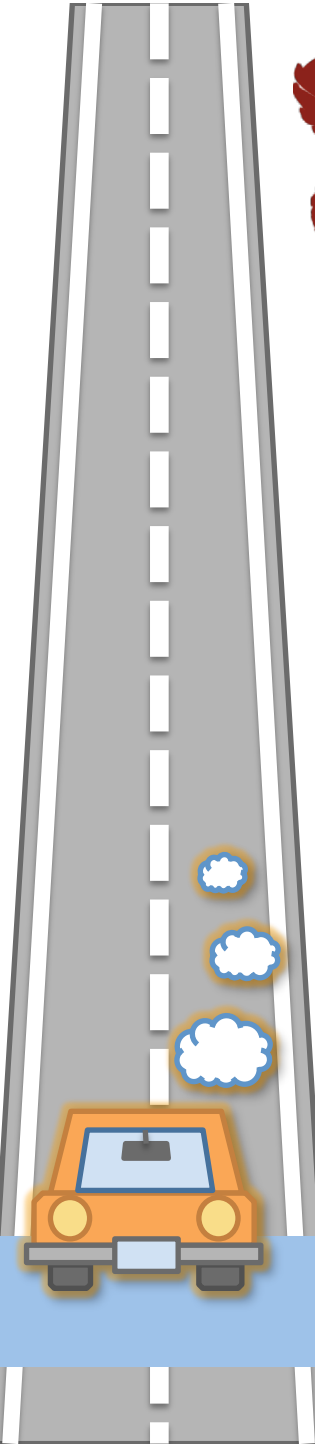
Fundamentals

- ✓ Non-linear (grey-box) models

Applications

- ✓ Epidemics
- ✓ Information diffusion  vs. 

– Online competition



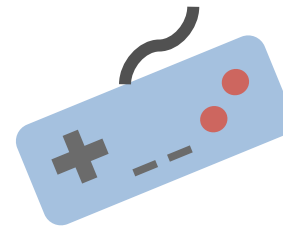


Online competition in social networks





Online competition in social networks



VS.



Q. How can we describe “virtual competition”?

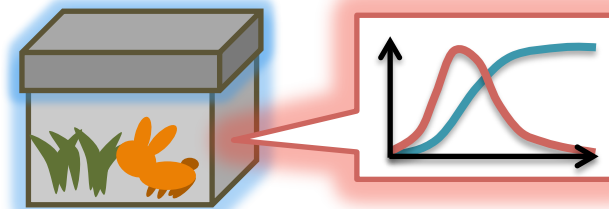


Online competition - roadmap



A. Non-linear (gray-box) modeling!

Solutions



- Winner-Takes-All [Prakash+ WWW'12]
- Co-existence of the two viruses [Beutel+ KDD'12]
- The Web as a Jungle [Matsubara+ WWW'15]

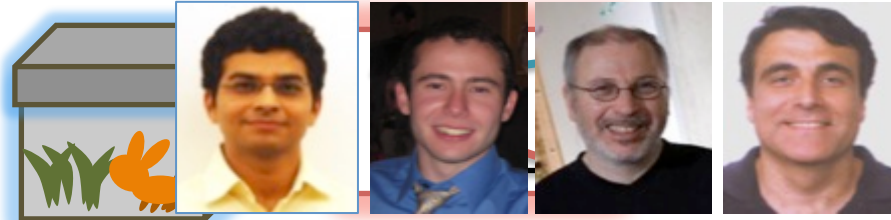


Online competition - roadmap



A. Non-linear (gray-box)
modeling!

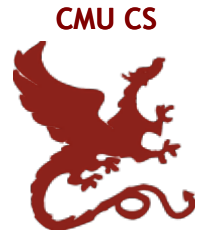
Solutions



- **Winner-Takes-All** [Prakash+ WWW'12]
- **Co-existence of the two viruses** [Beutel+ KDD'12]
- **The Web as a Jungle** [Matsubara+ WWW'15]



Competing contagions



[Prakash+ WWW'12]

Contagions: viruses, online activities



iPhone v Android



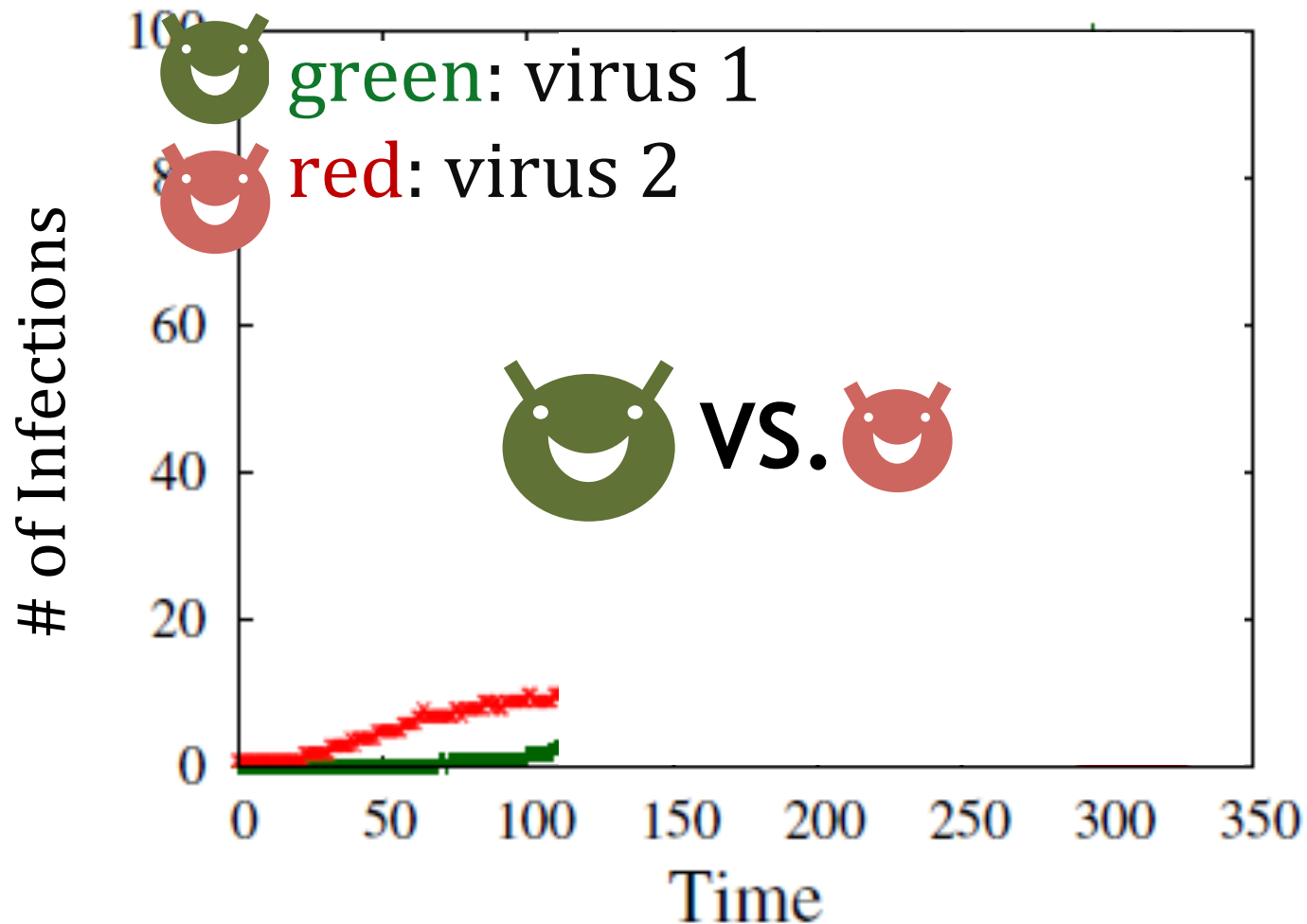
Blu-ray v HD-DVD

Q. What happen when two viruses compete?



Competing contagions

[Prakash+ WWW'12]

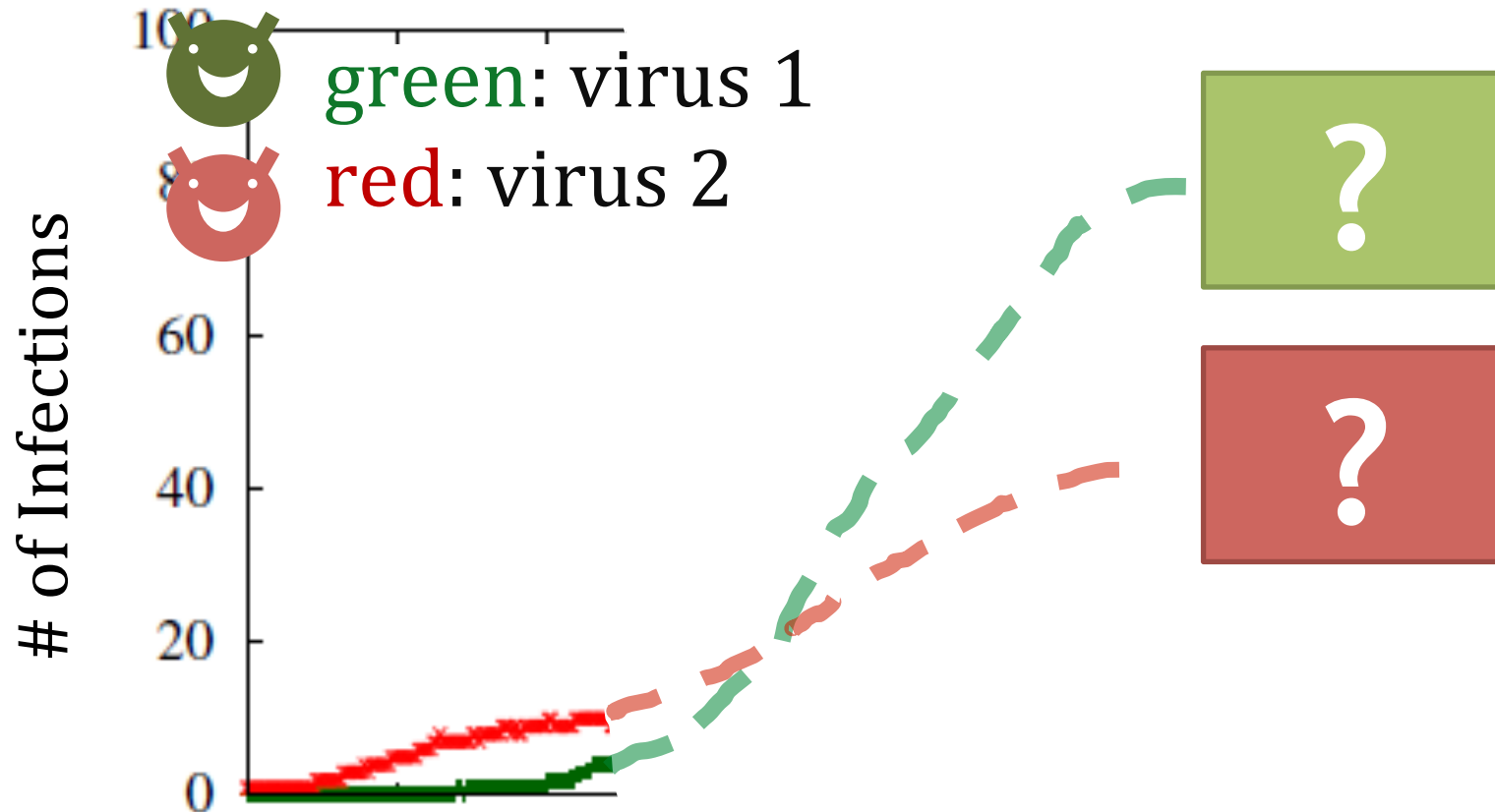


ASSUME: Virus 1 is stronger than Virus 2



Competing contagions

[Prakash+ WWW'12]



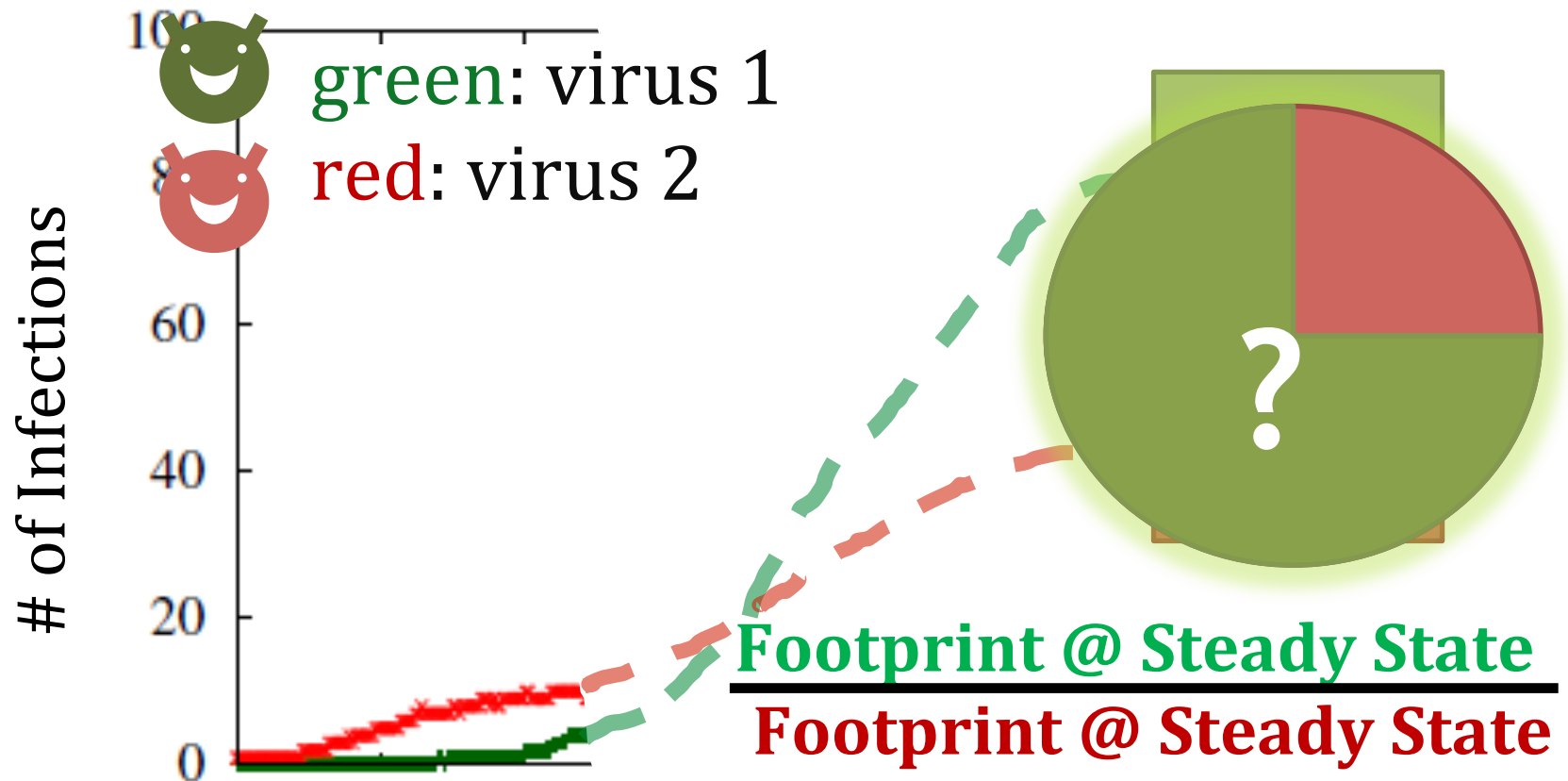
Q: What happens in the end?

ASSUME: virus 1 is stronger than virus 2
<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/> © 2015 Sakurai, Matsubara & Faloutsos



Competing contagions

[Prakash+ WWW'12]



Q: What happens in the end?

ASSUME: virus 1 is stronger than virus 2
<http://www.cs.kumamoto-u.ac.jp/~yasuko/TALKS/15-SIGMOD-tut/>
 © 2015 Sakurai, Matsubara & Faloutsos

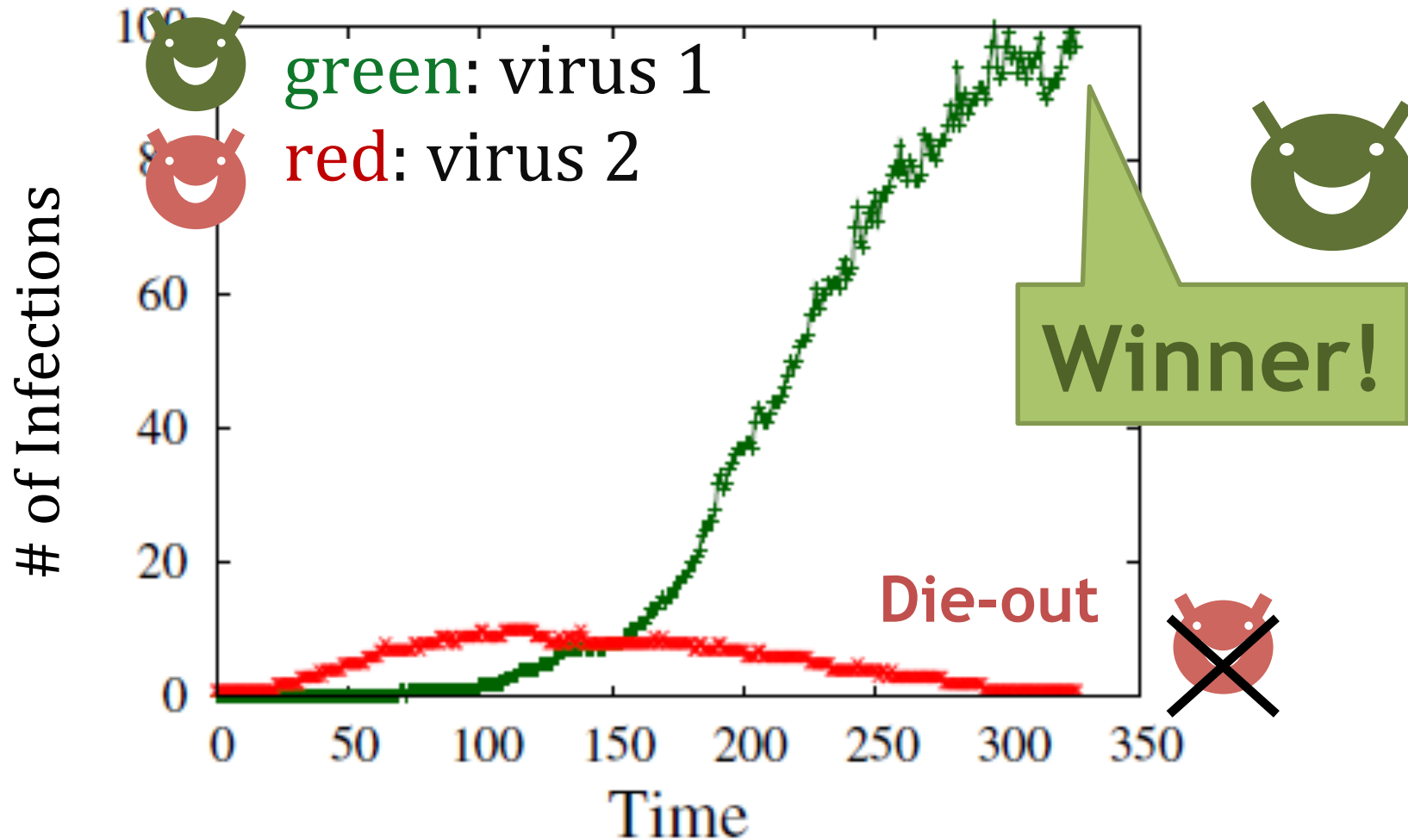


Answer:

Winner-Takes-All!



[Prakash+ WWW'12]



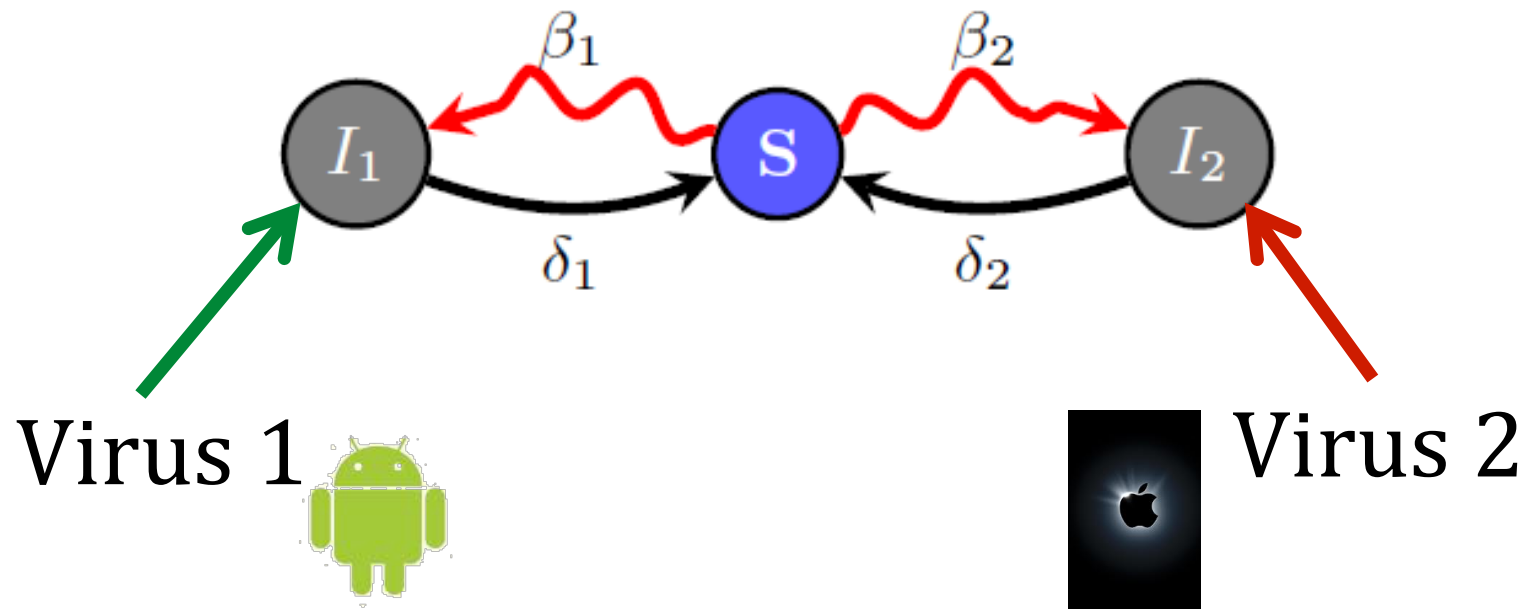
ASSUME: Virus 1 is stronger than Virus 2



A simple model

[Prakash+ WWW'12]

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



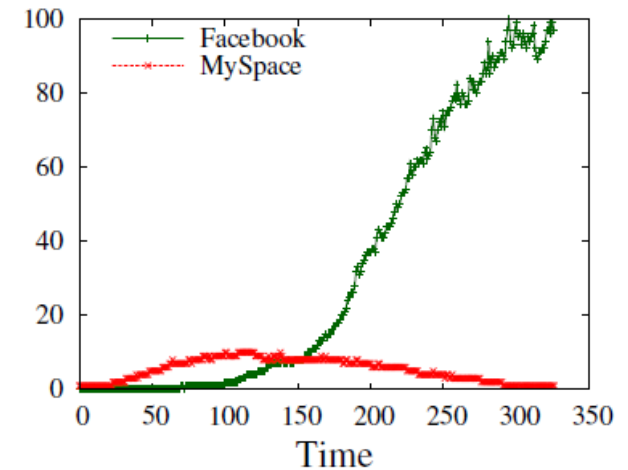


Result: Winner-Takes-All



[Prakash+ WWW'12]

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



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!

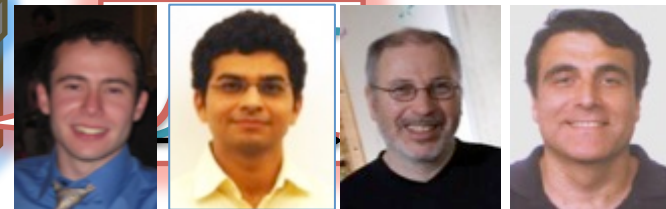
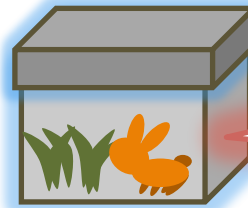


Online competition in social networks



A. Non-linear (gray-box)
modeling!

Solutions



- Winner-Takes-All [Prakash+ WWW'12]
- **Co-existence of the two viruses** [Beutel+ KDD'12]
- The Web as a Jungle [Matsubara+ WWW'15]

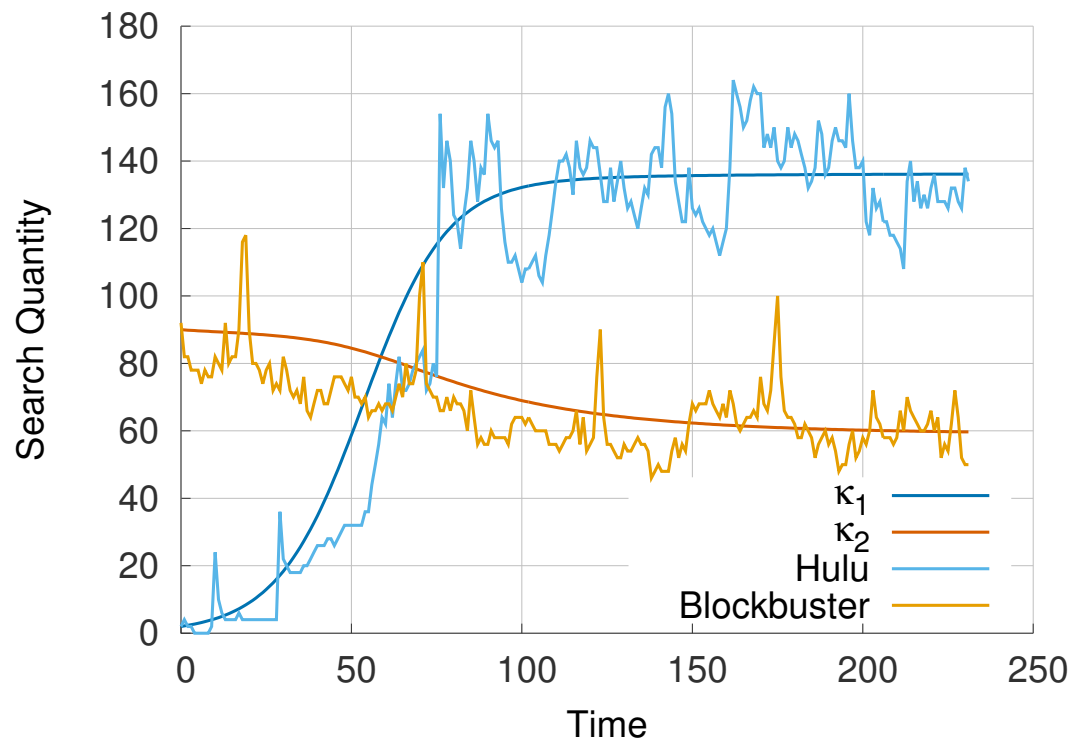


Interacting Viruses: Can Both Survive?



Real example of “co-existence”

[Google Search Trends data]



Hulu v Blockbuster

hulu



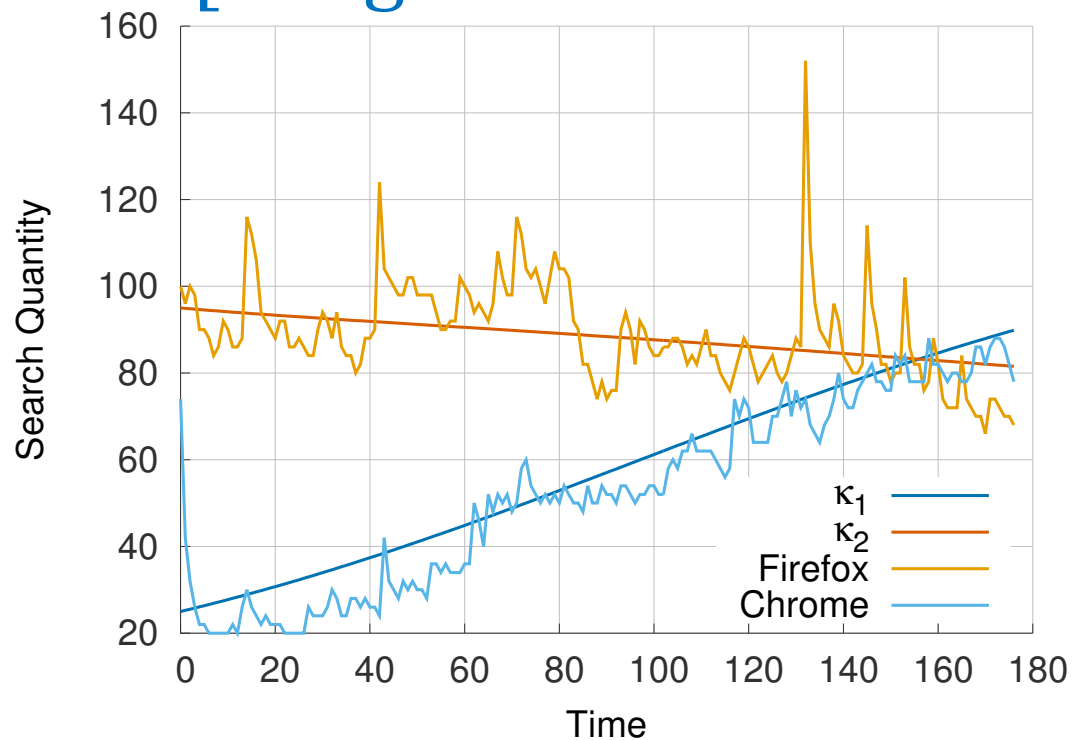


Interacting Viruses: Can Both Survive?



Real example of “co-existence”

[Google Search Trends data]



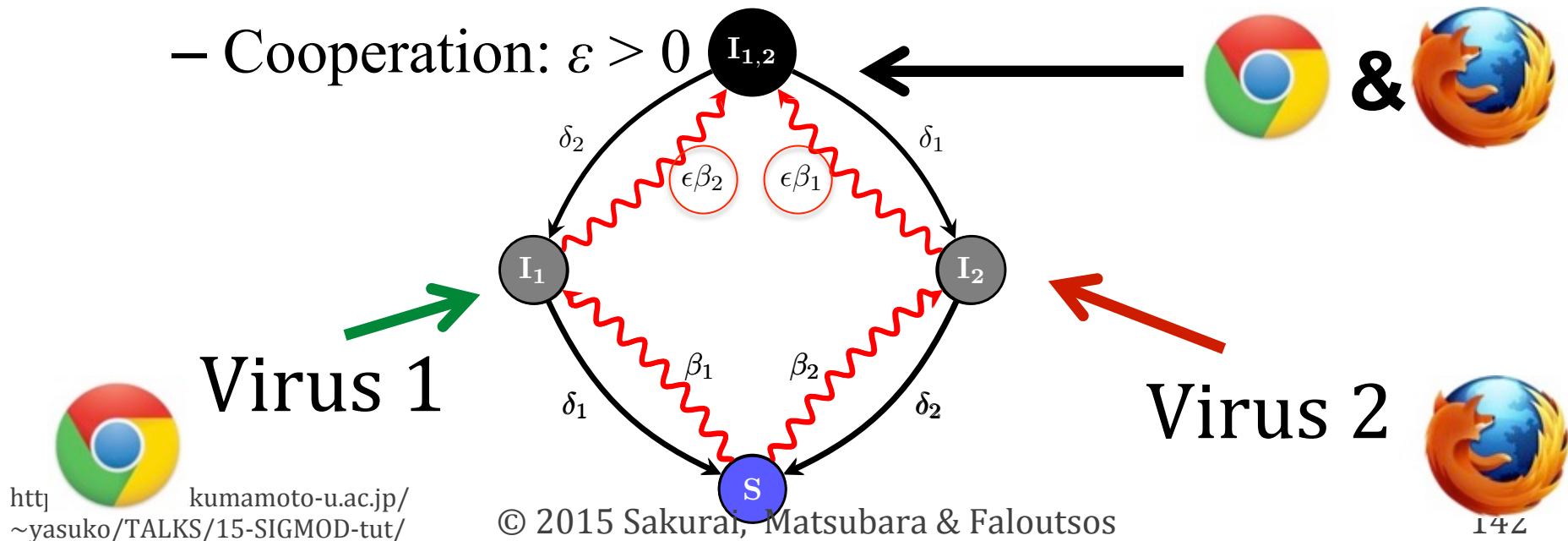
Chrome v Firefox





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

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





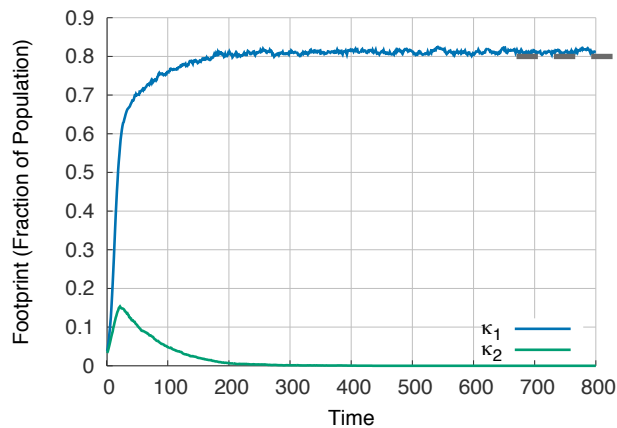
Question:



What happens in the end?

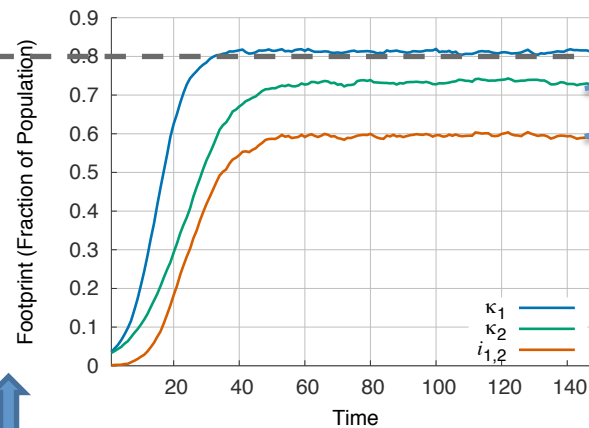
$\epsilon = 0$

Winner takes all



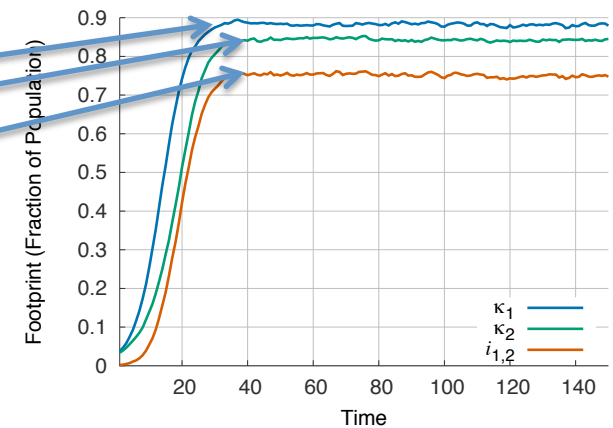
$\epsilon = 1$

Co-exist independently



$\epsilon = 2$

Viruses cooperate



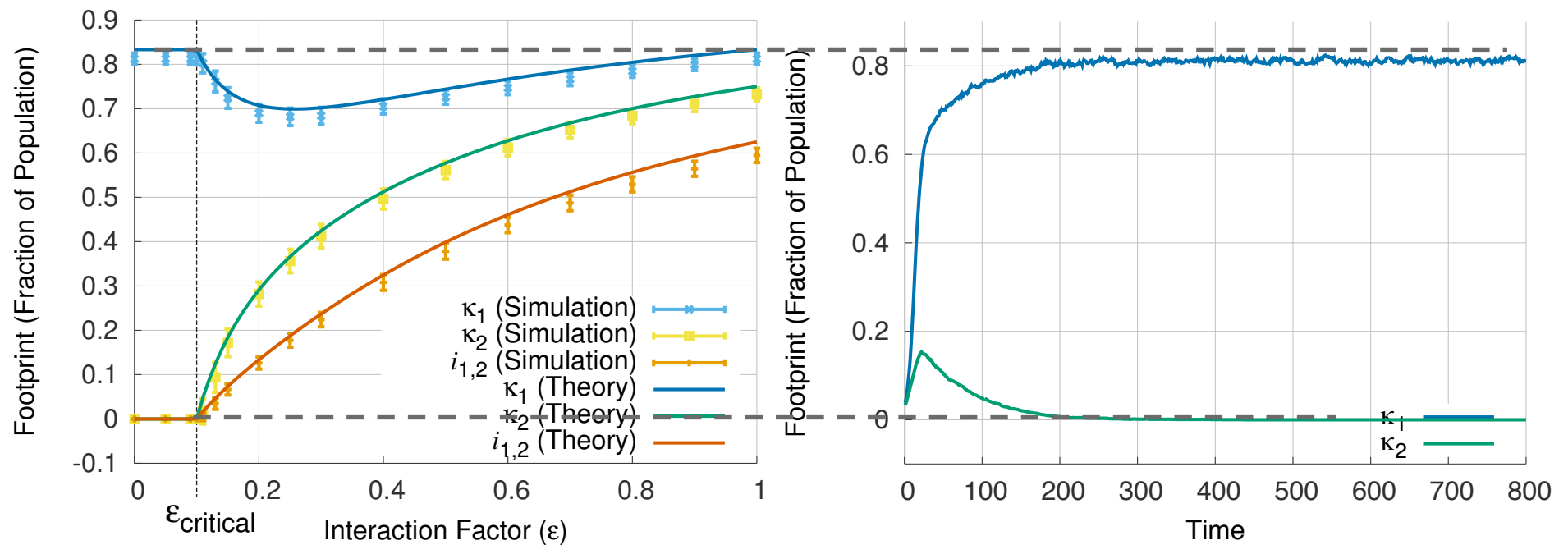
What about for $0 < \epsilon < 1$?
Is there a point at which both viruses
can *co-exist*?

ASSUME: Virus 1 is stronger than virus 2



Answer: Yes!

There is a phase transition

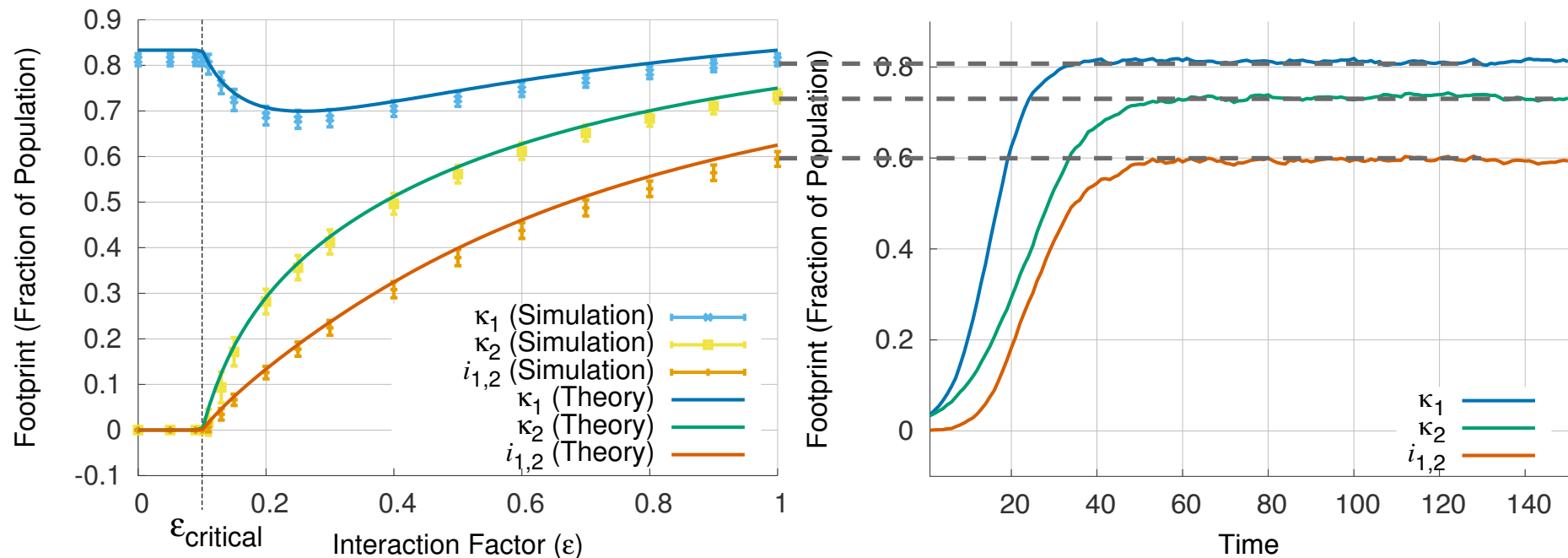


ASSUME: Virus 1 is stronger than Virus 2



Answer: Yes!

There is a phase transition

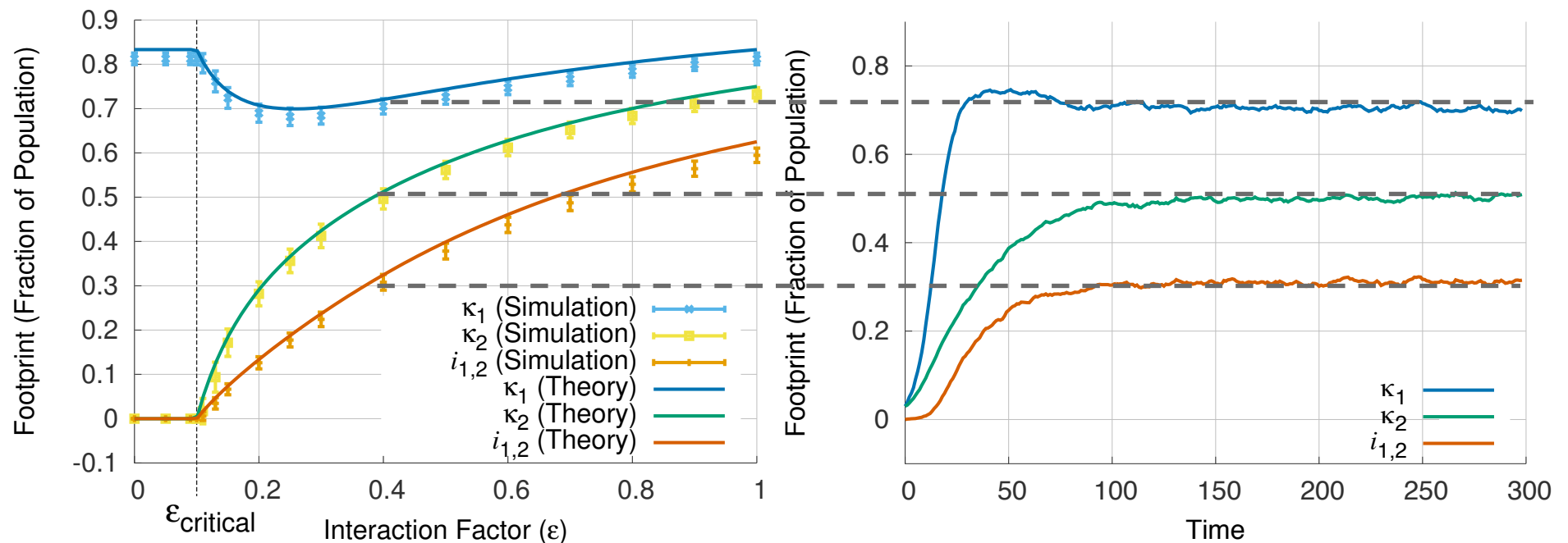


ASSUME: Virus 1 is stronger than Virus 2



Answer: Yes!

There is a phase transition



ASSUME: Virus 1 is stronger than Virus 2



Result:

Viruses can Co-exist



Given this 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.

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. $\text{Strength}(\text{Virus}) \sigma = N \beta / \delta$

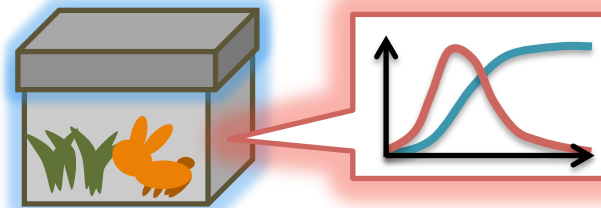


Online competition in social networks



A. Non-linear (gray-box)
modeling!

Solutions



- Winner-Takes-All [Prakash+ WWW'12]
- Co-existence of the two viruses [Beutel+ KDD'12]
- **The Web as a Jungle** [Matsubara+ WWW'15]



[Matsubara+ WWW'15]

The Web as a Jungle: Non-Linear Dynamical Systems for Co-evolving Online Activities

Yasuko Matsubara (Kumamoto University)

Yasushi Sakurai (Kumamoto University)

Christos Faloutsos (CMU)



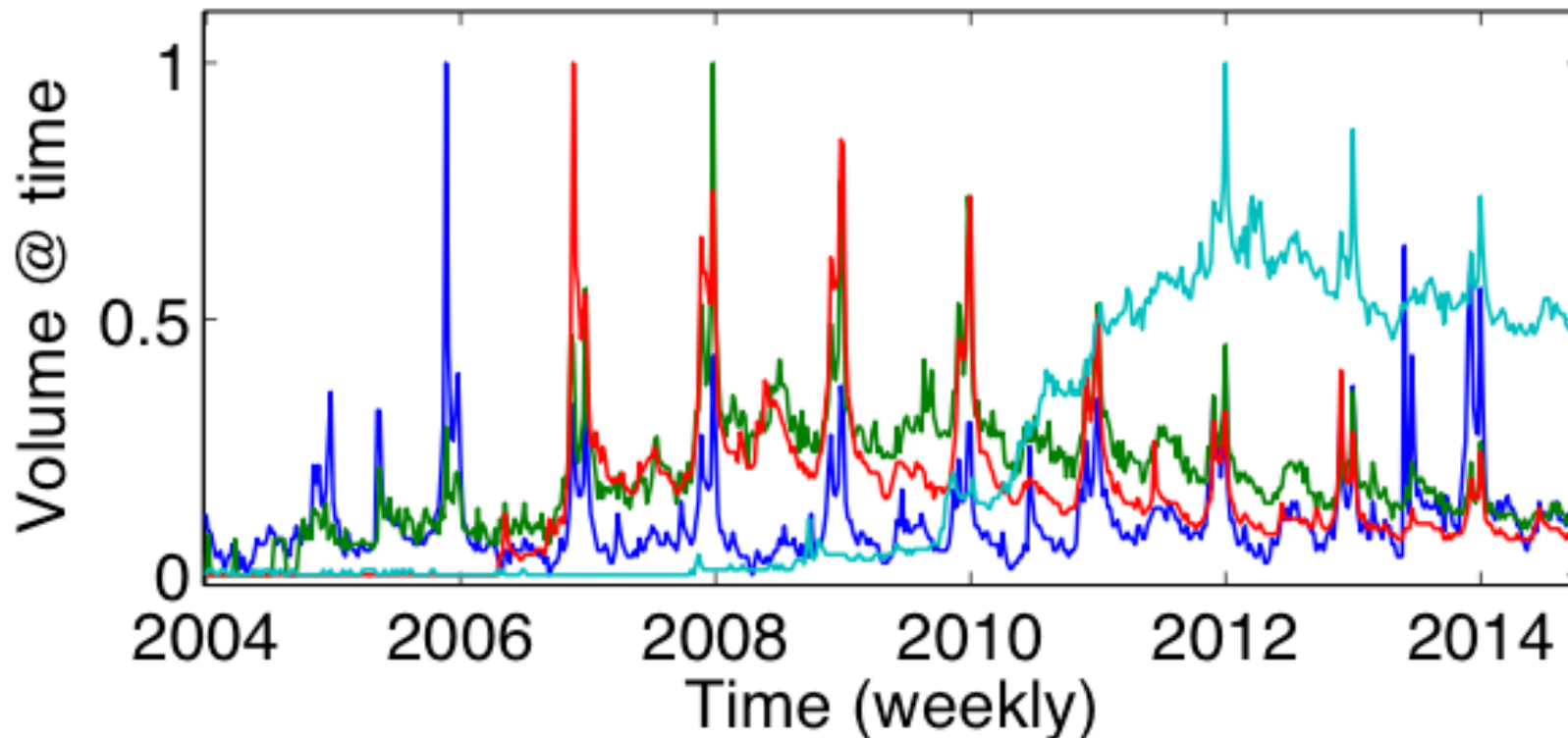


Given: online user activities



e.g., Google search volumes for

Xbox, **PlayStation**, **Wii**, **Android**



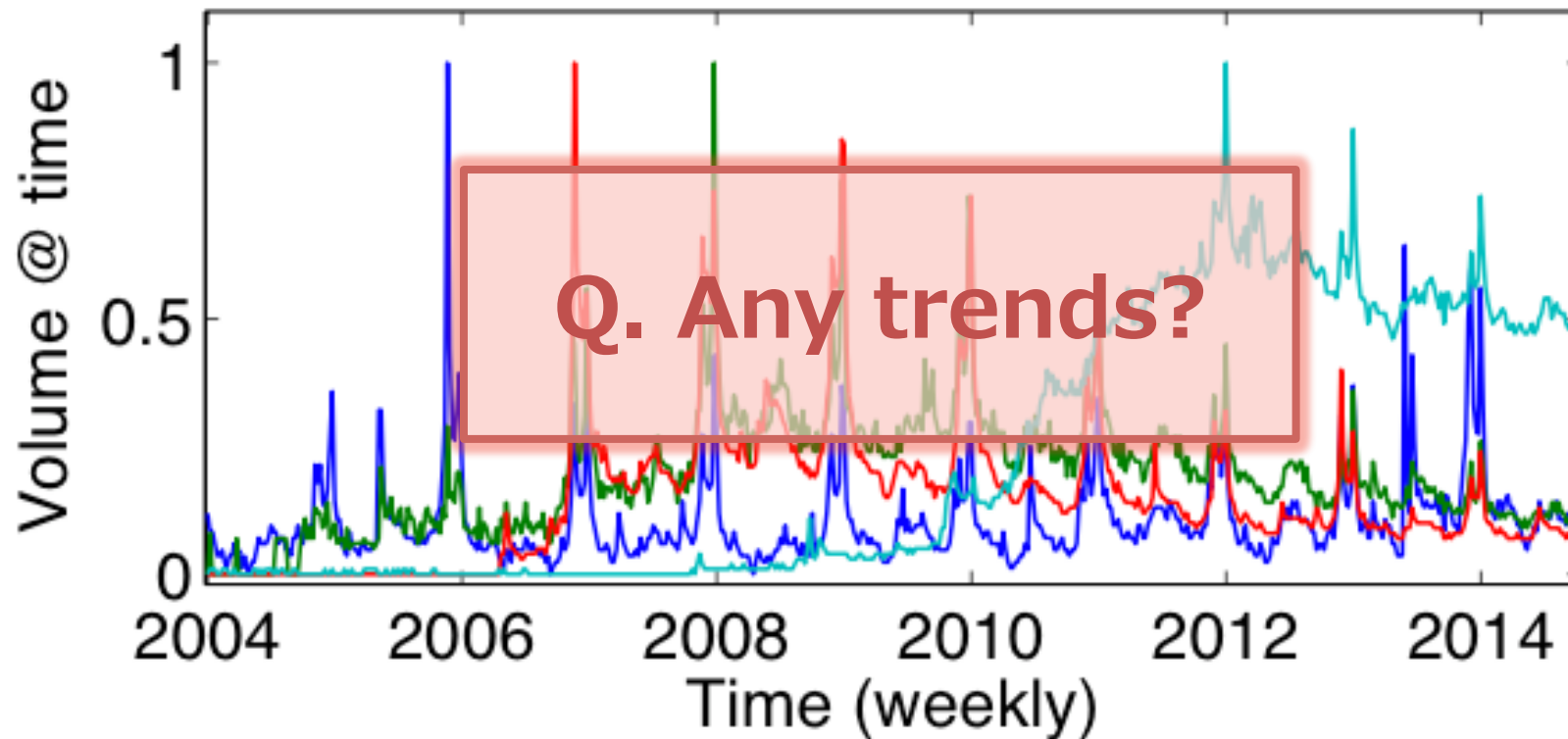


Given: online user activities



e.g., Google search volumes for

Xbox, **PlayStation**, **Wii**, **Android**



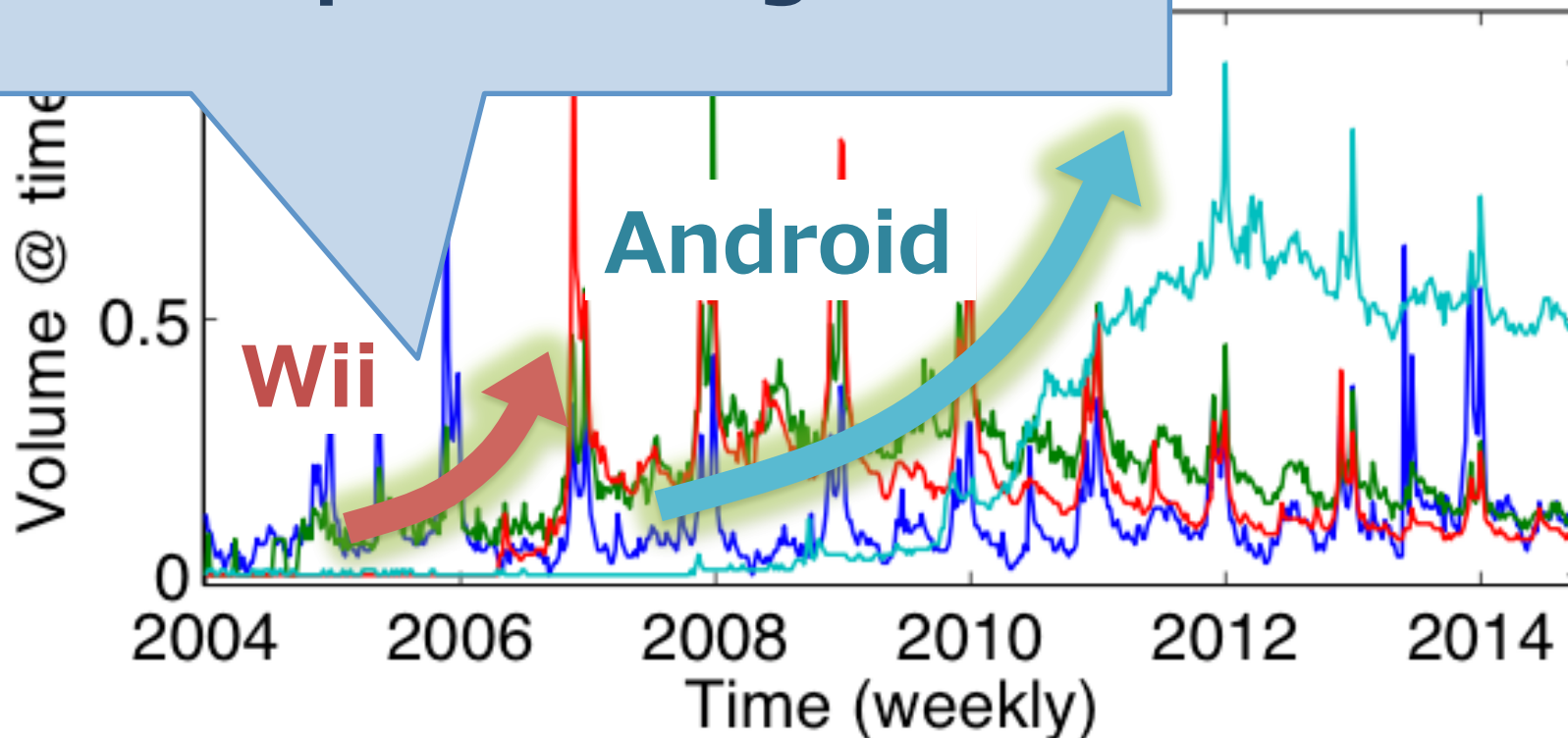


Given: online user activities



e.g., Google search volumes for

1. Exponential growth, Android





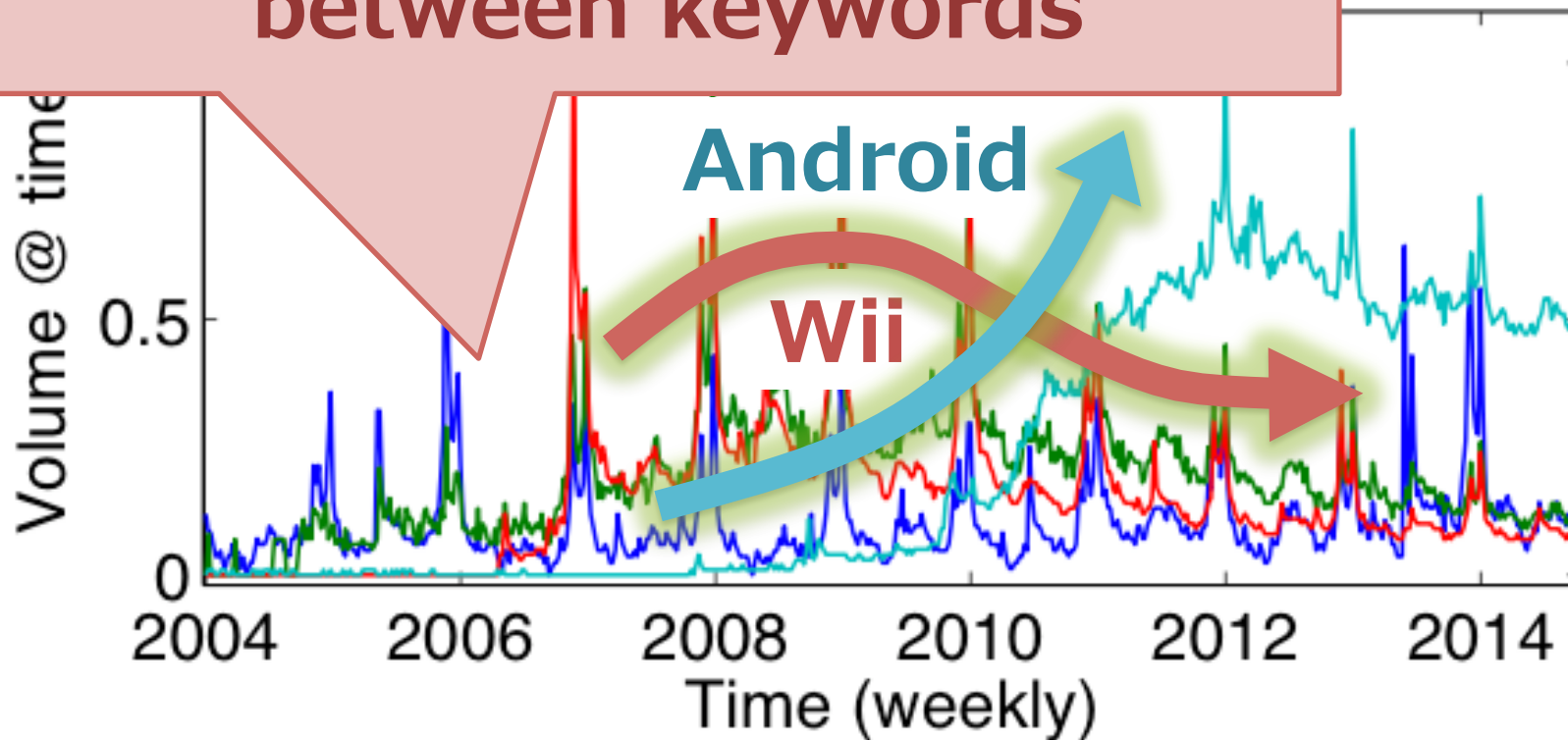
Given: online user activities



e.g., Google search volumes for

2. (Hidden) interaction between keywords

droid





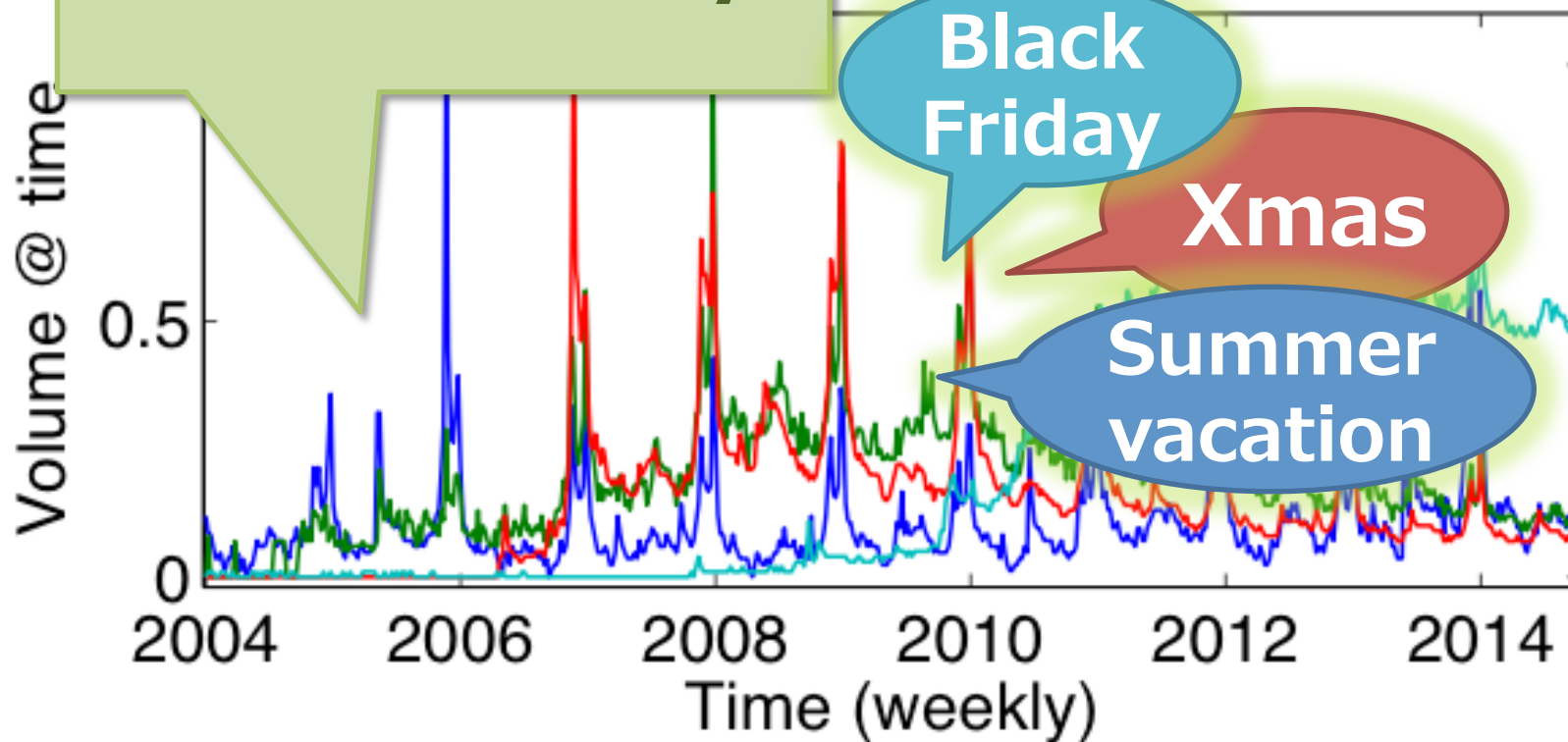
Given: online user activities



e.g., Google search volumes for

3. Seasonality

iPhone, Wii, Android



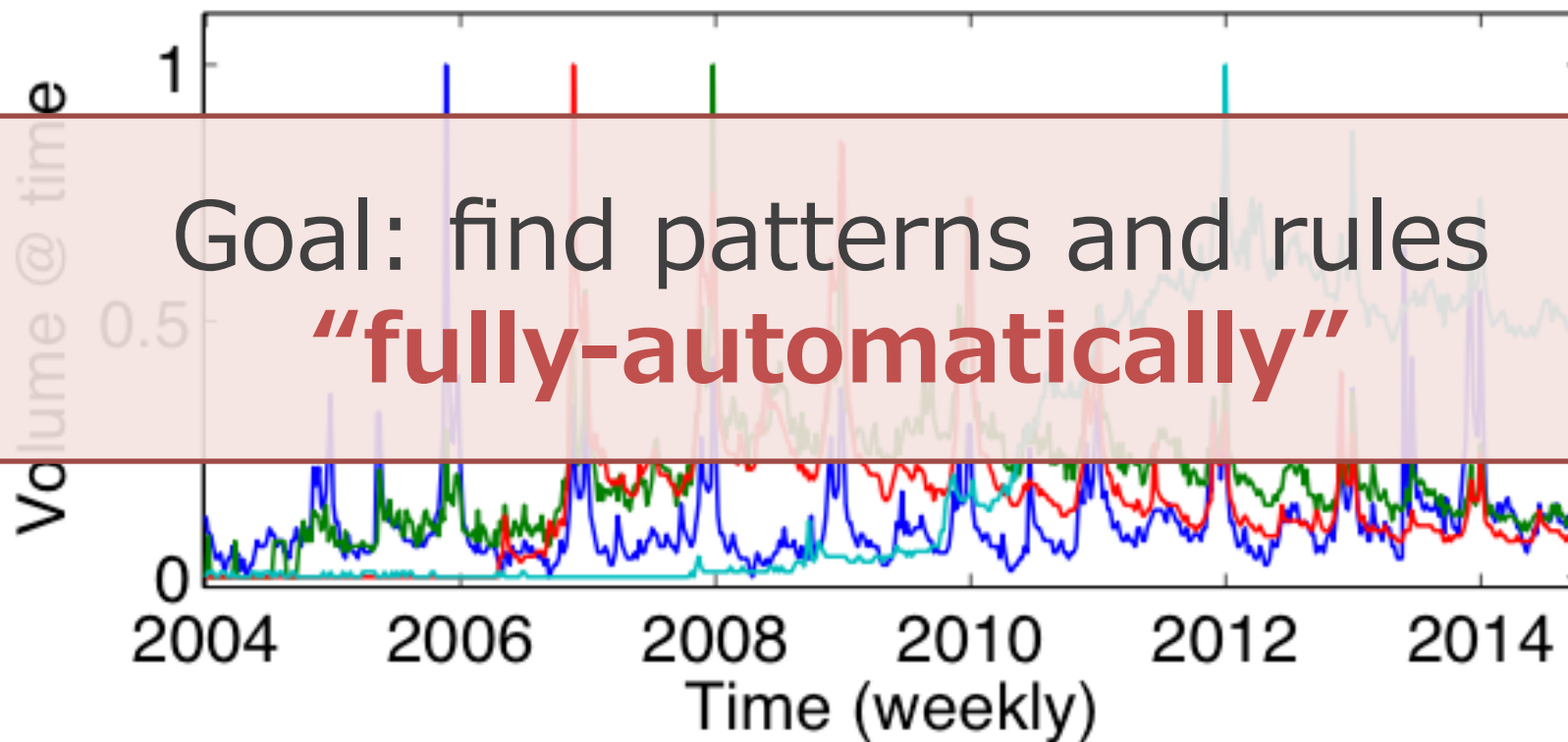


Given: online user activities



e.g., Google search volumes for

Xbox, **PlayStation**, **Wii**, **Android**

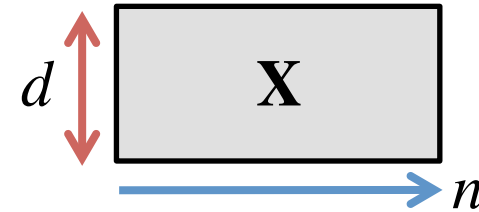




Problem definition

Given: Co-evolving online activities

X (activity x time)

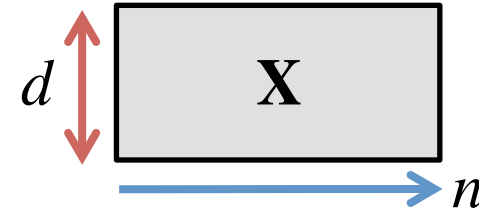




Problem definition

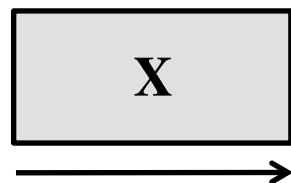
Given: Co-evolving online activities

X (activity \times time)

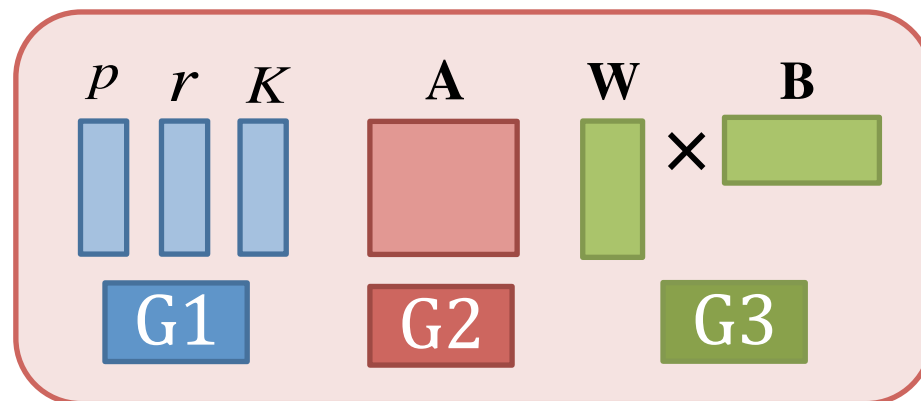


Find: Compact description of X

EcoWeb

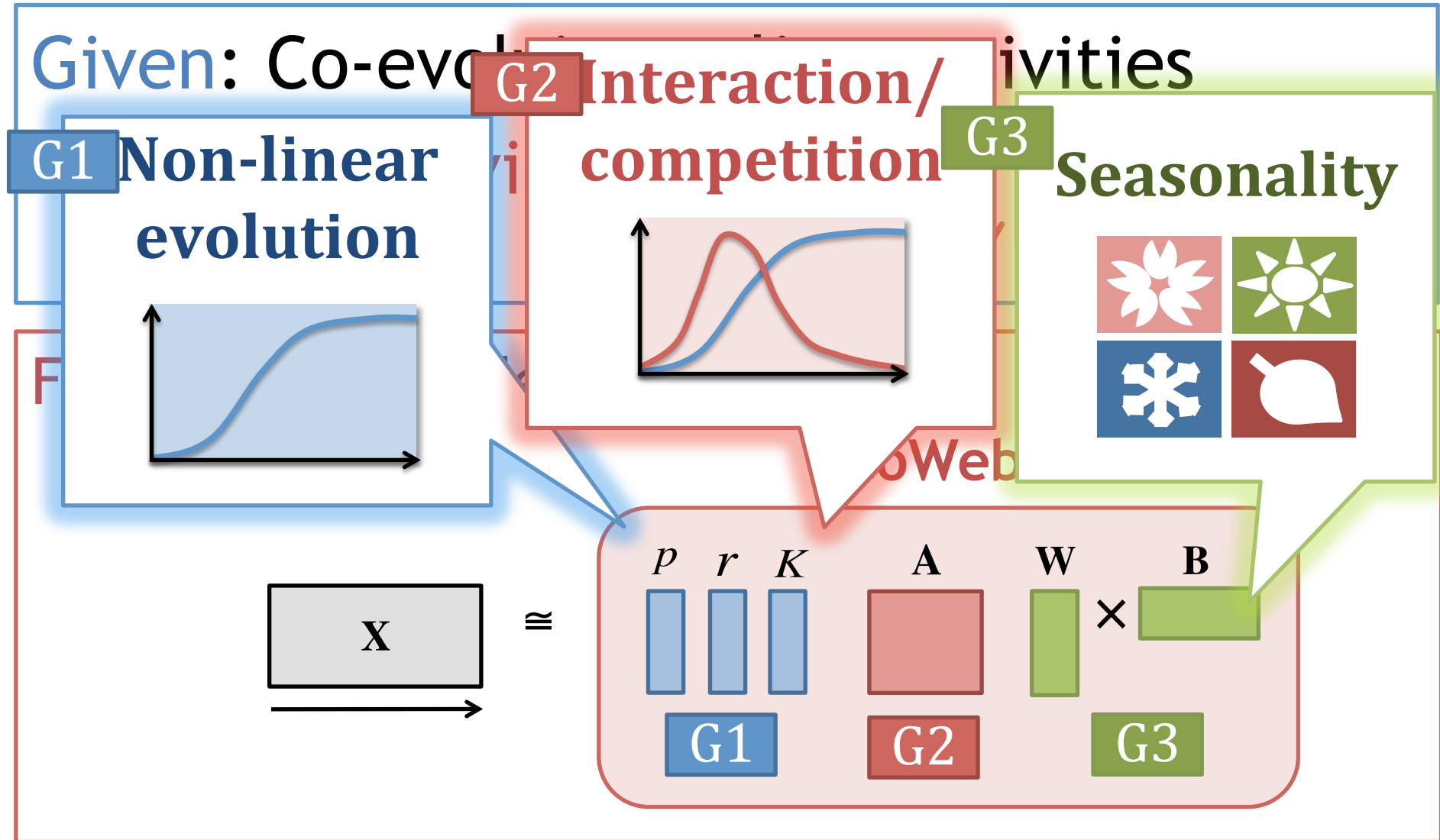


\mathbb{R}



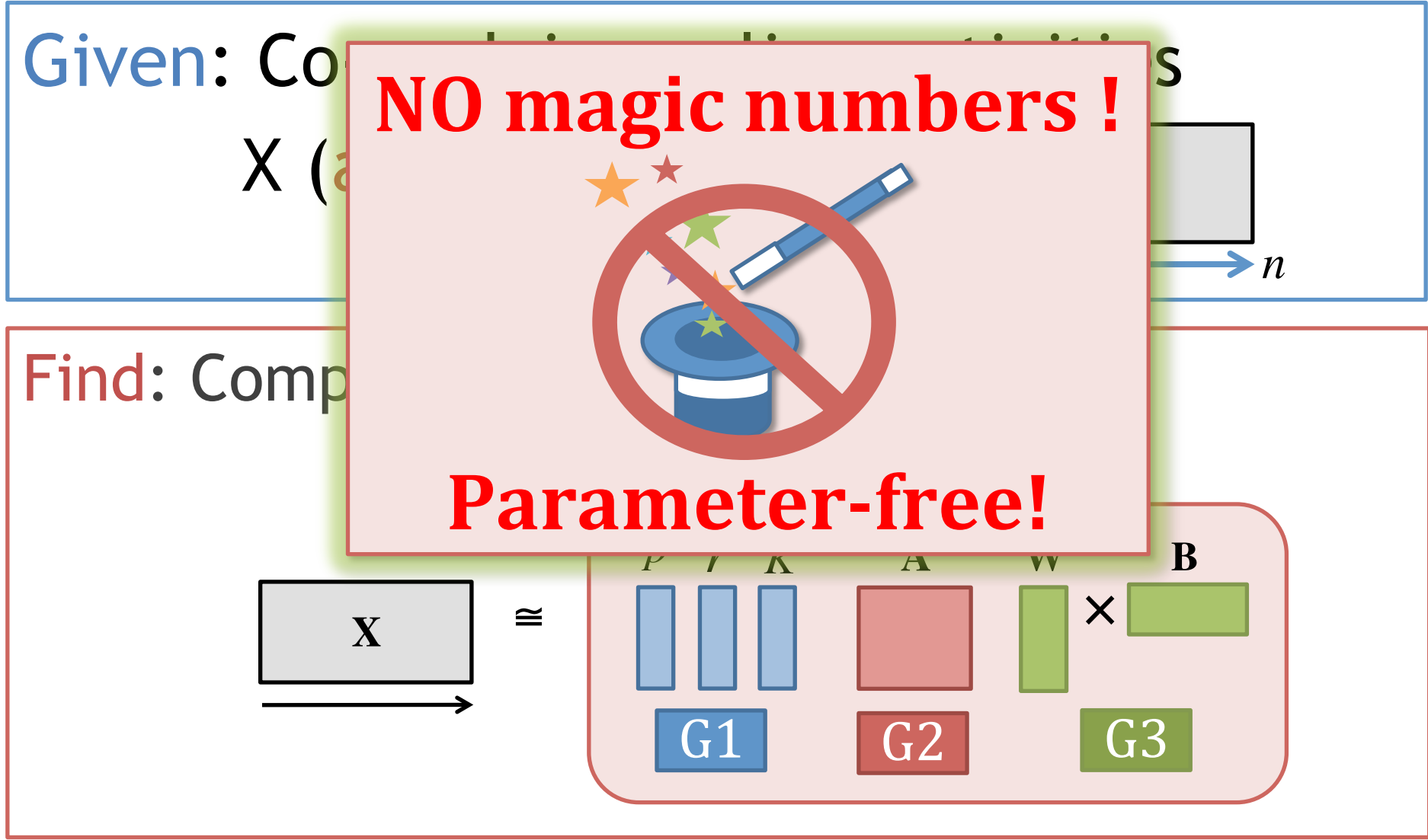


Problem definition





Problem definition



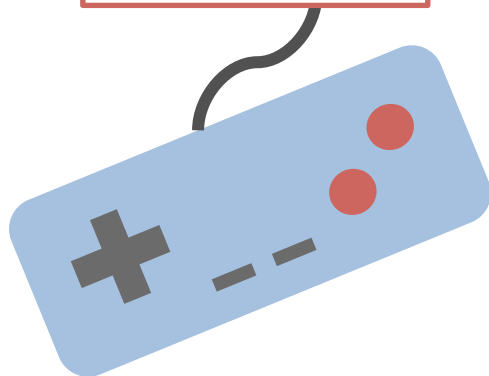


Modeling power of EcoWeb



Questions

Q1



Q2



Q3





Modeling power of EcoWeb



Q1 (games)

Who is the competitor?



VS.

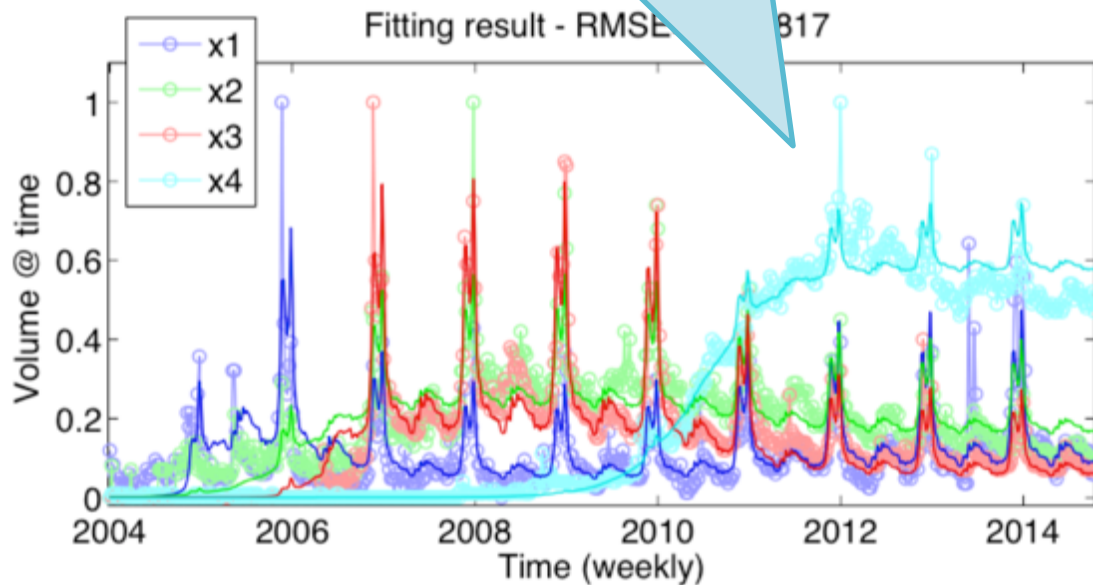




Modeling power of EcoWeb



A. Android!



EcoWeb-Fit

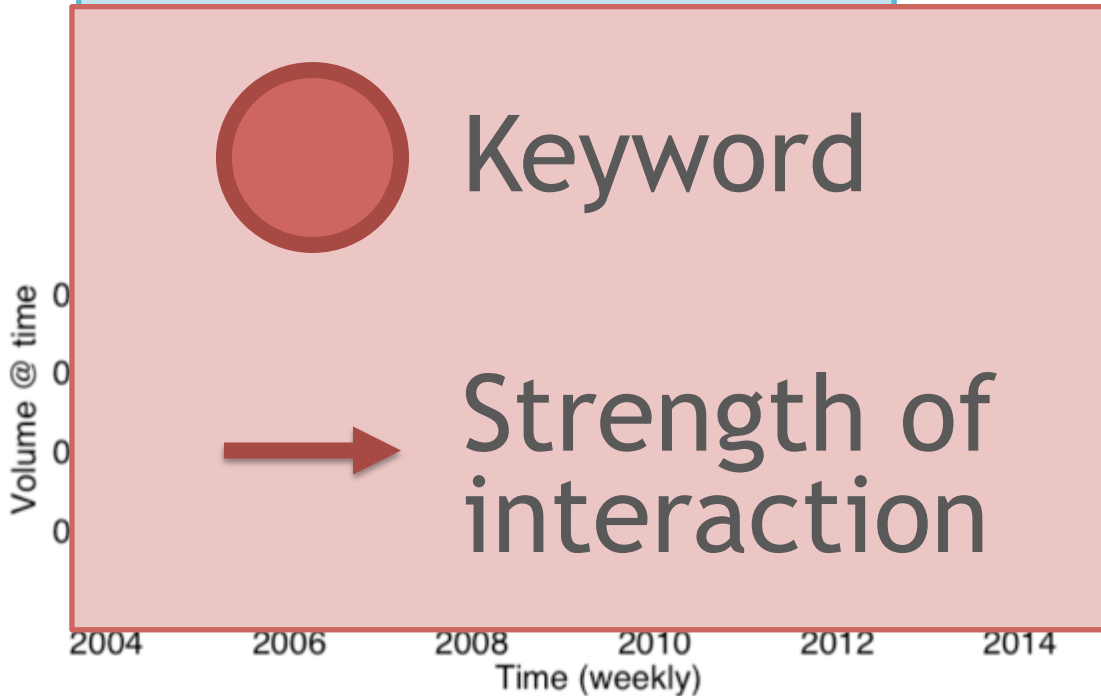
Interaction network (latent)



Modeling power of EcoWeb



A. Android!



Interaction network (latent)

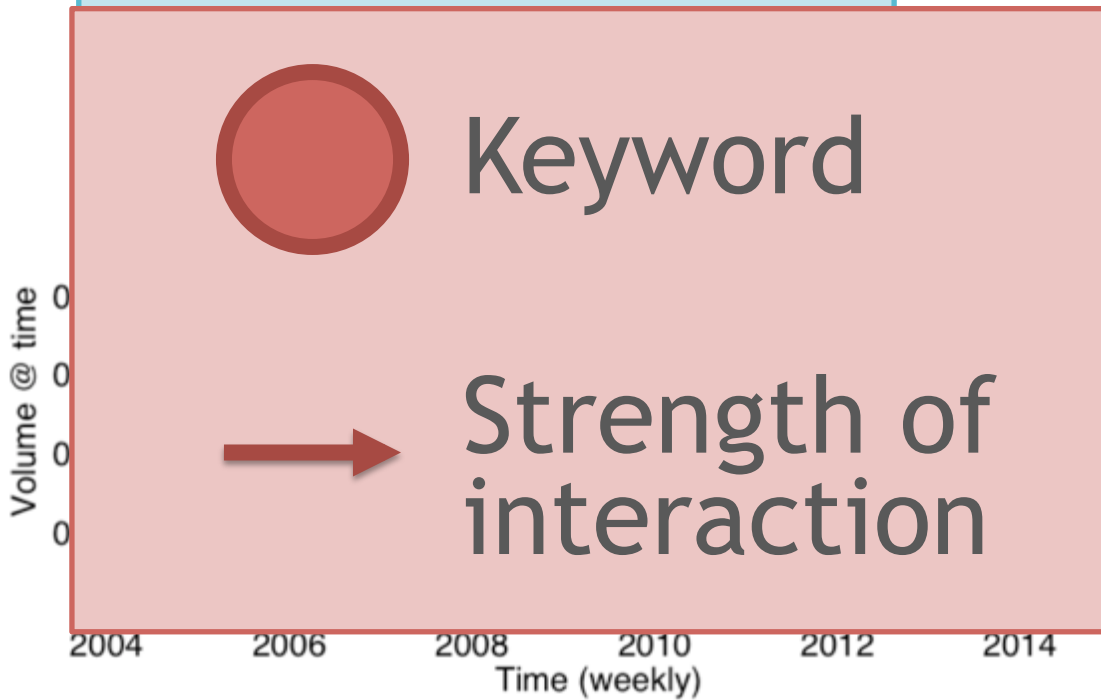
EcoWeb-Fit



Modeling power of EcoWeb



A. Android!

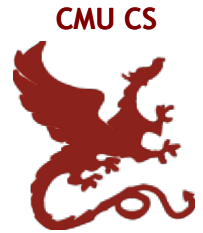


EcoWeb-Fit

Interaction network (latent)

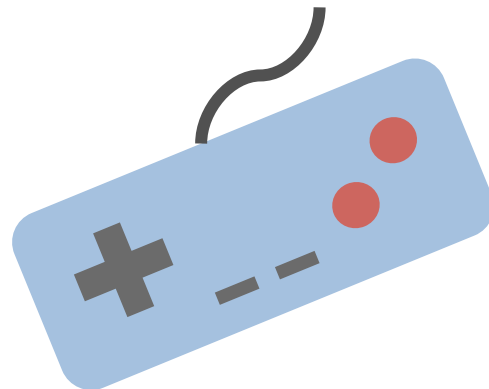


Modeling power of EcoWeb



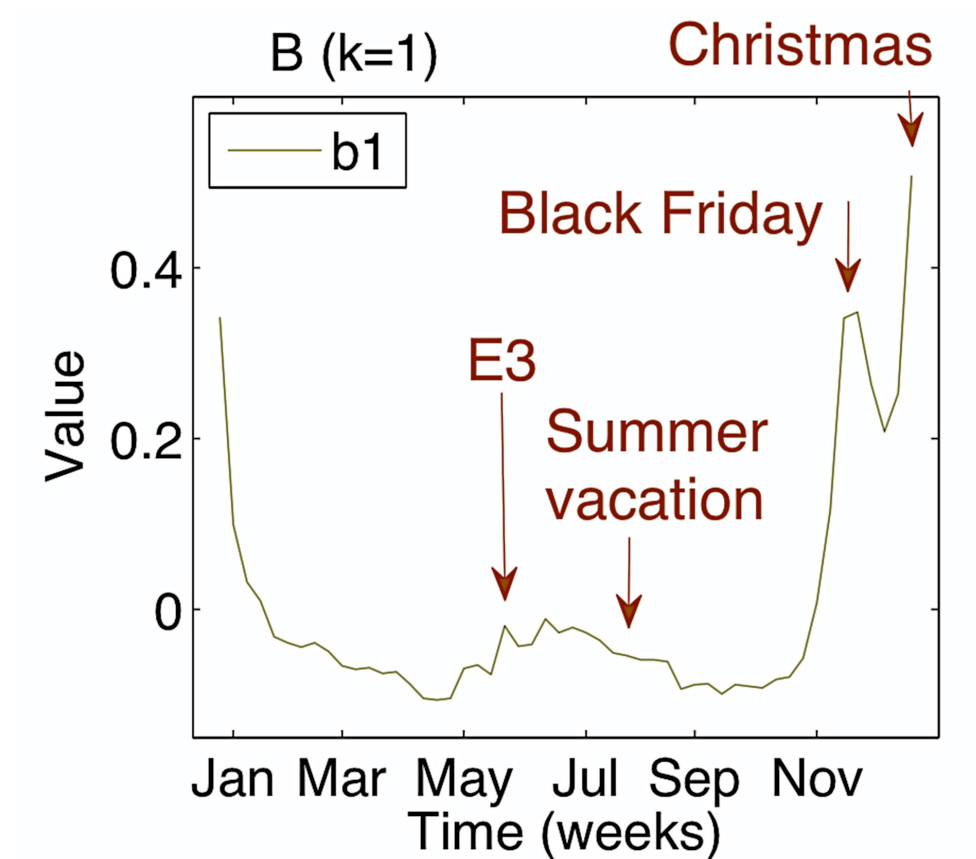
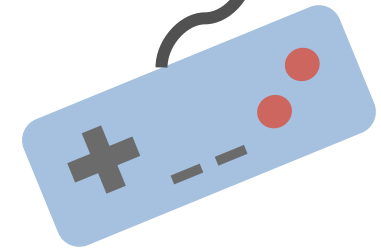
Q1 (games)

Any seasonal events?





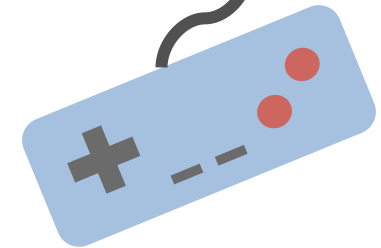
Modeling power of EcoWeb



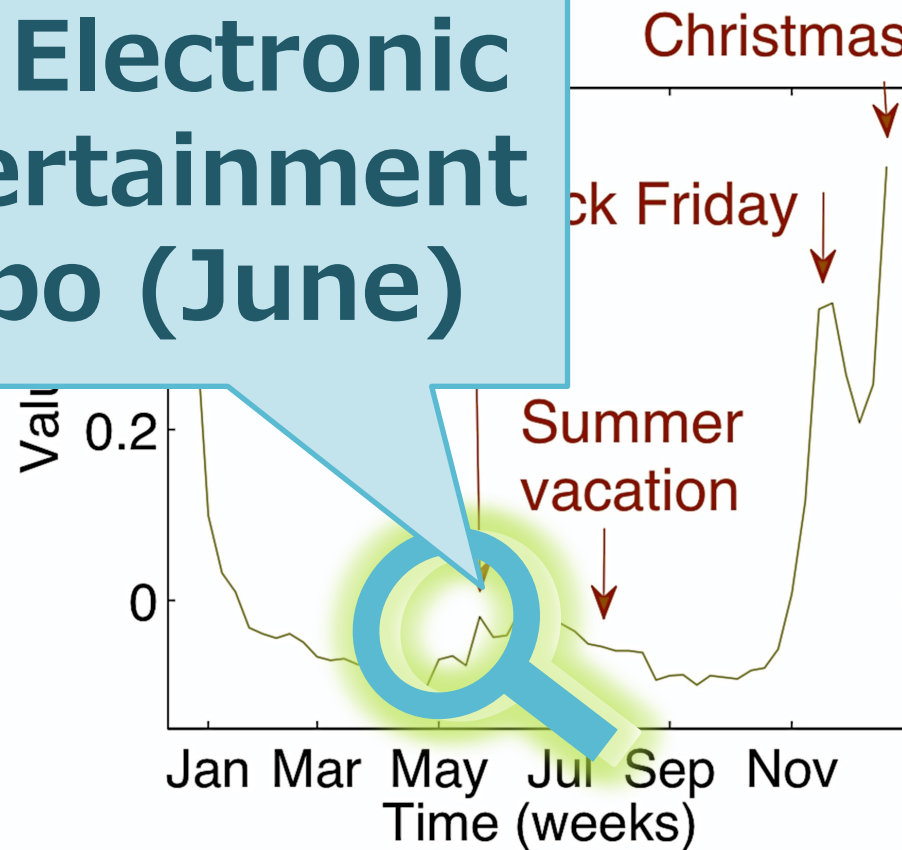
EcoWeb: seasonal component



Modeling power of EcoWeb



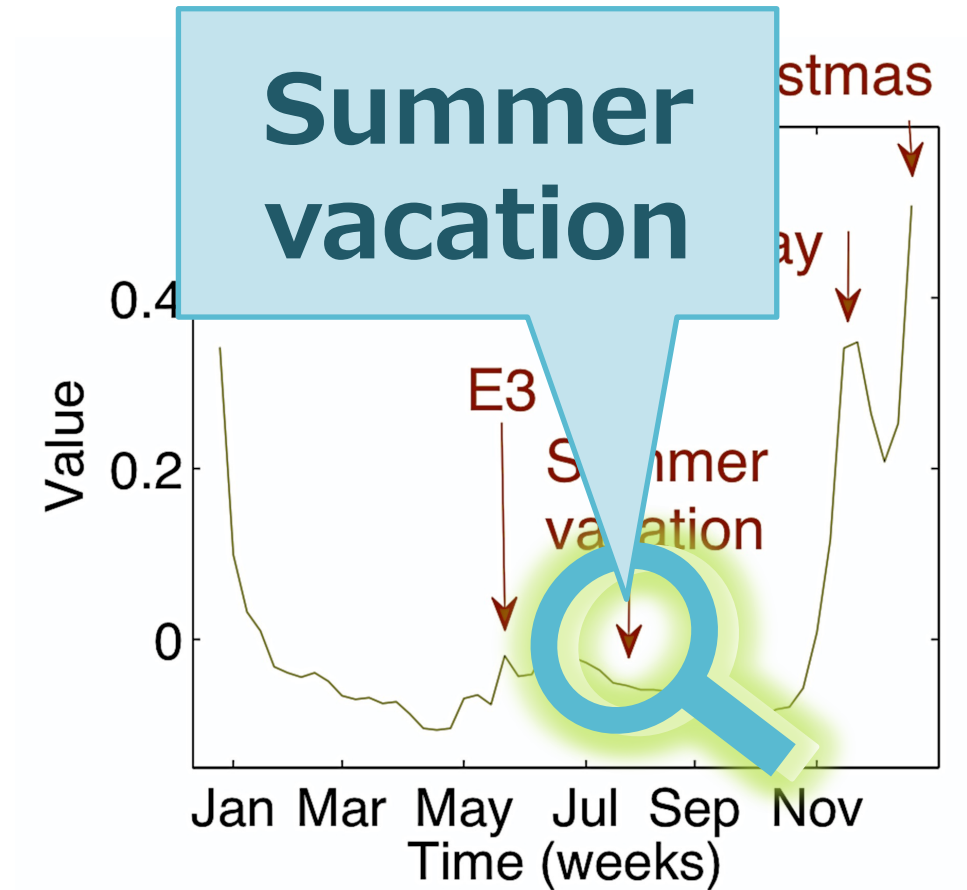
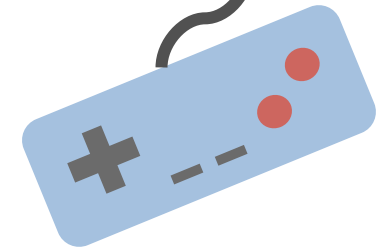
E3: Electronic Entertainment Expo (June)



EcoWeb: seasonal component



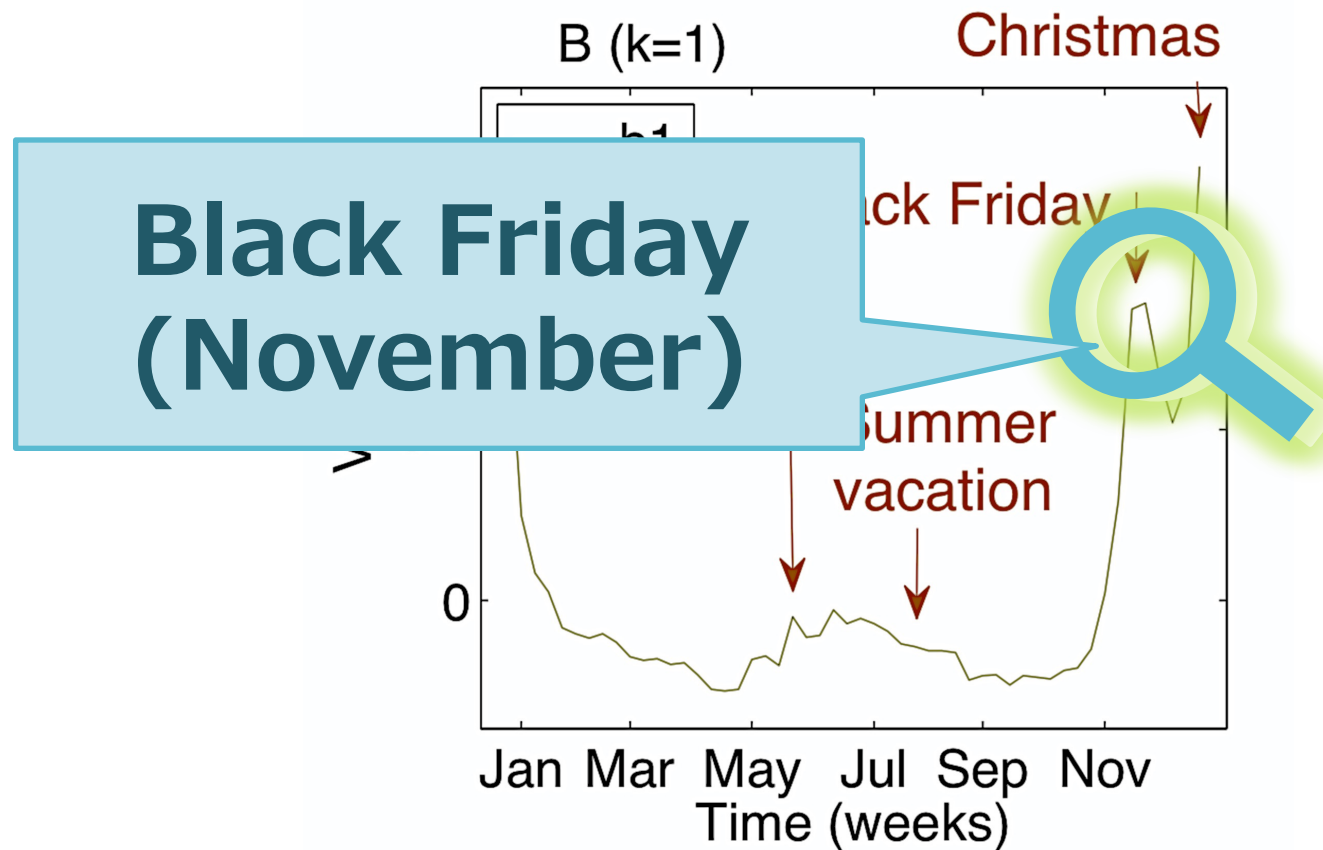
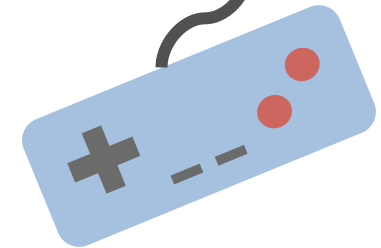
Modeling power of EcoWeb



EcoWeb: seasonal component



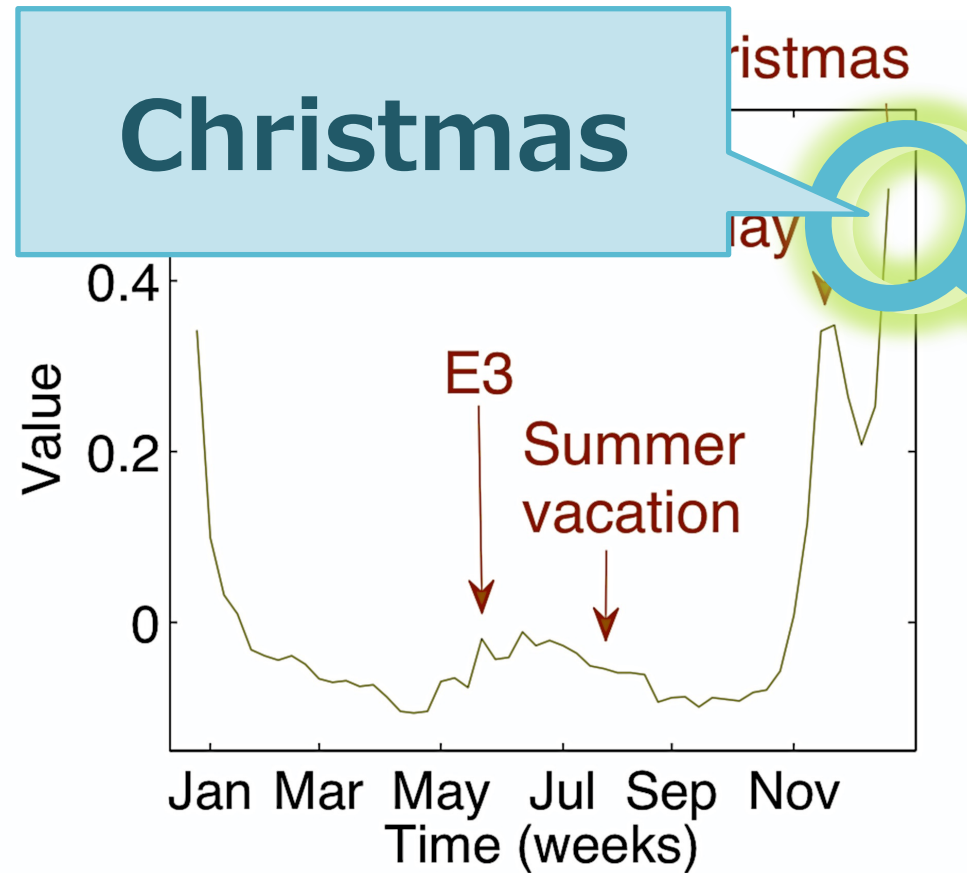
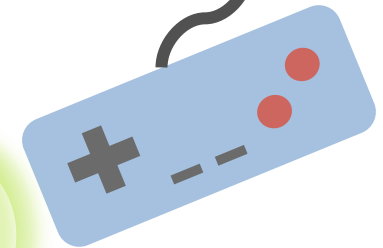
Modeling power of EcoWeb



EcoWeb: seasonal component



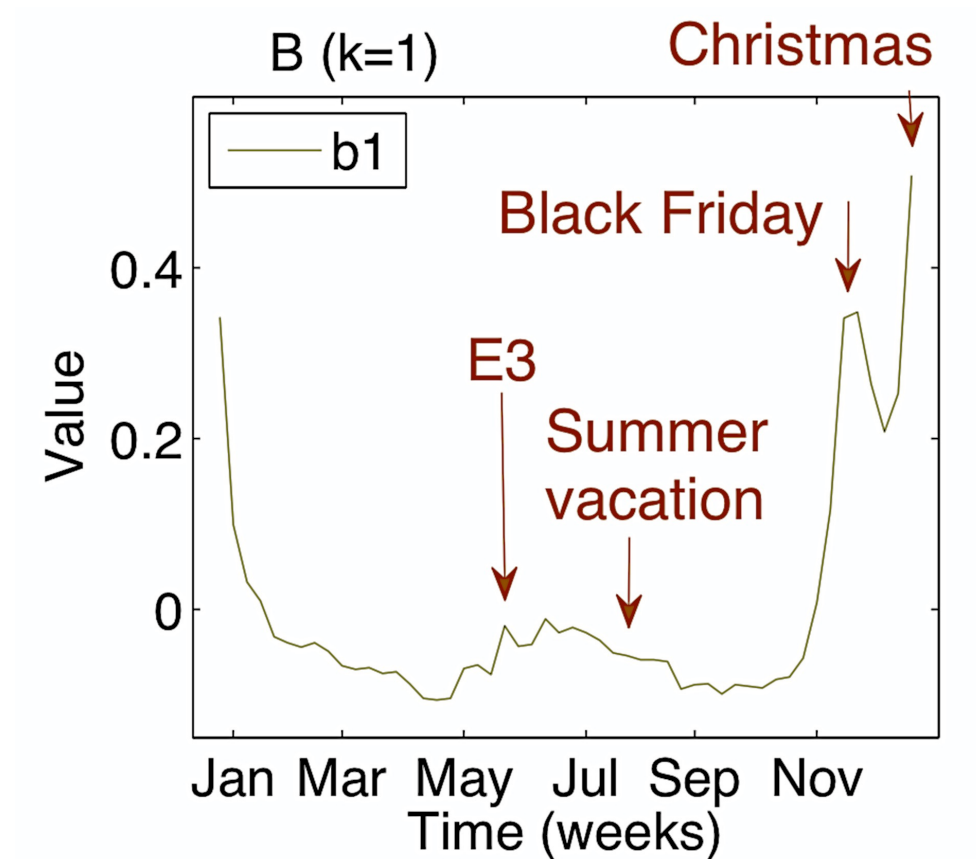
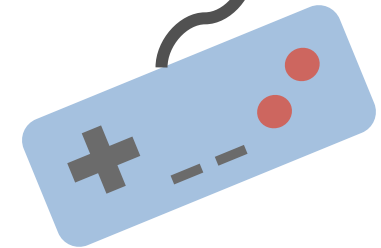
Modeling power of EcoWeb



EcoWeb: seasonal component



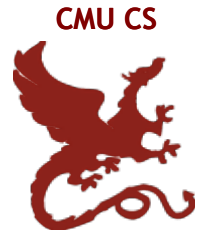
Modeling power of EcoWeb



EcoWeb: seasonal component

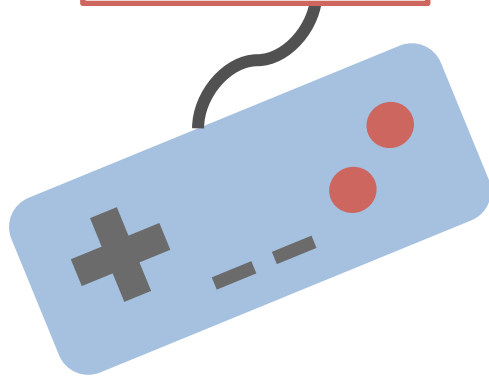


Modeling power of EcoWeb



Questions

Q1



Q2



Q3



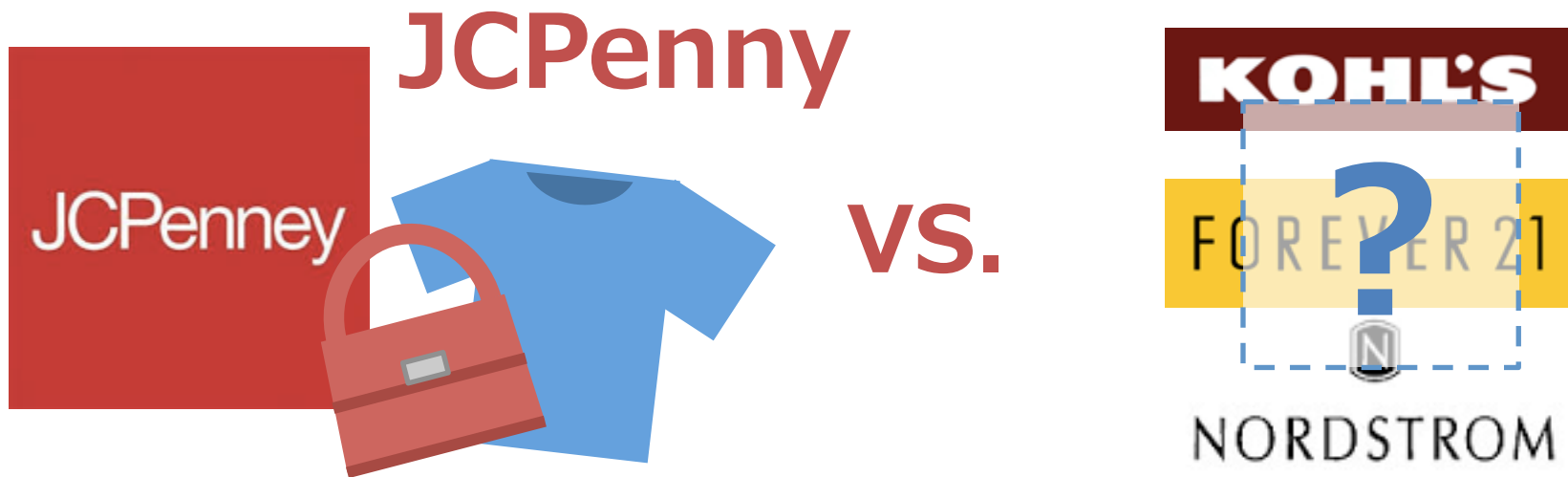


Modeling power of EcoWeb



Q2 (apparels)

Who is the competitor?

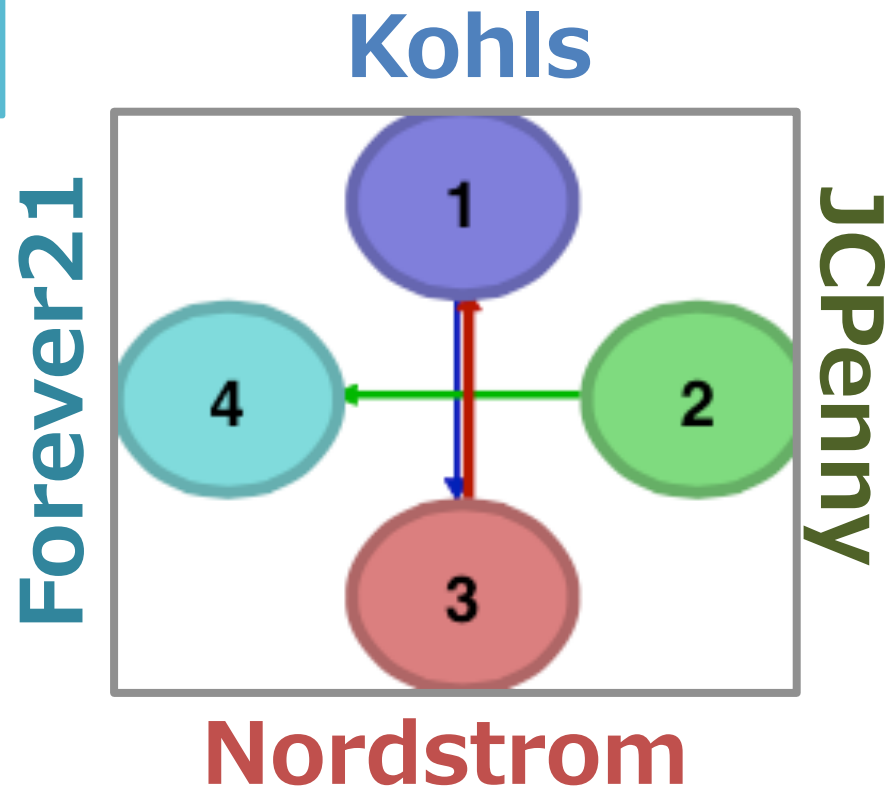
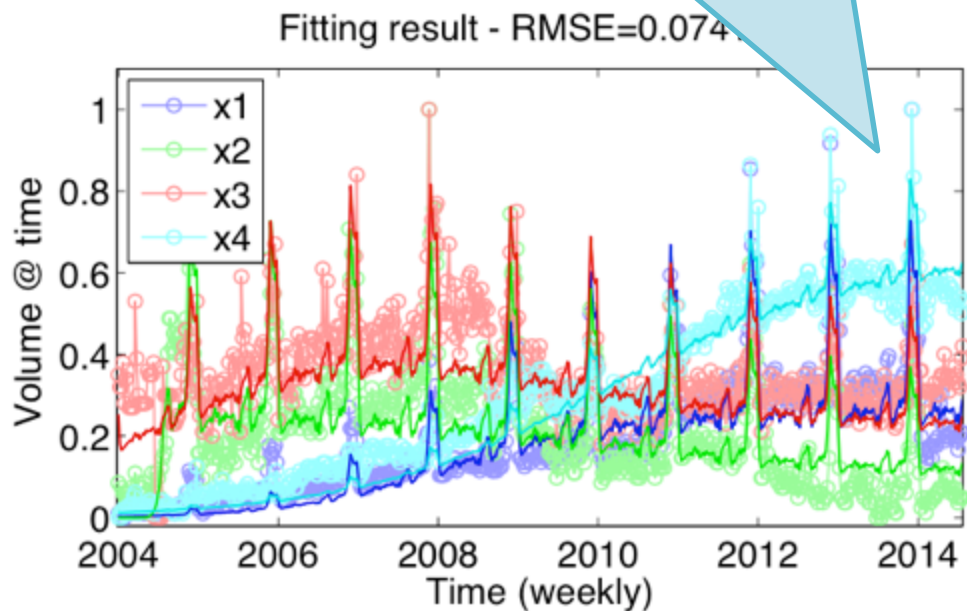




Modeling power of EcoWeb



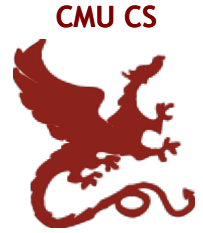
A2. Forever21!



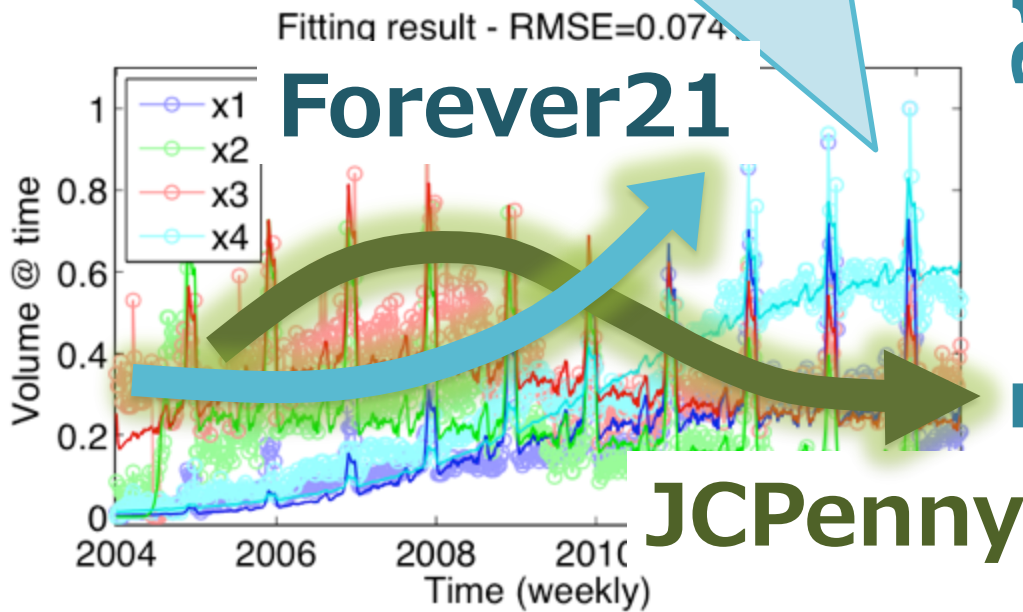
EcoWeb: Interaction network



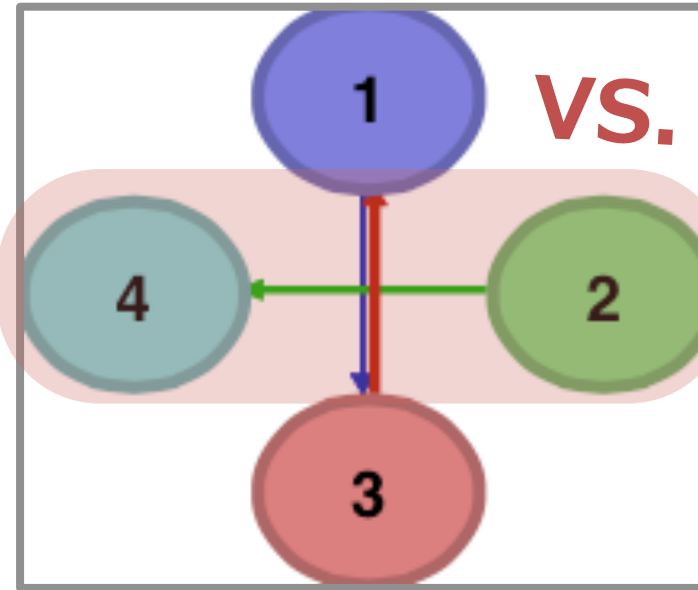
Modeling power of EcoWeb



A2. Forever21!



Forever21



Kohls

VS.

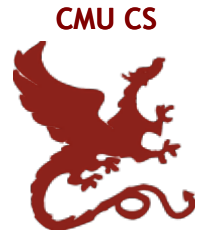
JCPenny

Nordstrom

EcoWeb: Interaction network



Modeling power of EcoWeb



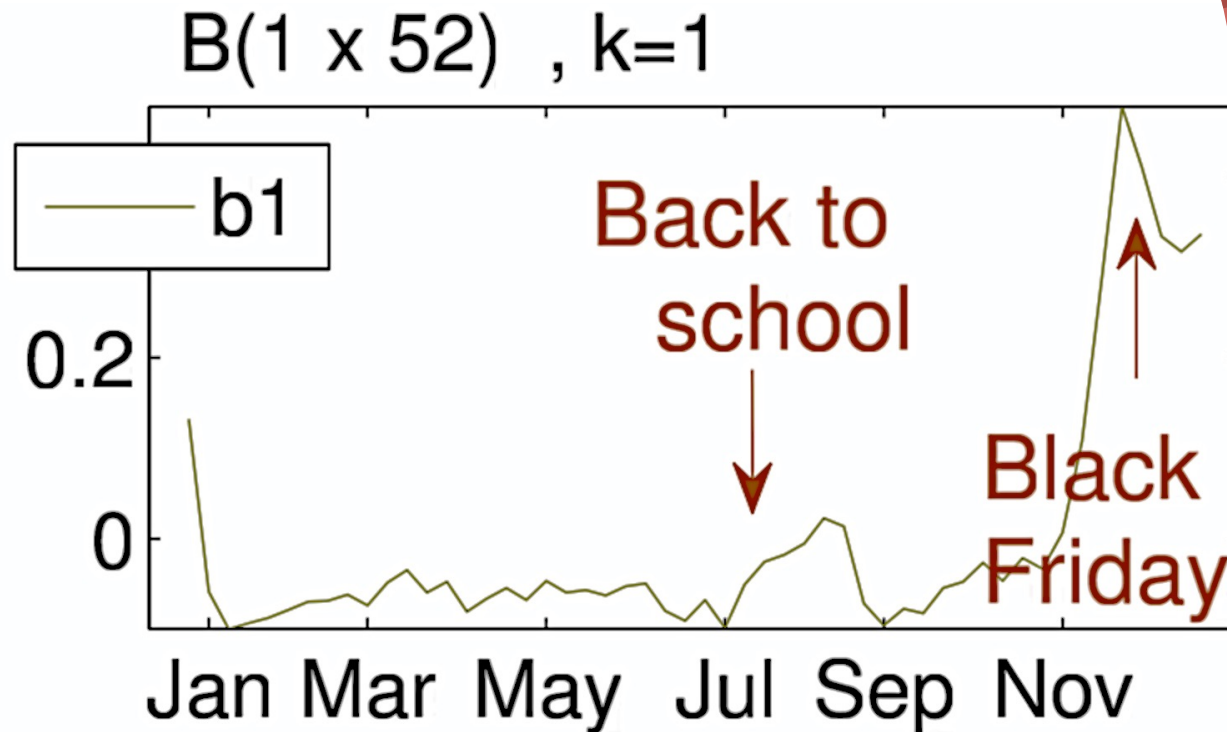
Q2 (apparels)

Any seasonal events?





Modeling power of EcoWeb



EcoWeb: seasonal component

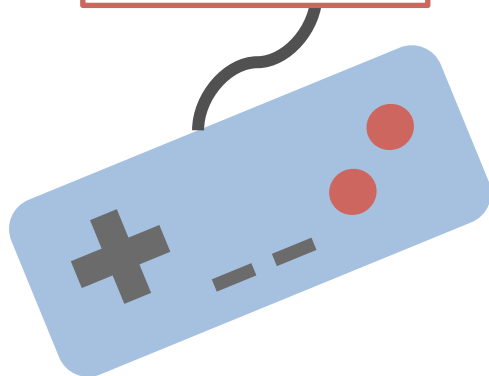


Modeling power of EcoWeb



Questions

Q1



Q2



Q3





Modeling power of EcoWeb



Q3 (retails)

Any patterns/trends?

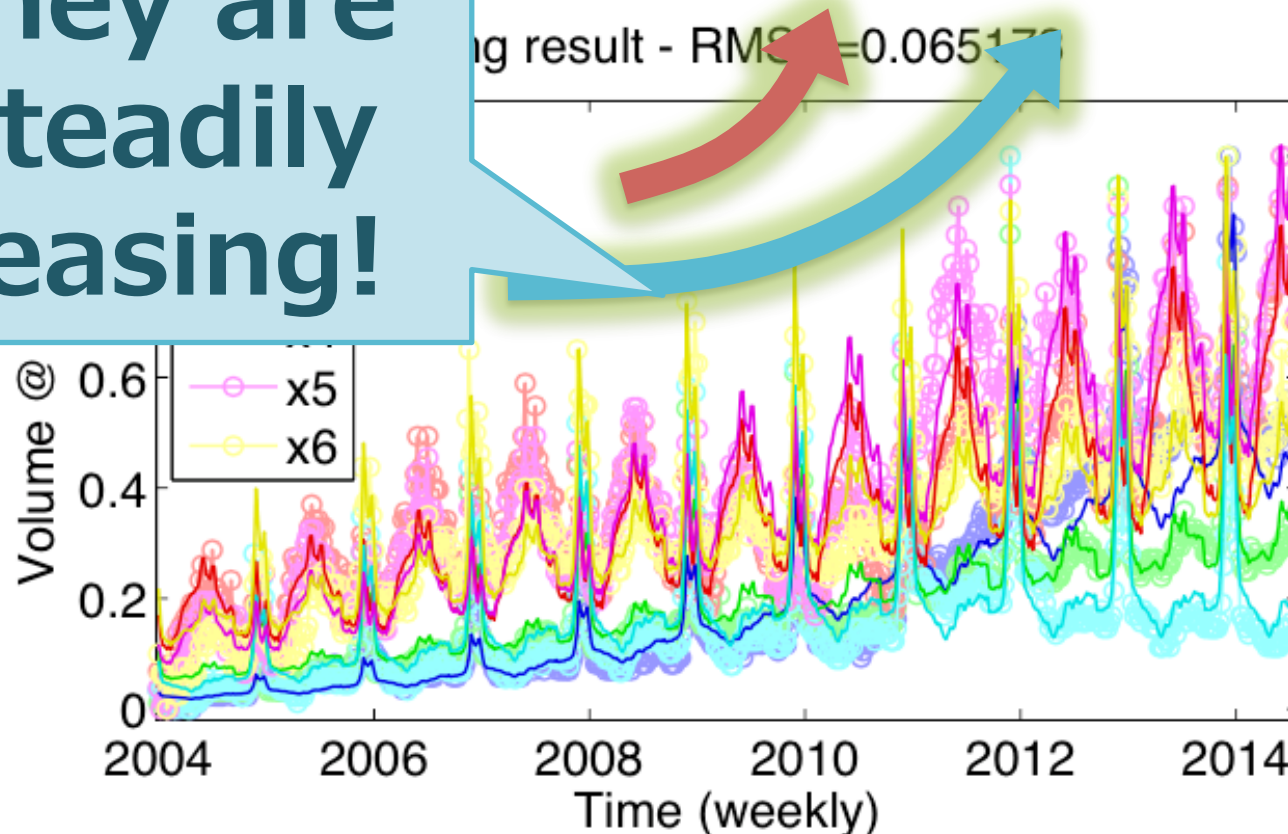




Modeling power of EcoWeb



A. They are all steadily increasing!

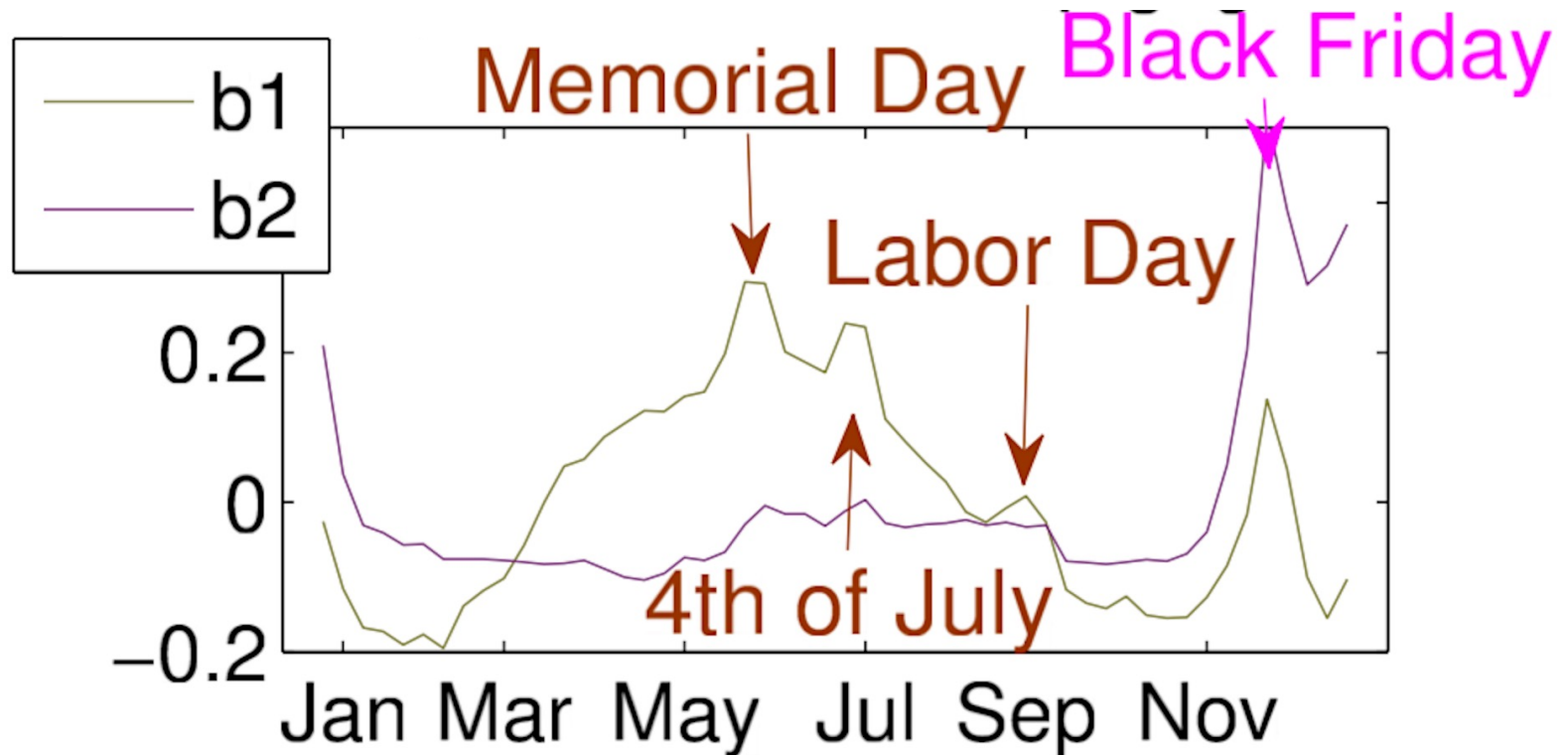


Amazon, Walmart, Home Depot, Best buy, ...



Modeling power of EcoWeb

2 seasonal components





Modeling power of EcoWeb



Black Friday sale



Black Friday



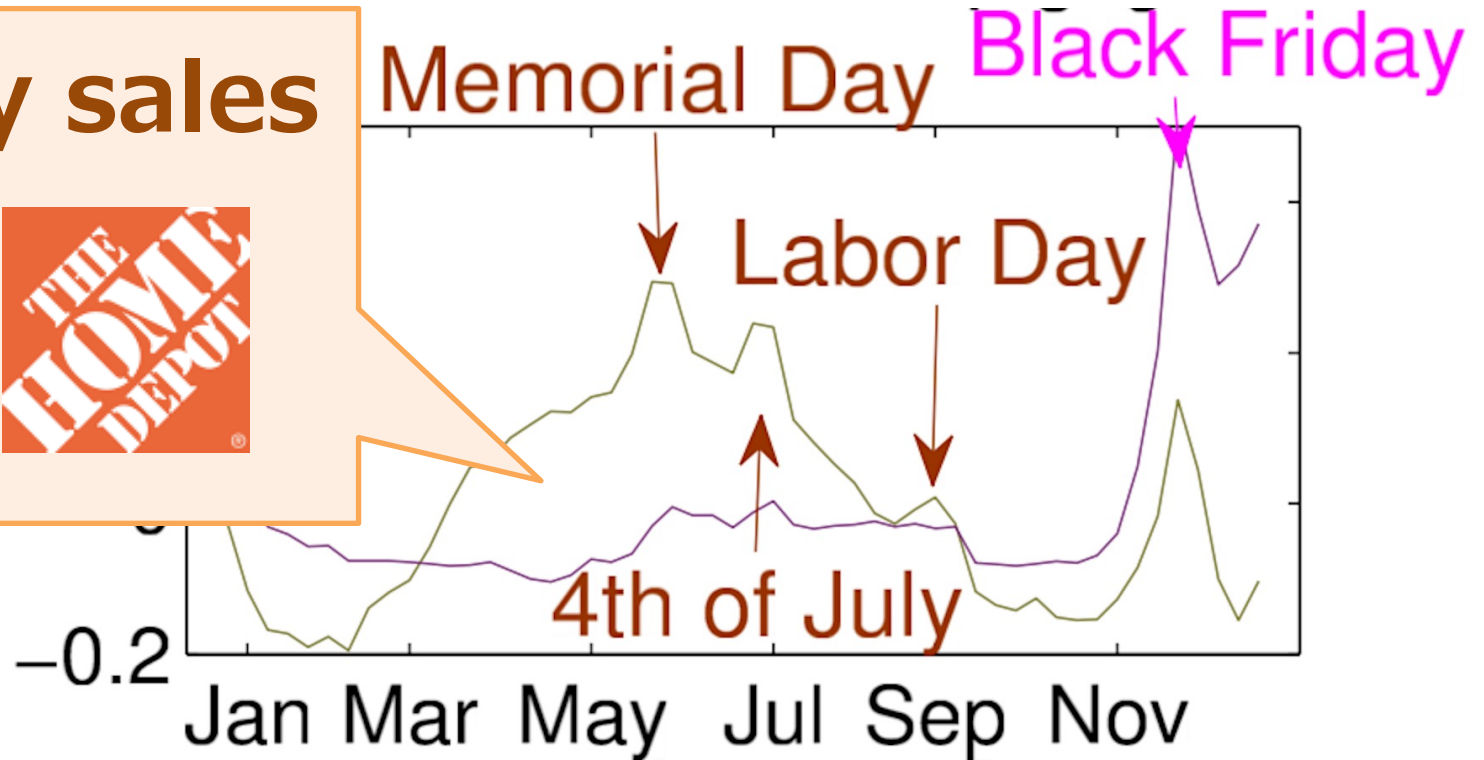


Modeling power of EcoWeb

2 seasonal components



Holiday sales

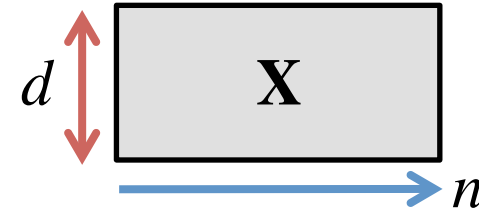





Problem definition

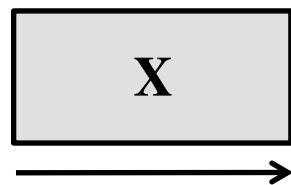
Given: Co-evolving online activities

X (activity \times time)

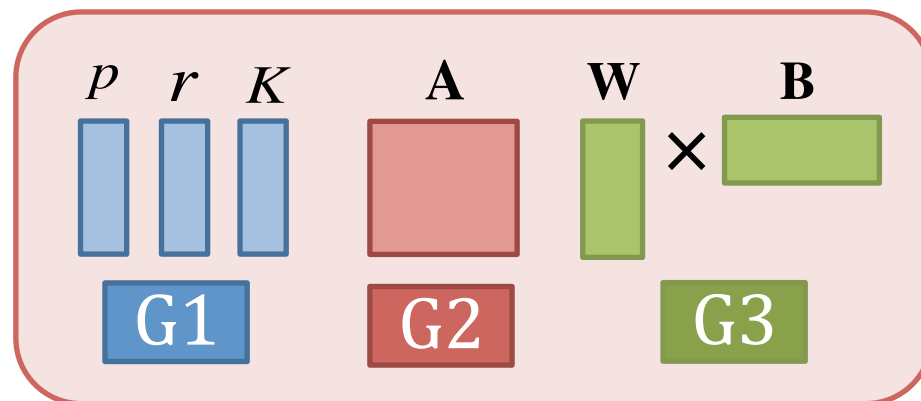


Find: Compact description of X

EcoWeb



\mathbb{R}



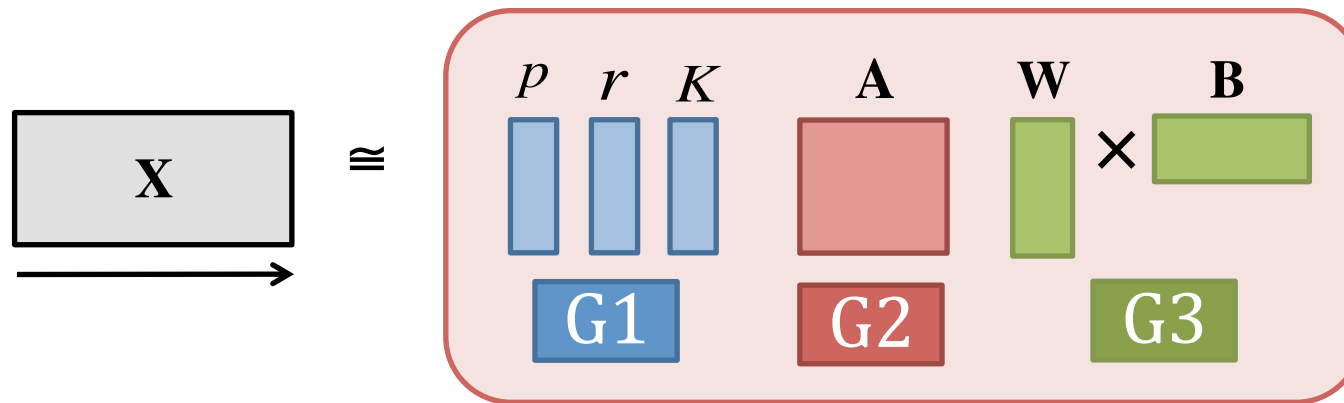


EcoWeb: Main idea



Q. How can we describe the evolutions of X ?

EcoWeb

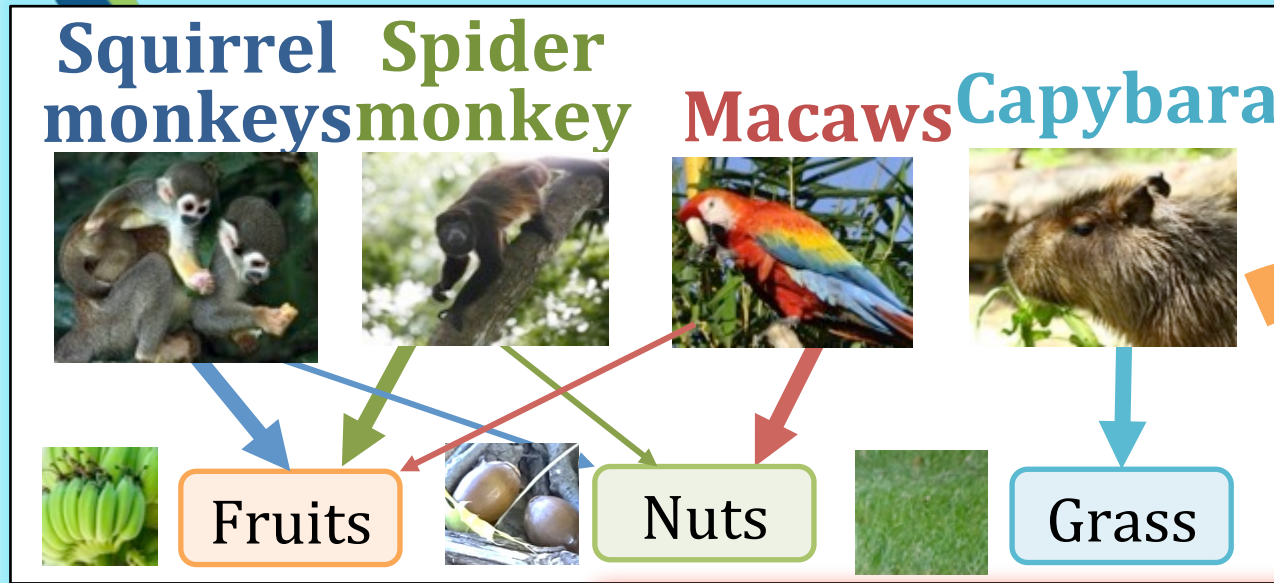
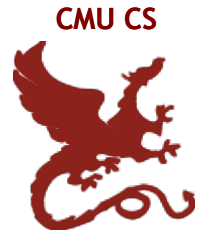


A. The Web as a jungle!

- “Virtual species” living on the Web
- Interacting with other species (activities)

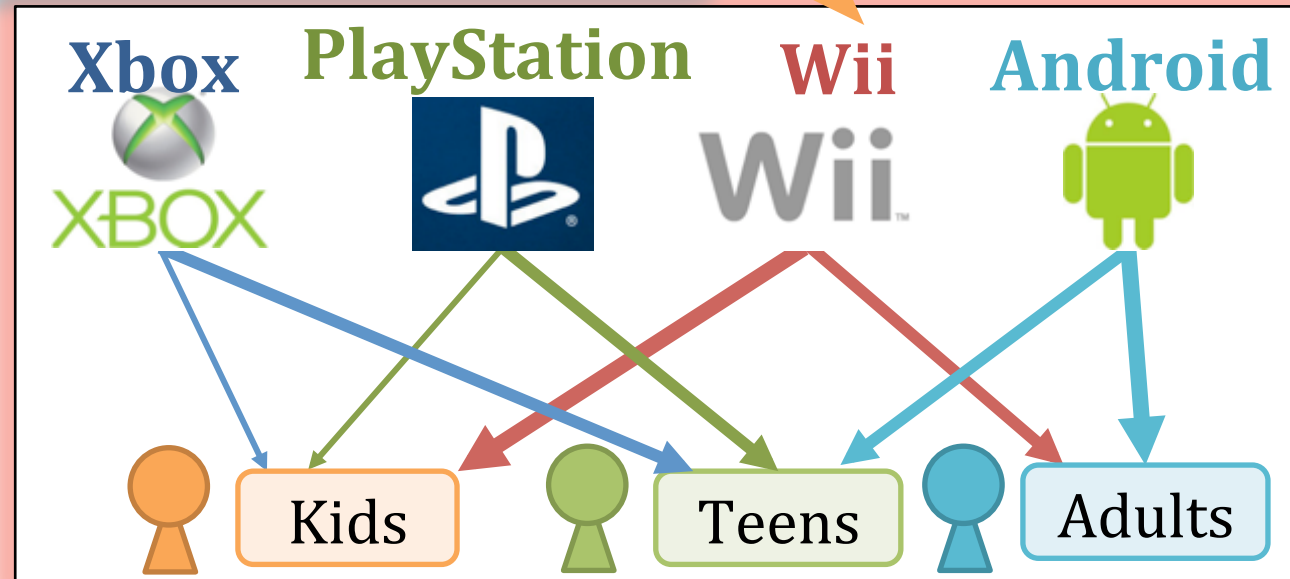


The Web as a jungle



Ecosystem on the Web

Ecosystem in the Jungle



Ecosystem on the Web

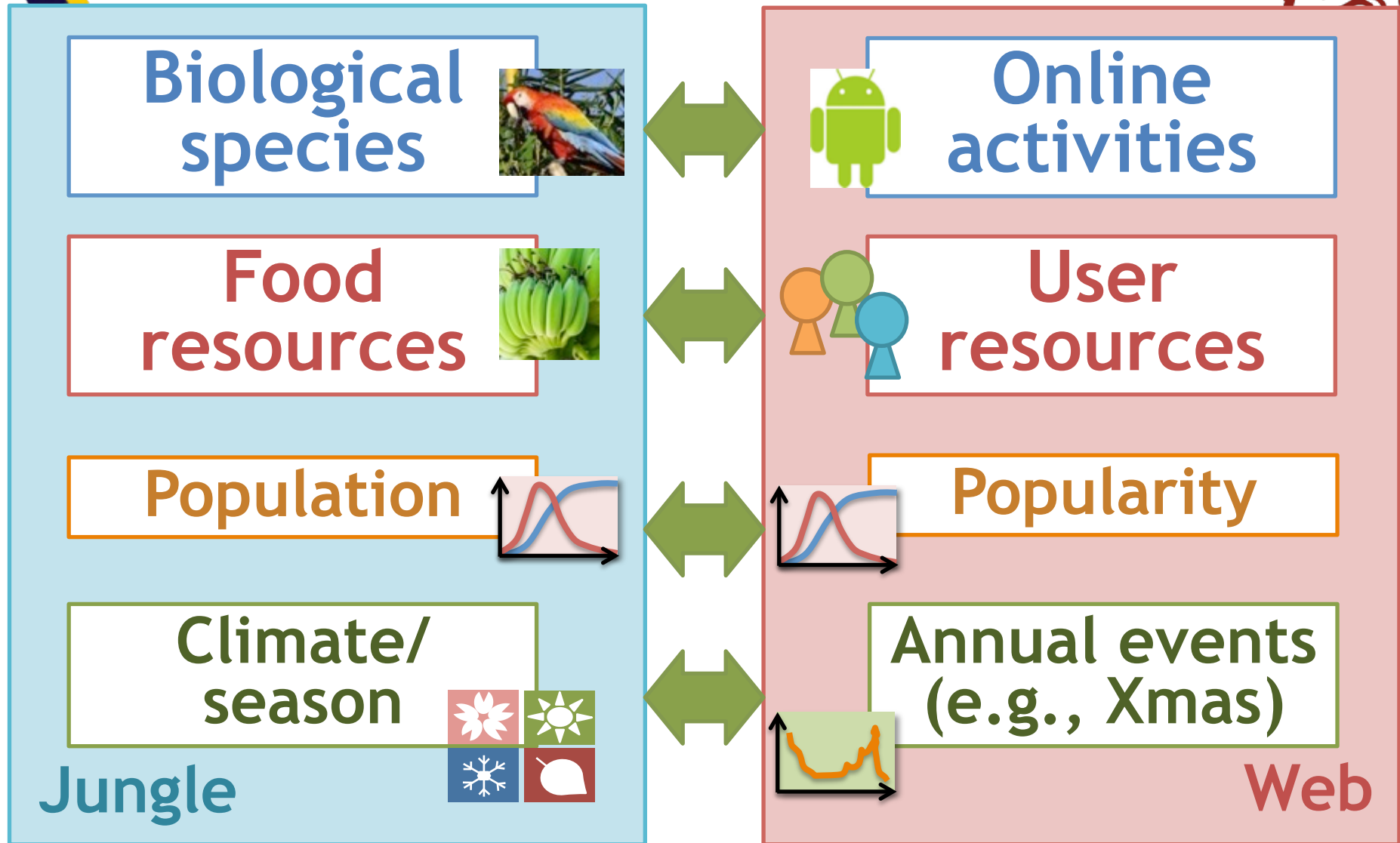


Image courtesy of xura, criminalatt, David Castillo Dominici, happykanppy at FreeDigitalPhotos.net.

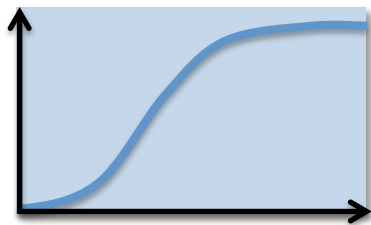


EcoWeb: Main idea

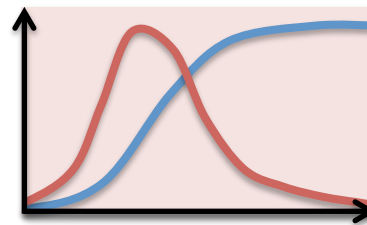


Q. How can we describe the evolutions of X ?

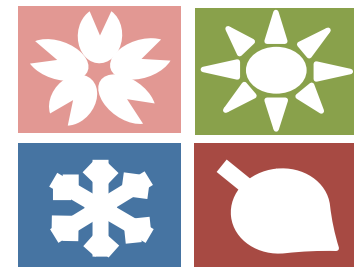
Non-linear evolution



Interaction/competition



Seasonality



A. Web as a jungle!

G1

G2

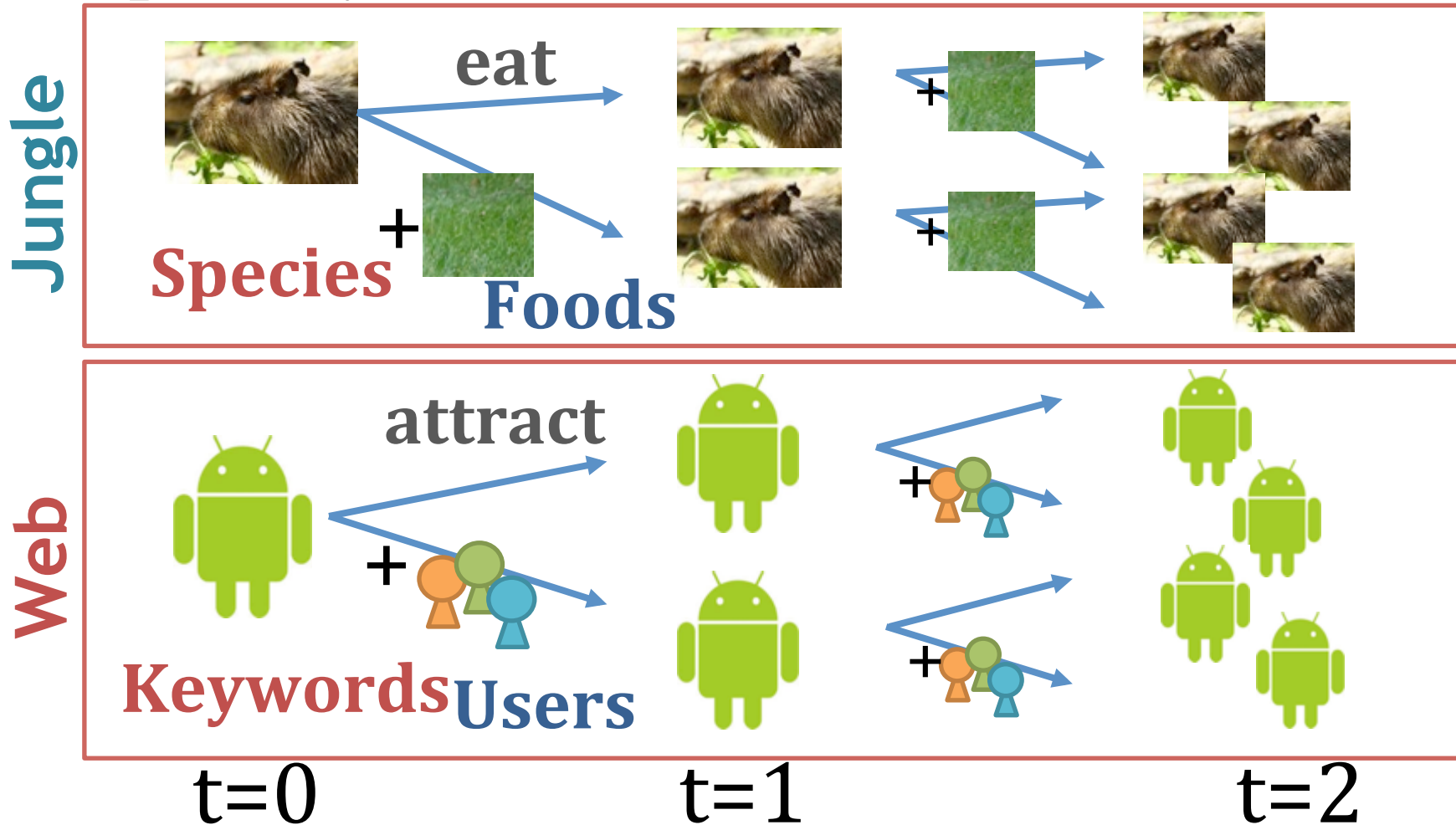
G3



G1: EcoWeb-individual



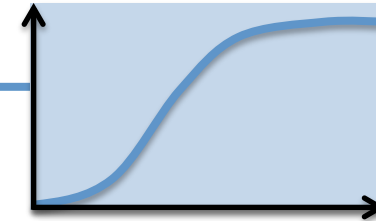
Popularity size increases over time



G1: EcoWeb-individual



Non-linear evolution of a single keyword



Popularity size

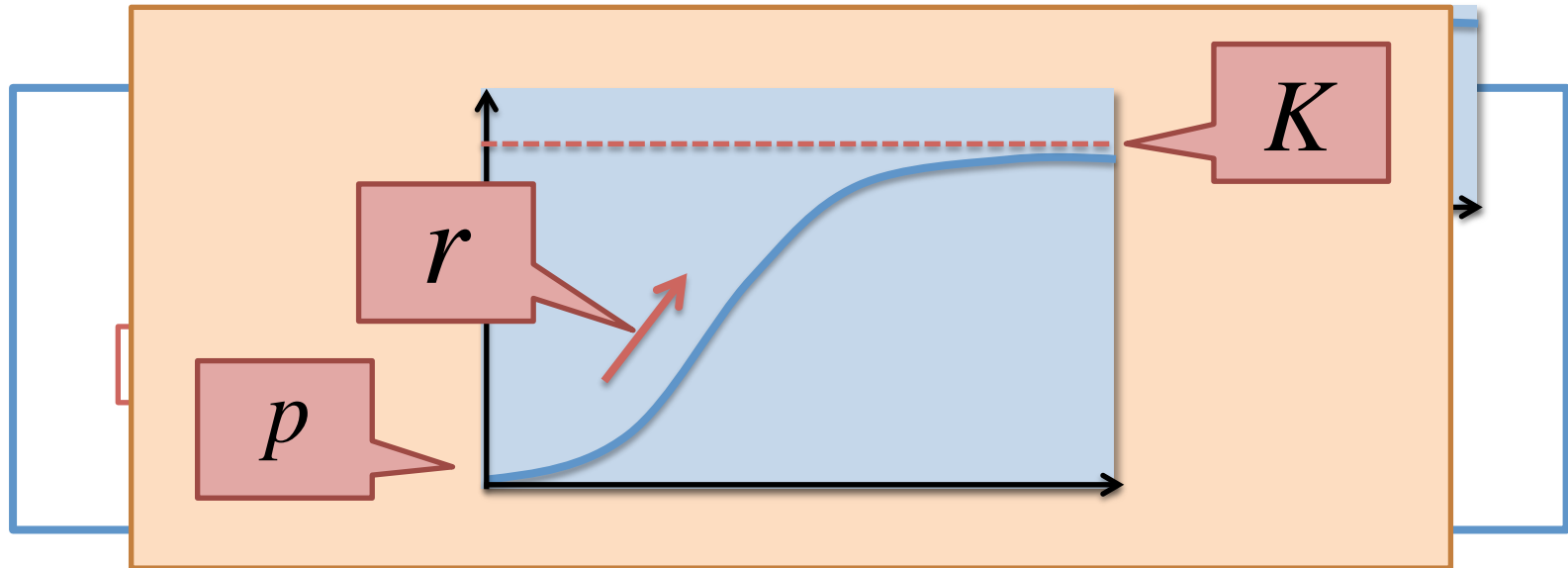
$$P(t + 1) = P(t) \left[1 + r \left(1 - \frac{P(t)}{K} \right) \right],$$

- p – Initial condition (i.e., $P(0) = p$)
- r – Growth rate, attractiveness
- K – Carrying capacity (=available user resources)

G1: EcoWeb-individual



Non-linear evolution of a single keyword



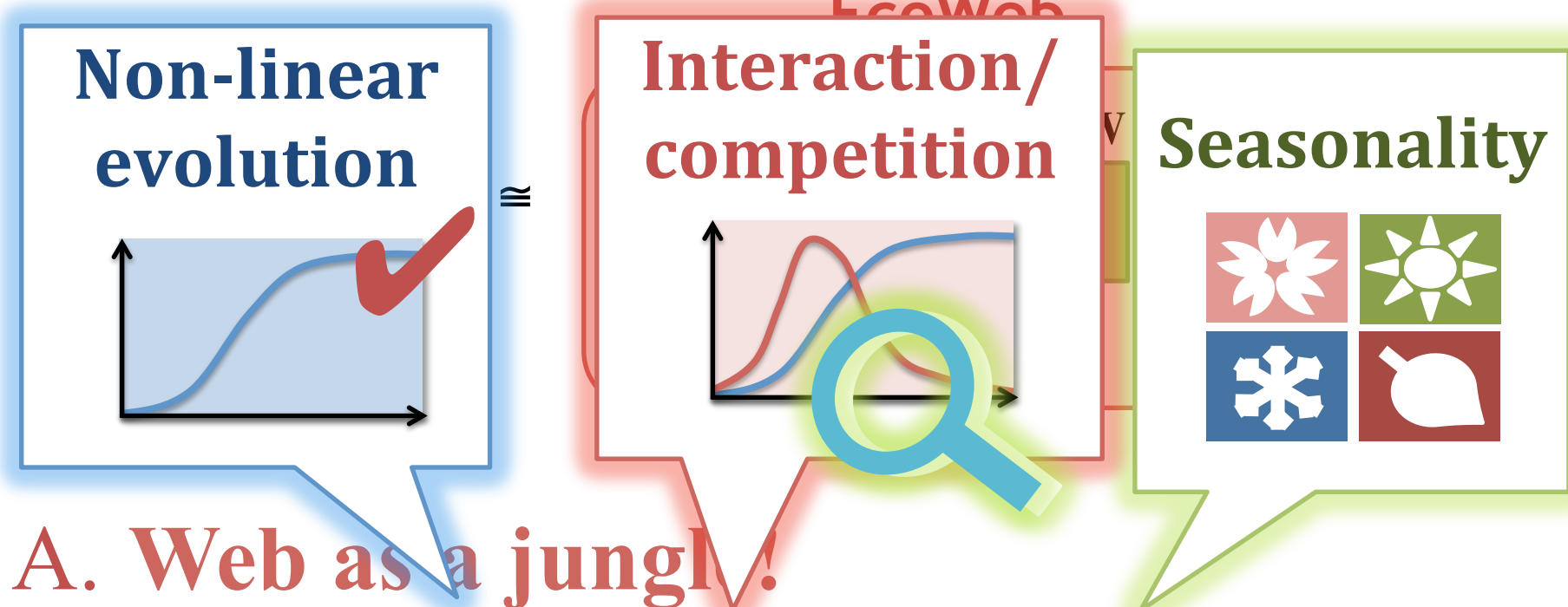
- p – Initial condition (i.e., $P(0) = p$)
- r – Growth rate, attractiveness
- K – Carrying capacity (=available user resources)



EcoWeb: Main idea



Q. How can we describe the evolutions of X ?



G1

G2

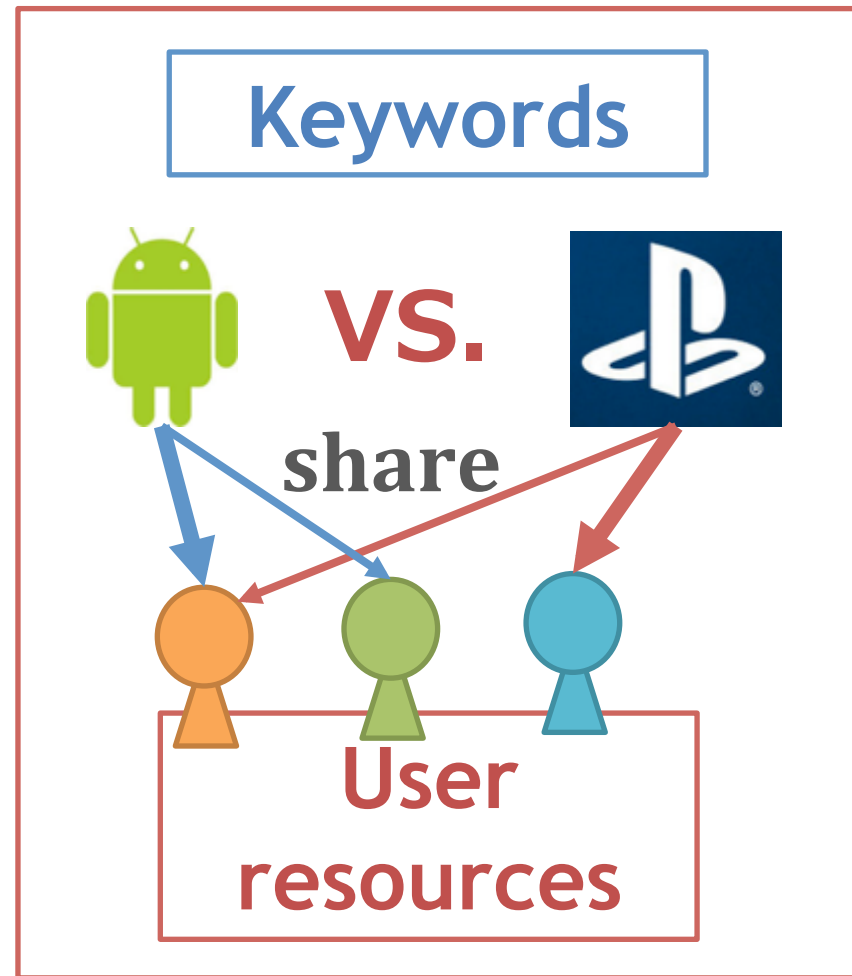
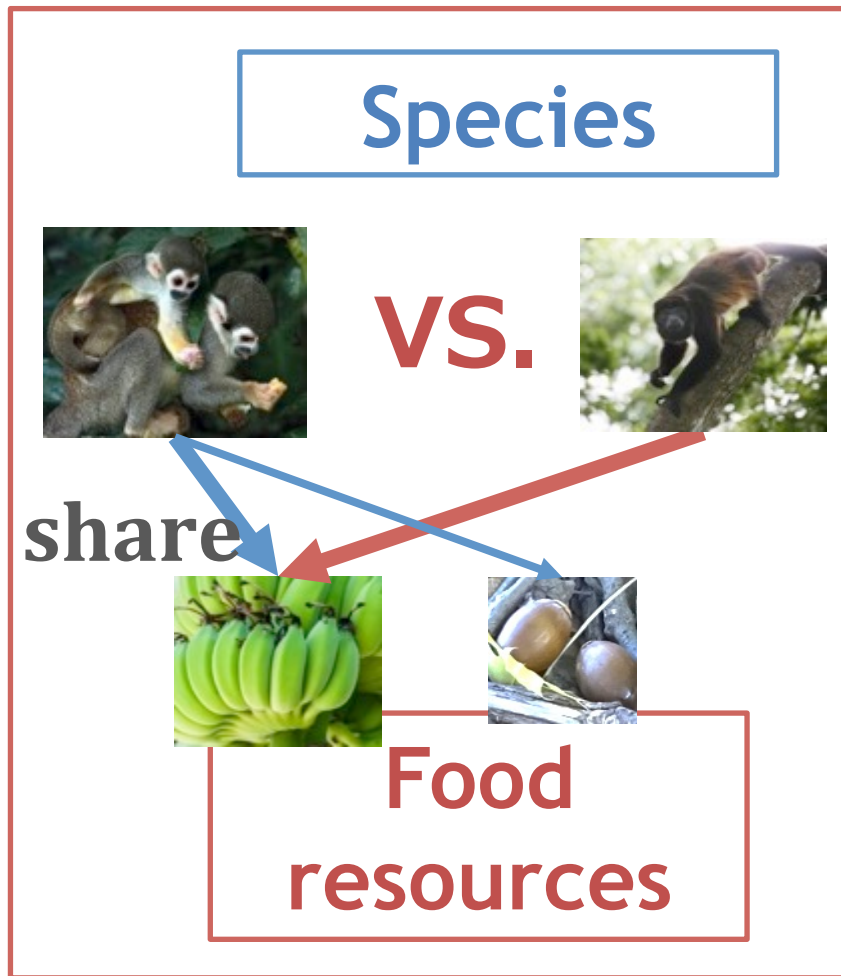
G3



G2: EcoWeb-interaction



Interaction between multiple keywords



G2: EcoWeb-interaction



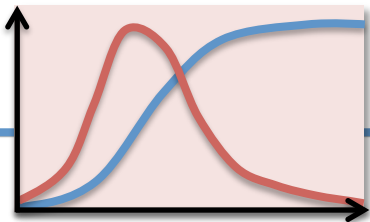
Interaction between multiple keywords

Popularity of keyword i

Popularity of j

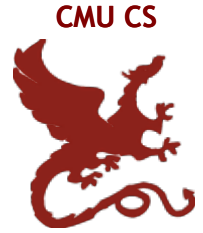
$$P_i(t + 1) = P_i(t) \left[1 + r_i \left(1 - \frac{\sum_{j=1}^d a_{ij} P_j(t)}{K_i} \right) \right],$$

$(i = 1, \dots, d), \quad (3)$

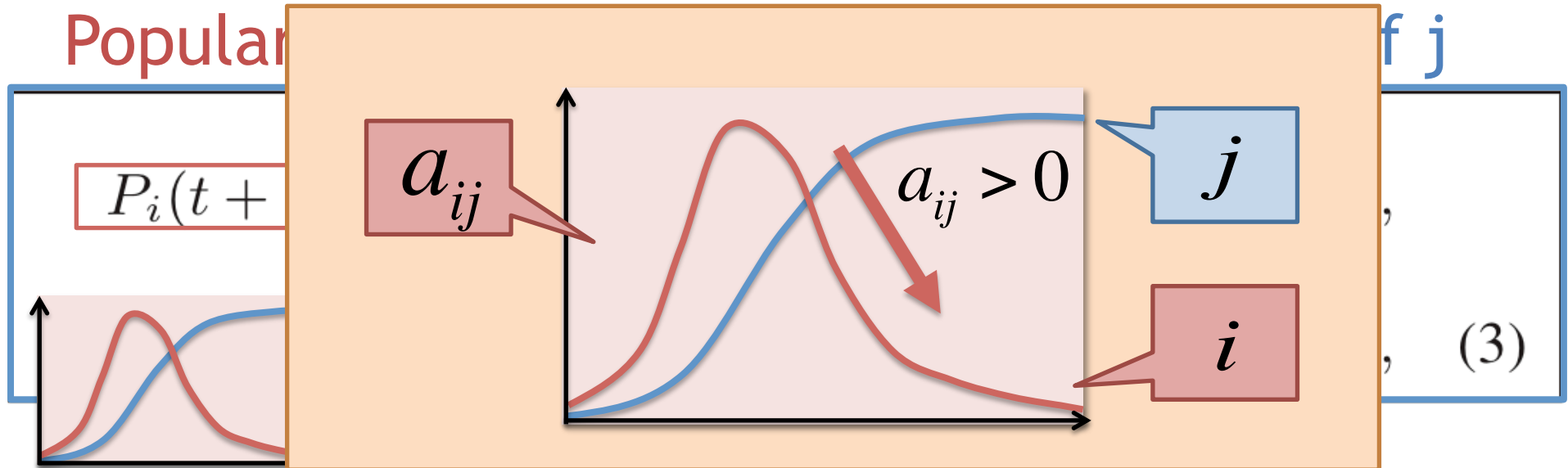


- a_{ij} – Interaction coefficient
- i.e., effect rate of keyword j on i

G2: EcoWeb-interaction



Interaction between multiple keywords



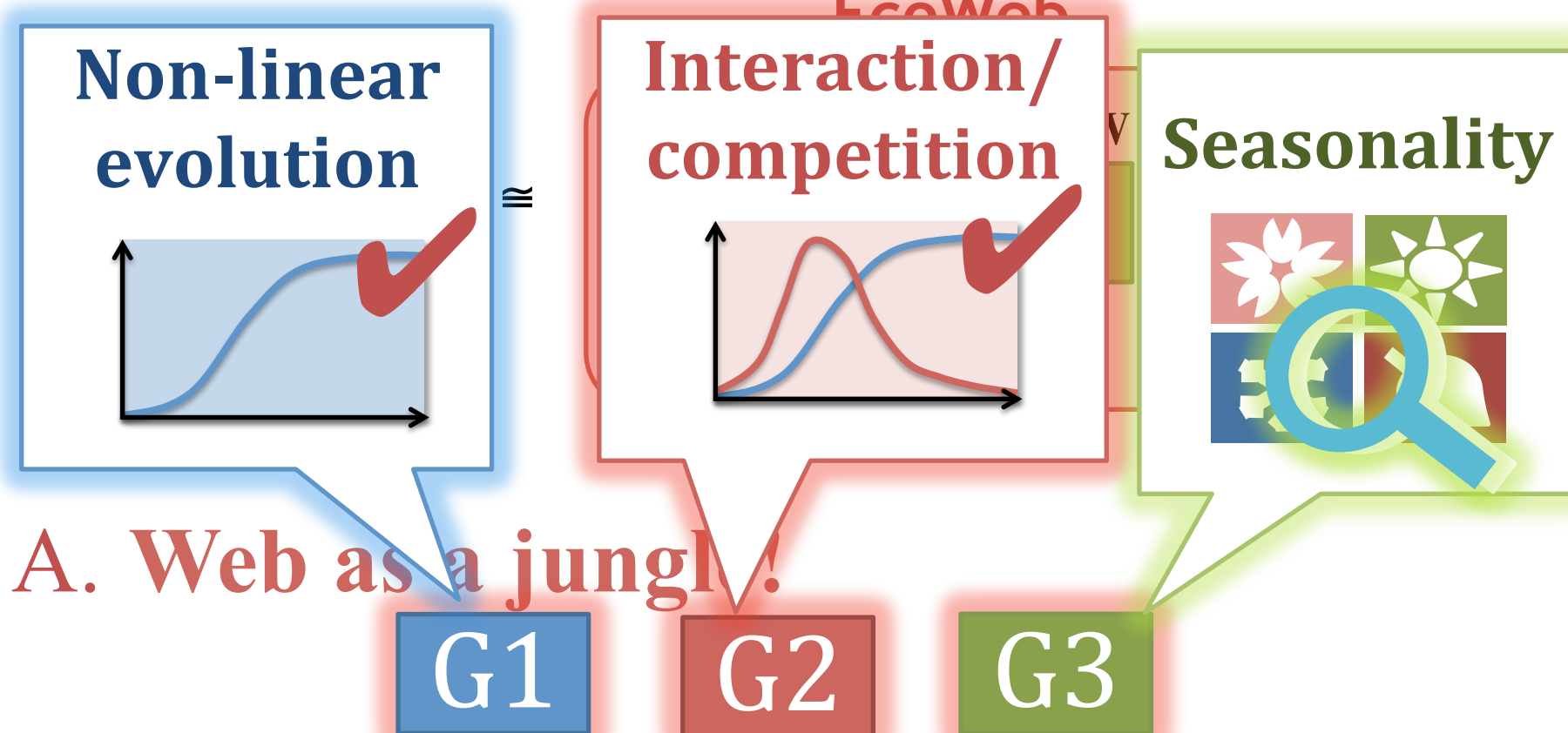
a_{ij} – Interaction coefficient
 – i.e., effect rate of keyword j on i



EcoWeb: Main idea



Q. How can we describe the evolutions of X ?





G3: EcoWeb-seasonality



“Hidden” seasonal activities

**Season/
Climate**

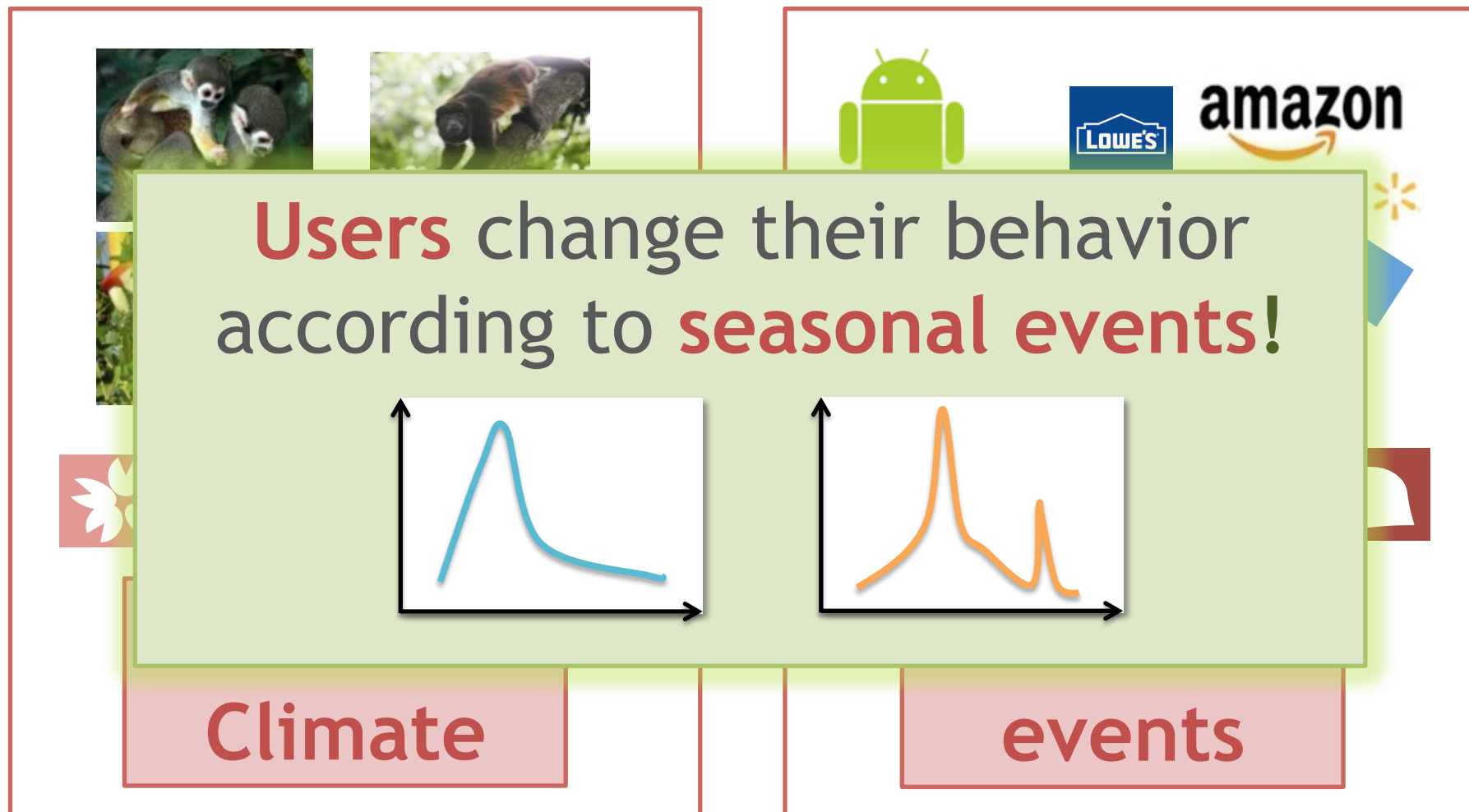
**Seasonal
events**



G3: EcoWeb-seasonality



“Hidden” seasonal activities



G3: EcoWeb-seasonality



“Hidden” seasonal activities

Estimated volume of keyword i

$$C_i(t) = P_i(t) [1 + e_i(t)] \quad (i = 1, \dots, d),$$

$$e_i(t) \simeq f(i, t | \mathbf{W}, \mathbf{B}) = \sum_{j=1}^k w_{ij} b_j(\tau) \quad (\tau = [t \bmod n_p])$$

Seasonal activities of i

W – Participation (weight) matrix

B – Seasonality matrix

G3: EcoWeb-seasonality



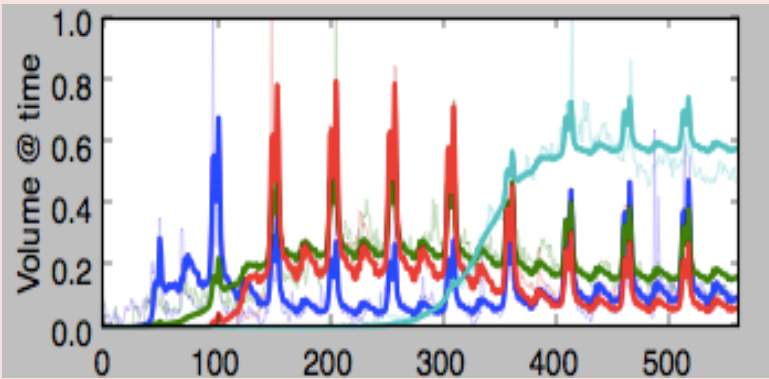
“Hidden” seasonal activities

Estimated volume of keyword i

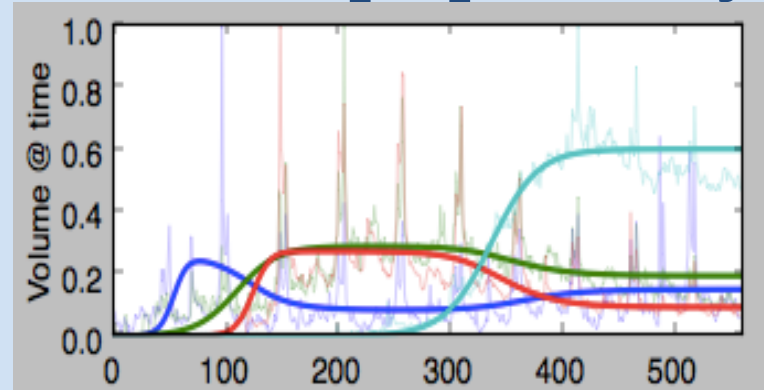
$$C_i(t) = P_i(t) [1 + e_i(t)] \quad (i = 1, \dots, d),$$

$$e_i(t) = f(i, t | \mathbf{W}, \mathbf{B}) = \sum_{k=1}^K w_k \cdot \text{[unclear]} [t - \text{[unclear]}]$$

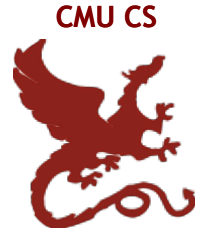
C: volume



P: latent popularity



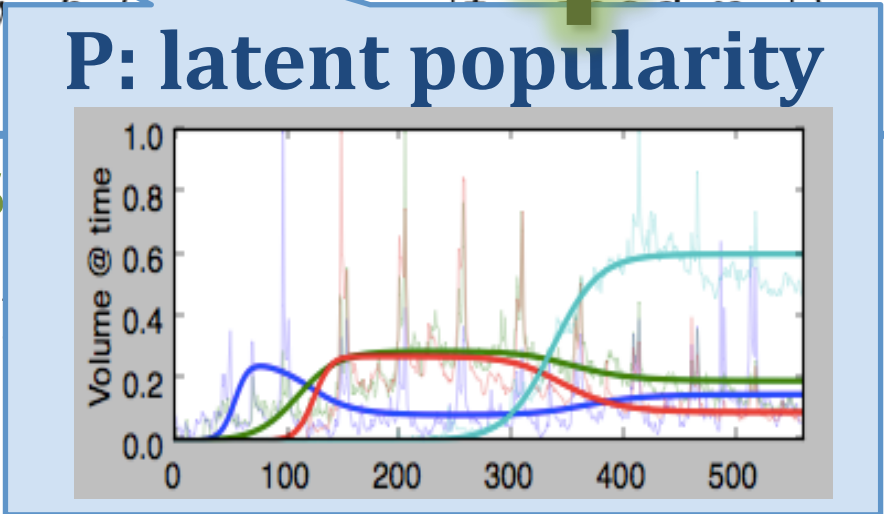
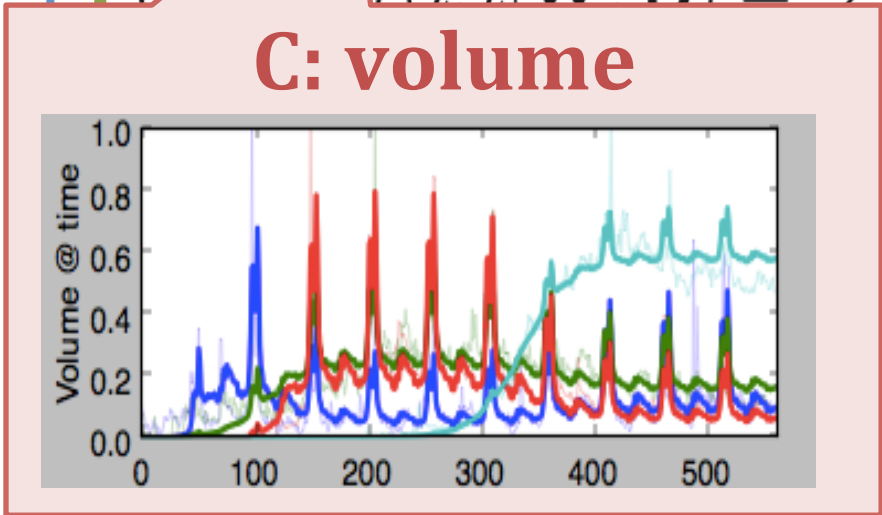
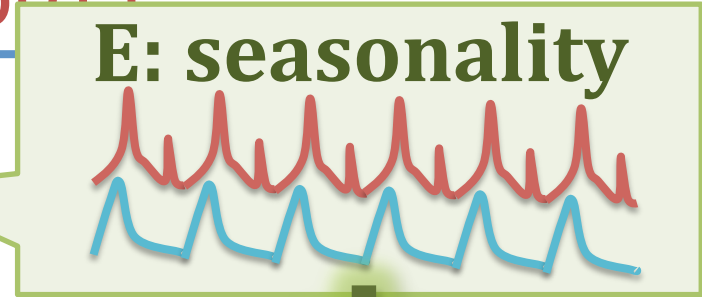
G3: EcoWeb-seasonality



“Hidden” seasonal activities

Estimated volume of keyword i

$$C_i(t) = P_i(t) [1 + e_i(t)]$$



G3: EcoWeb-seasonality



“Hidden” seasonal activities

Estimated volume of keyword i

$$C_i(t) = P_i(t) [1 + e_i(t)] \quad (i = 1, \dots, d),$$

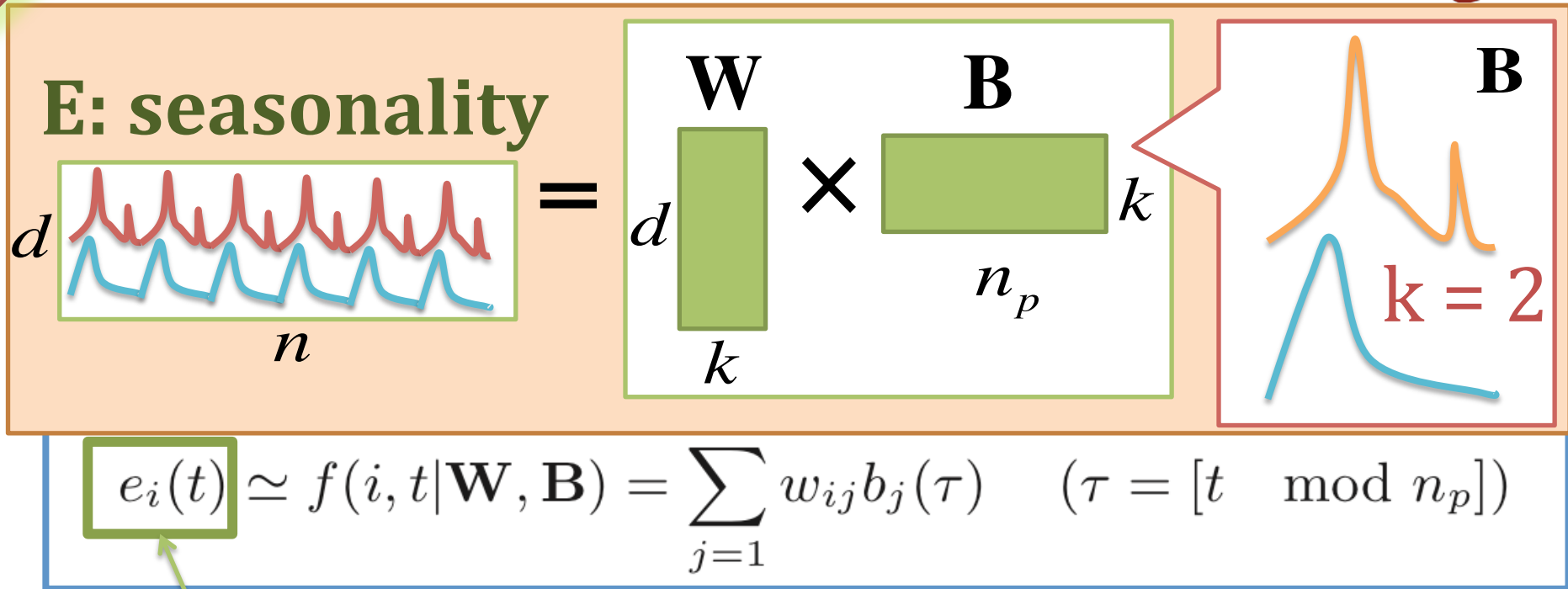
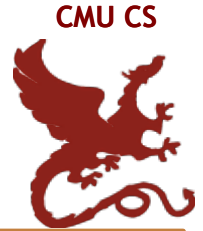
$$e_i(t) \simeq f(i, t | \mathbf{W}, \mathbf{B}) = \sum_{j=1}^k w_{ij} b_j(\tau) \quad (\tau = [t \bmod n_p])$$

Seasonal activities of keyword i

W – Participation (weight) matrix

B – Seasonality matrix

G3: EcoWeb-seasonality



Seasonal activities of keyword i

W – Participation (weight) matrix

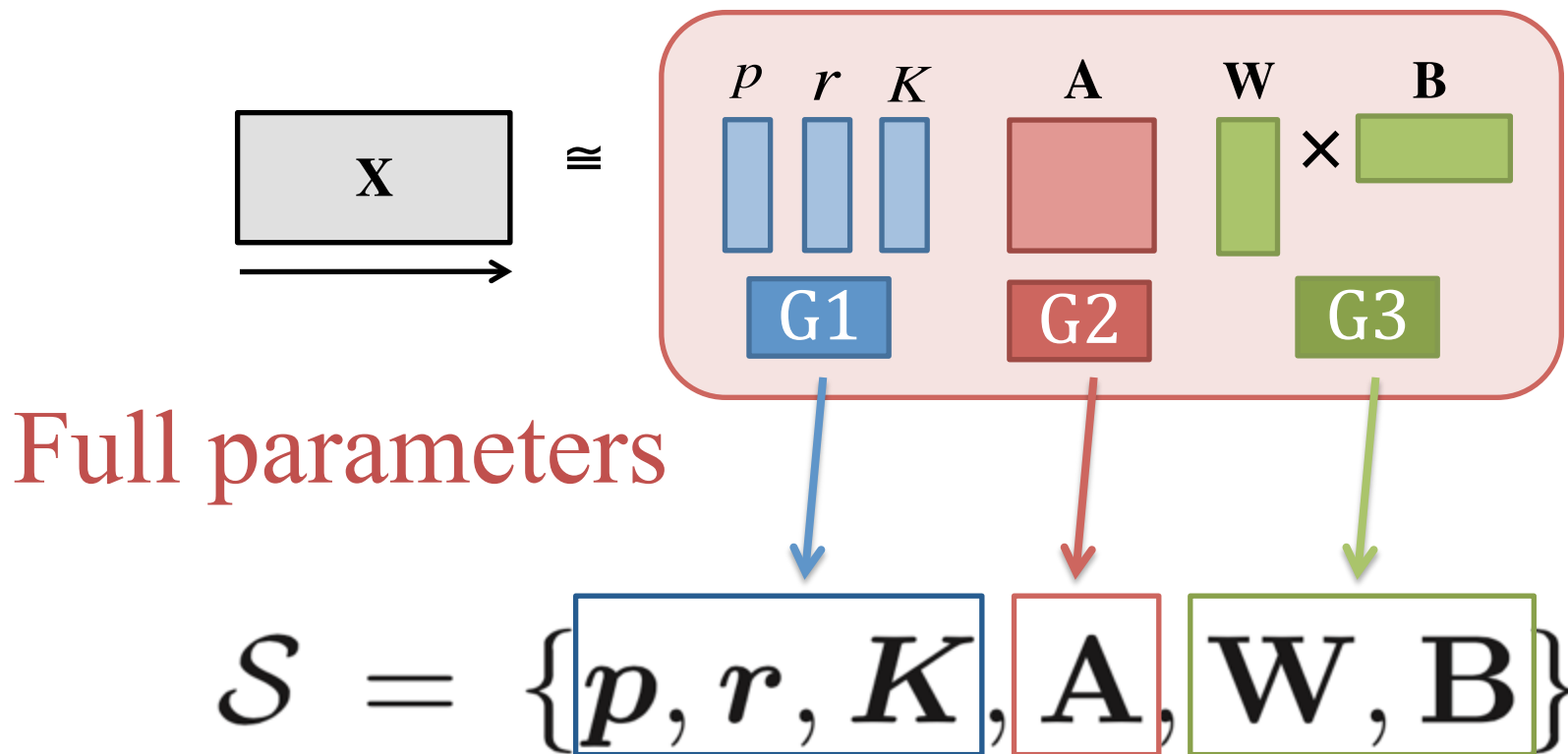
B – Seasonality matrix



EcoWeb: Main idea



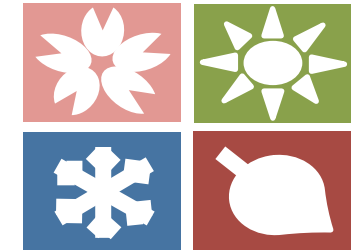
Q. How can we describe the evolutions of X ?
EcoWeb





Algorithms

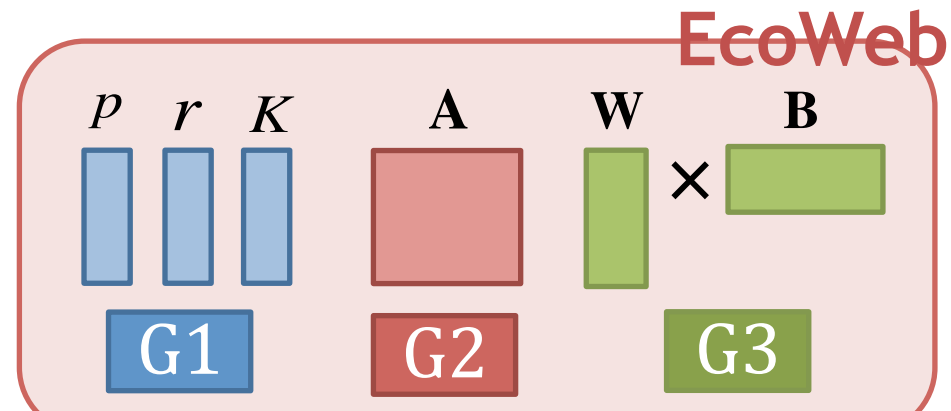
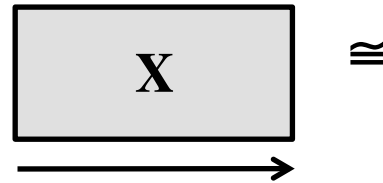
Q1. How can we automatically find “seasonal components” ?



Idea (1) : Seasonal component analysis

Q2. How can we efficiently estimate

full-parameters ?



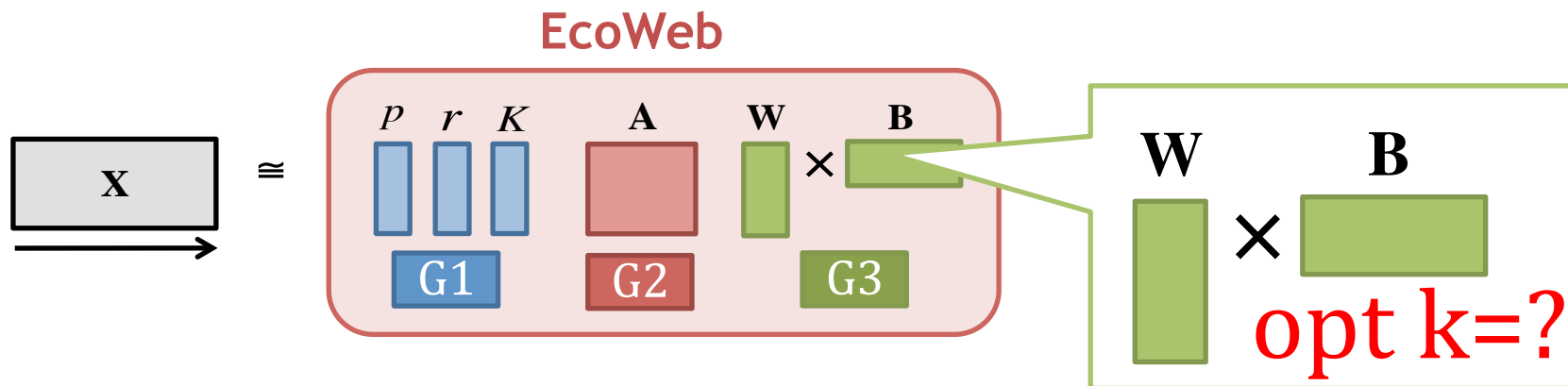
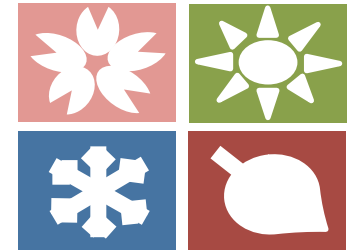
Idea (2): Multi-step fitting



Idea (1): Seasonal component analysis



Q1. How can we automatically find “k-seasonal components” ?



Idea (1) :

- a. Seasonal component detection
- b. Automatic component analysis

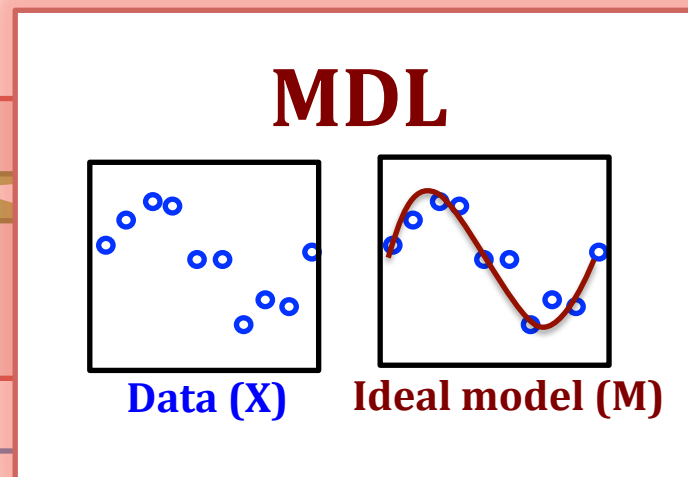
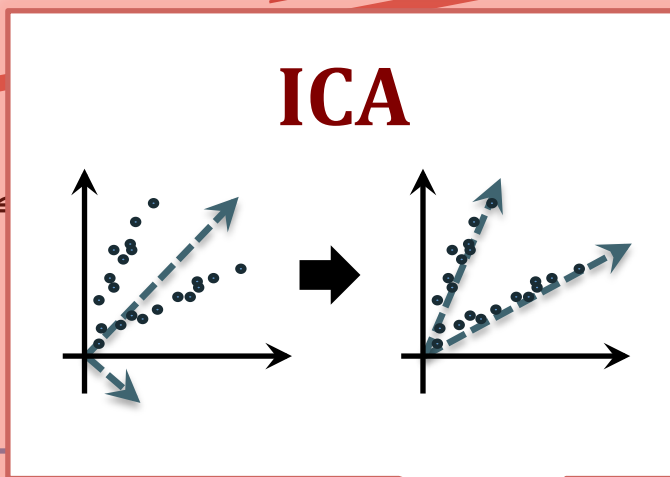
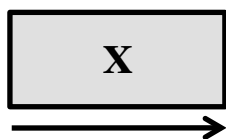
Idea (1): Seasonal component analysis

Q1. How can we automatically

find “

Details @ part1

components” ?



Idea (1) :

- a. Seasonal component detection
- b. Automatic component analysis

ICA

MDL

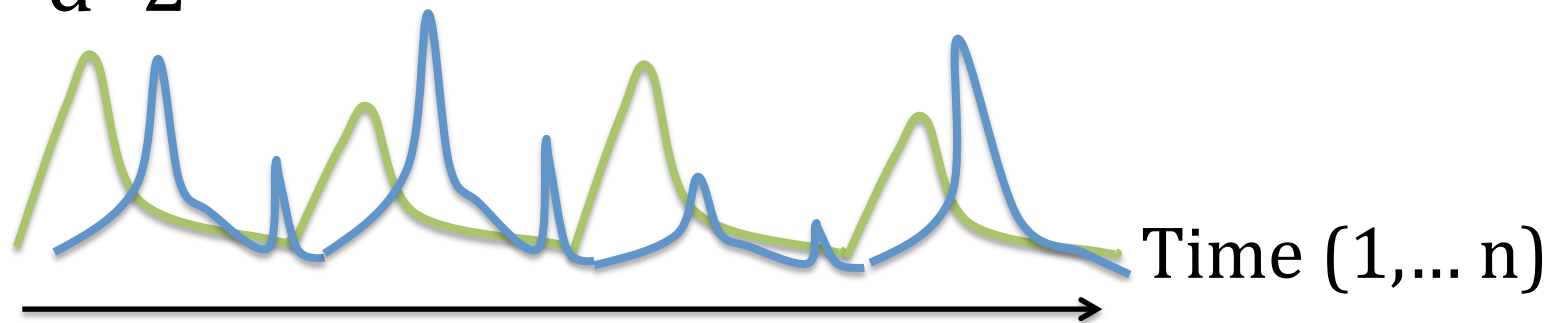


Idea (1): Seasonal component analysis



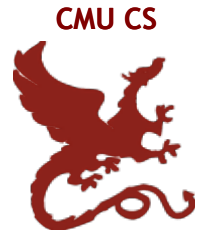
Idea(1-a) Seasonal component detection

E $d=2$

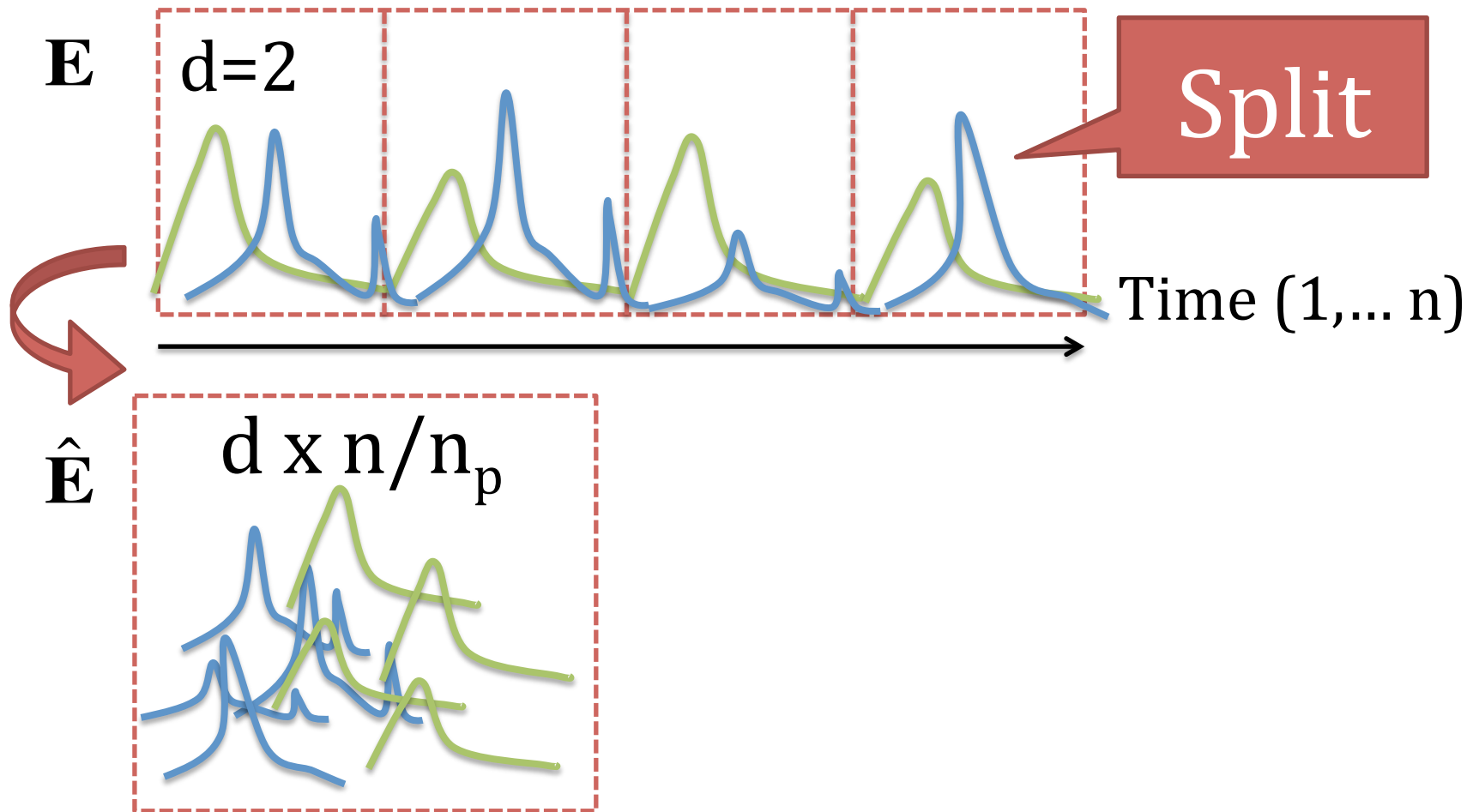




Idea (1): Seasonal component analysis



Idea(1-a) Seasonal component detection

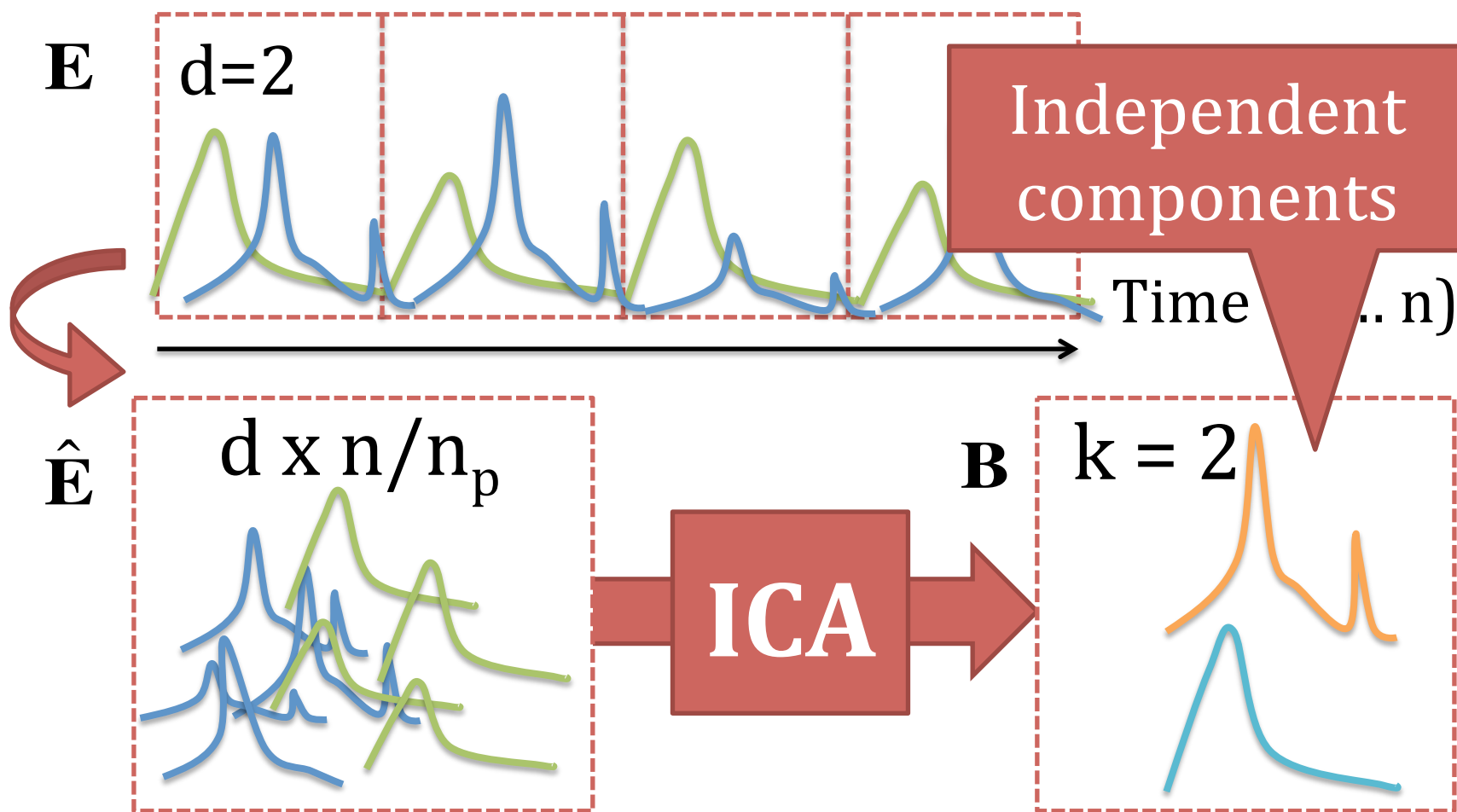




Idea (1): Seasonal component analysis

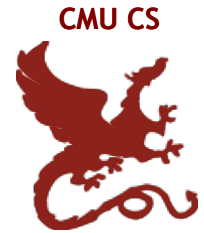


Idea(1-a) Seasonal component detection





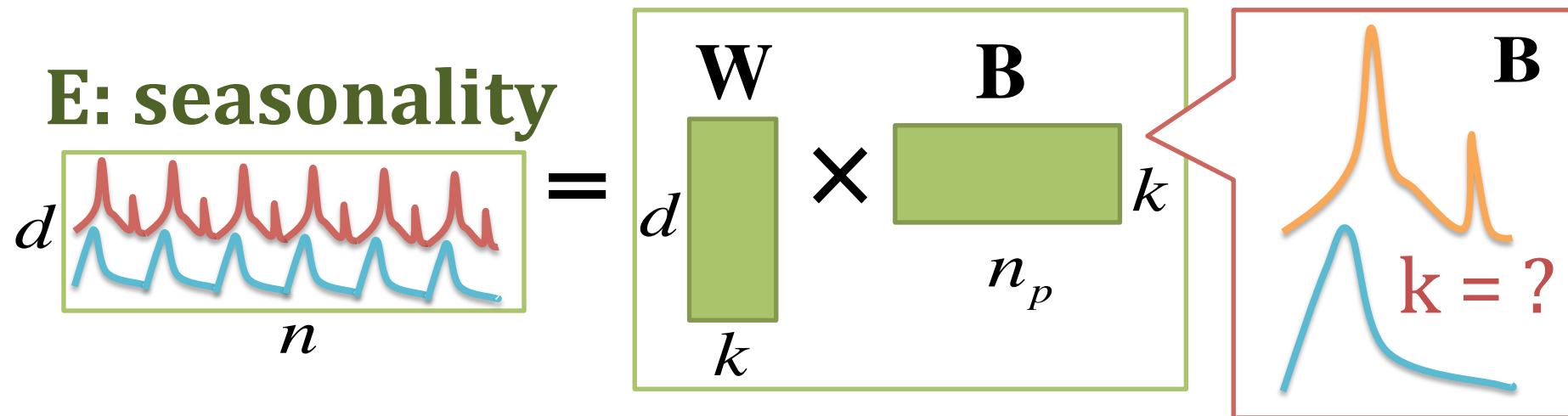
Idea (1): Seasonal component analysis



Idea(1-b) Automatic component analysis

Find optimal number k ($1 \leq k \leq d$)

d : dimension



opt $k = ?$



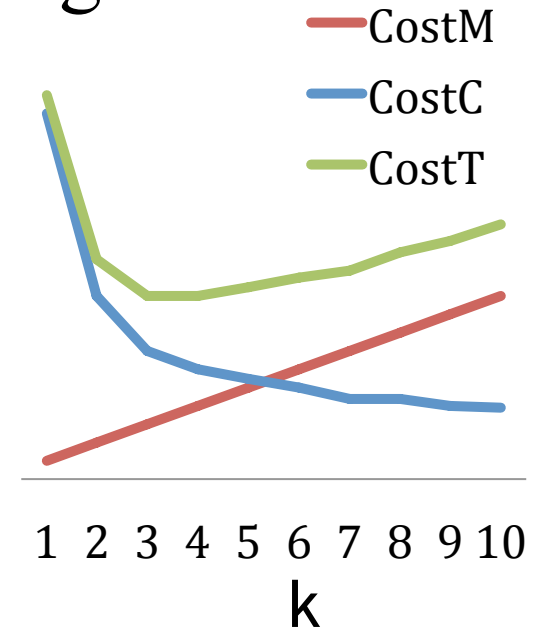
Idea (1): Seasonal component analysis



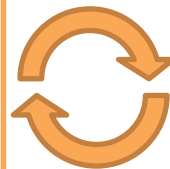
Idea(1-b) MDL -> Minimize encoding cost!

$$\min \left(\boxed{\text{Cost}_M(S)} + \boxed{\text{Cost}_c(X|S)} \right)$$

Model cost
Coding cost



Good
compression



Good
description



Idea (1): Seasonal component analysis



Idea(1-b) MDL -> Minimize encoding cost!

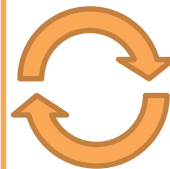
— CostM

— CostC

$$Cost_T(X; \mathcal{S}) = \log^*(d) + \log^*(n) + Cost_M(\mathbf{p}, \mathbf{r}, \mathbf{K}) \\ + Cost_M(\mathbf{A}) + Cost_M(k, \mathbf{W}, \mathbf{B}) + Cost_C(X|\mathcal{S})$$

$$k_{opt} = \arg \min_k Cost_T(X; \mathcal{S})$$

Good
compression



Good
description



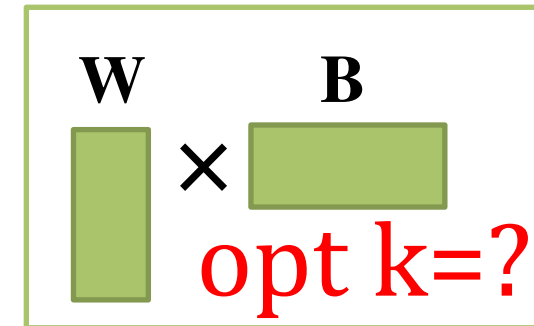
Idea (1): Seasonal component analysis



Idea(1-b) Automatic component analysis

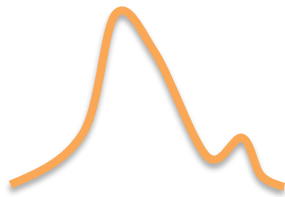
Find optimal number k ($1 \leq k \leq d$)

d : dimension



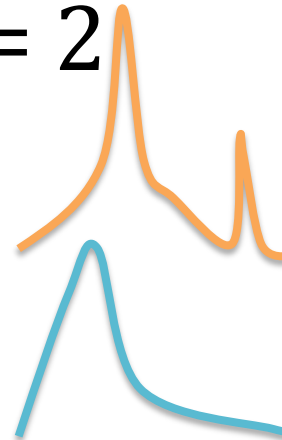
B

$k = 1$



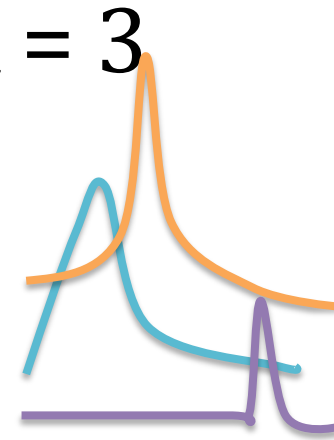
Cost(1) = \$\$

$k = 2$



Cost(2) = \$

$k = 3$



Cost(3) = \$\$\$



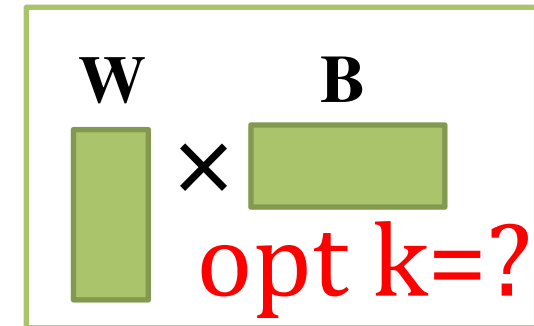
Idea (1): Seasonal component analysis



Idea(1-b) Automatic component analysis

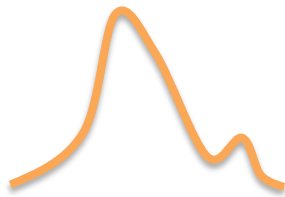
Find optimal number k ($1 \leq k \leq d$)

Optimal k



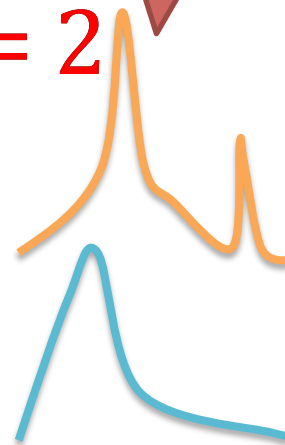
B

$k = 1$



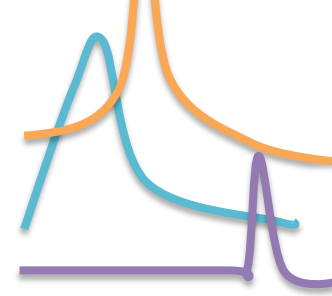
Cost(1) = \$\$

$k = 2$



Cost(2) = \$

$k = 3$

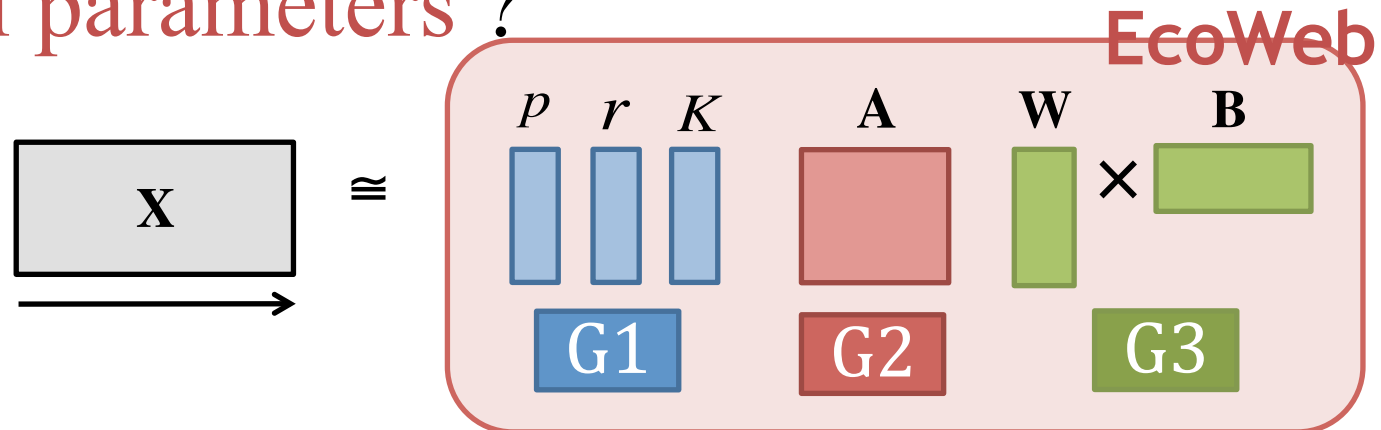


Cost(3) = \$\$\$



Idea (2): EcoWeb-Fit

Q2. How can we efficiently estimate model parameters ?



Idea (2): Multi-step fitting

a. **StepFit** (sub)

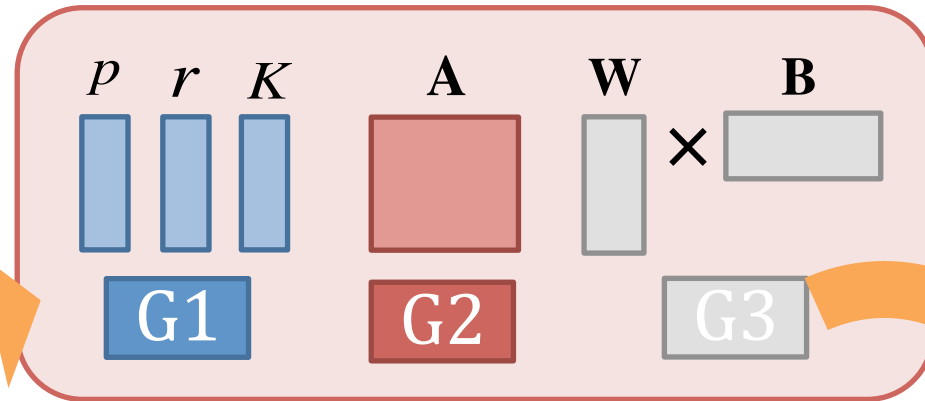
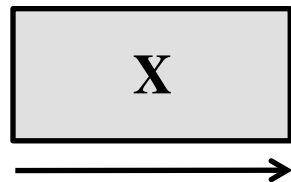
b. **EcoWeb-Fit** (full)



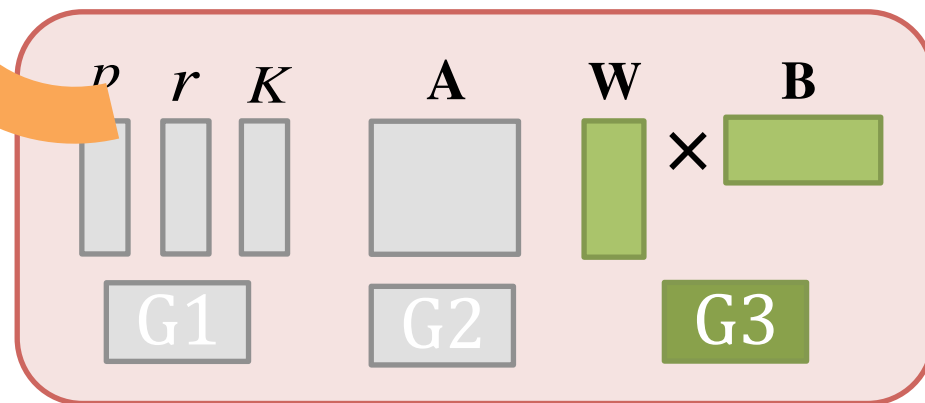
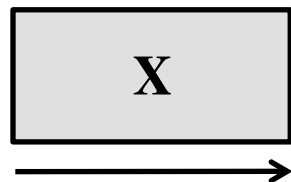
Idea (2): EcoWeb-Fit

(2-a). StepFit: Update parameters *alternately*

Step A



Step B



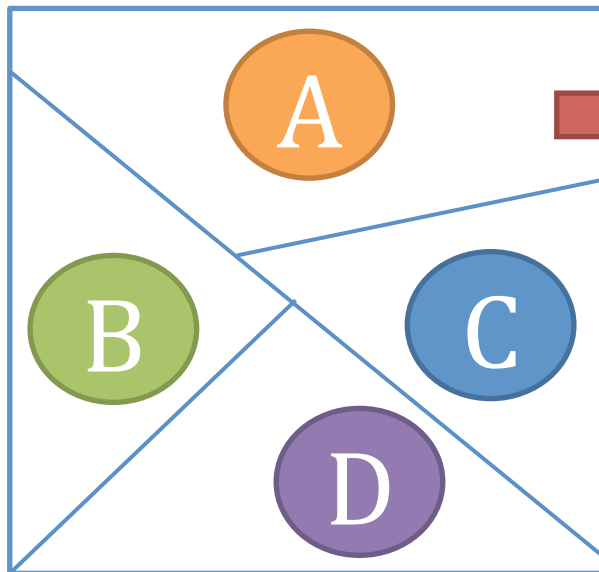


Idea (2): EcoWeb-Fit

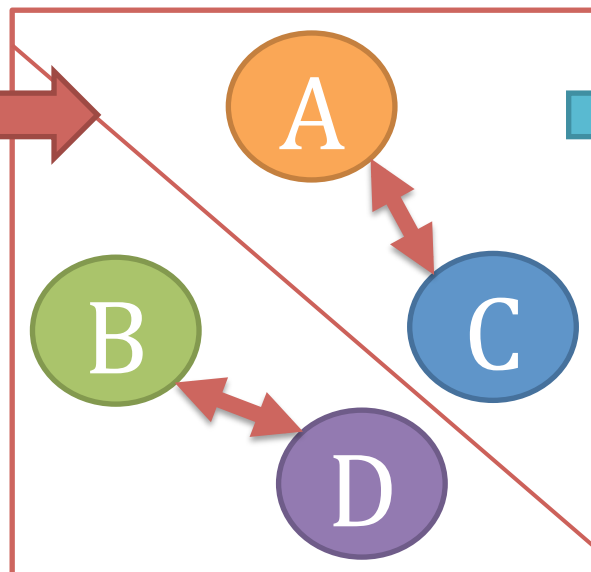
(2-b). EcoWeb-Fit: full algorithm

e.g., 4 keywords: A B C D

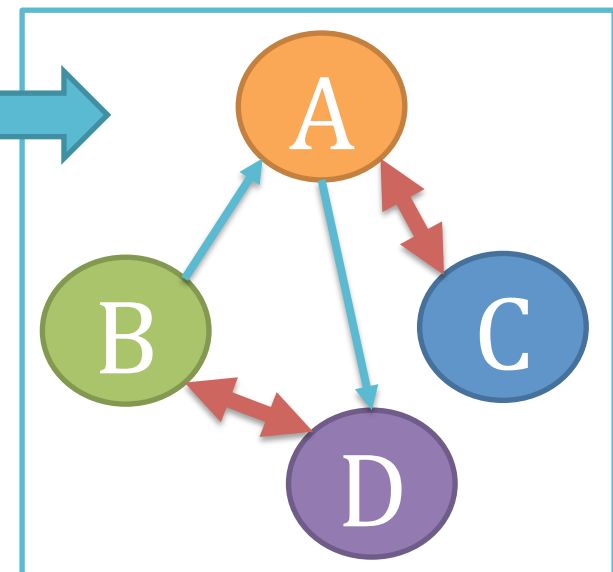
1. Individual-Fit



2. Pair-Fit



3. Full-Fit



EcoWeb-Fit updates parameters, separately



Experiments

We answer the following questions...

Q1. Effectiveness

How successful is it in spotting patterns?

Q2. Accuracy

How well does it match the data?

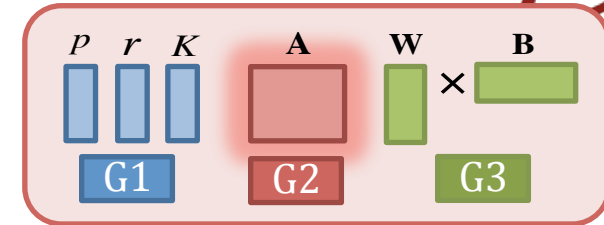
Q3. Scalability

How does it scale in terms of computational time?



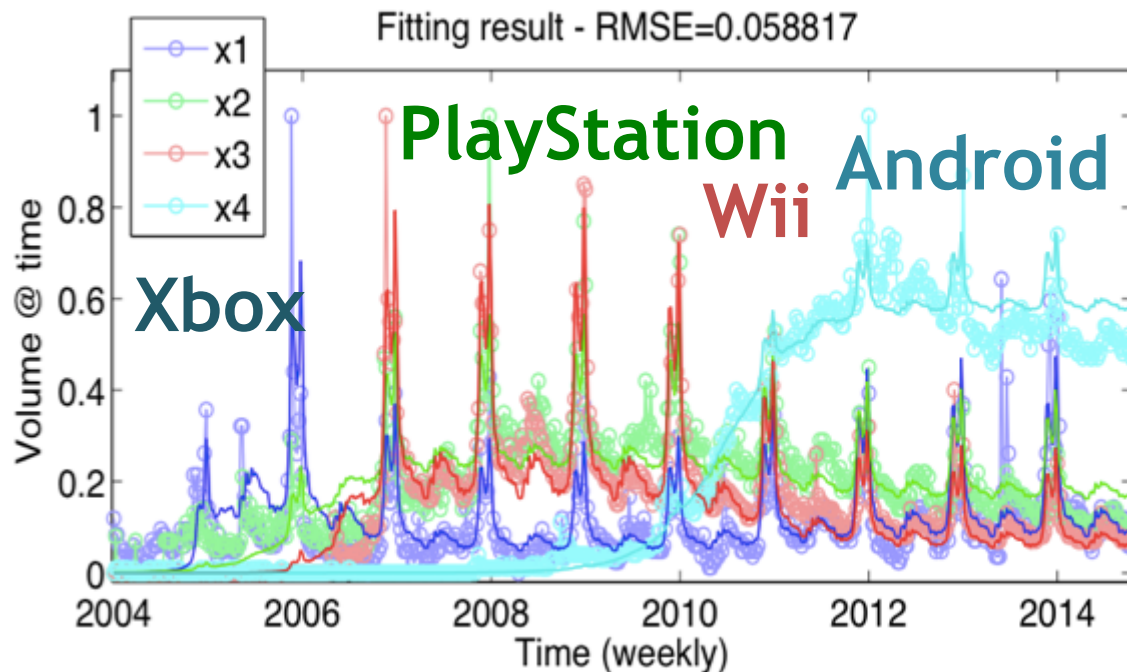
Q1. Effectiveness

(#1) Video games



Interactions

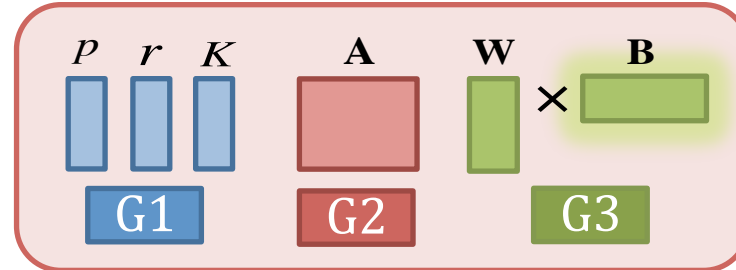
between keywords



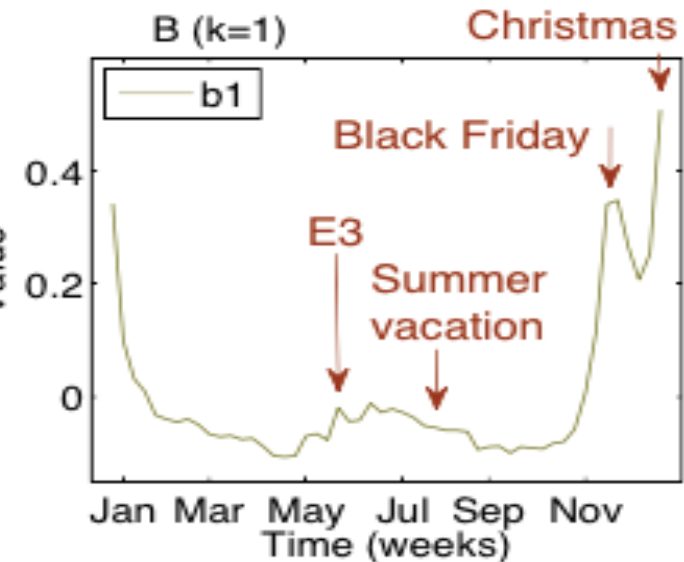
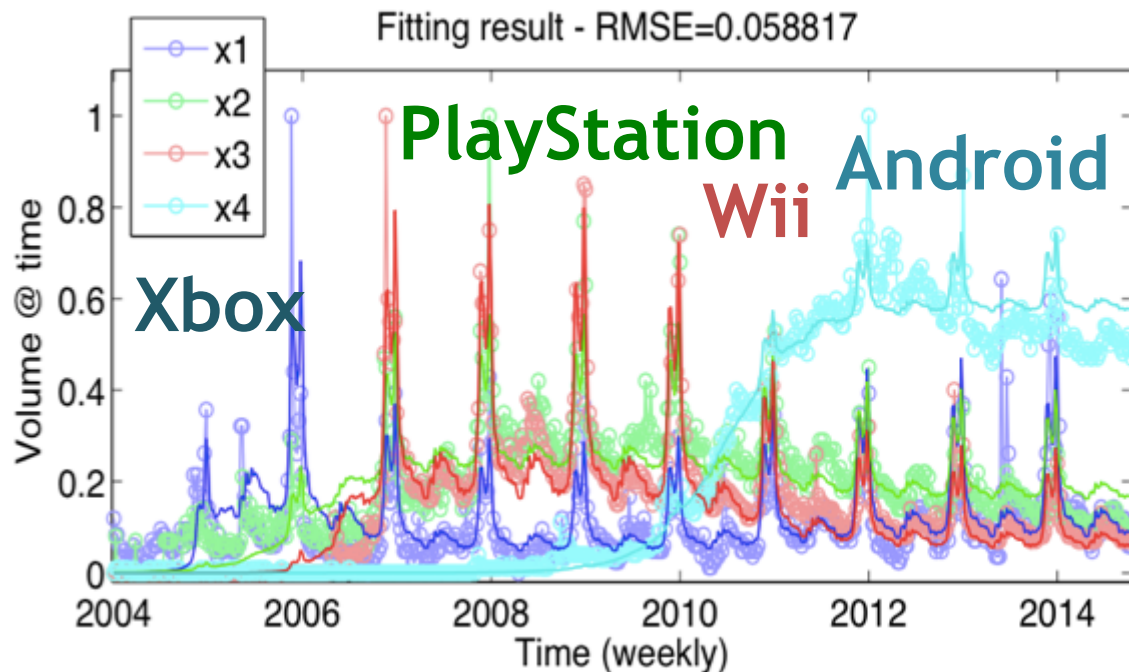


Q1. Effectiveness

(#1) Video games



Seasonality



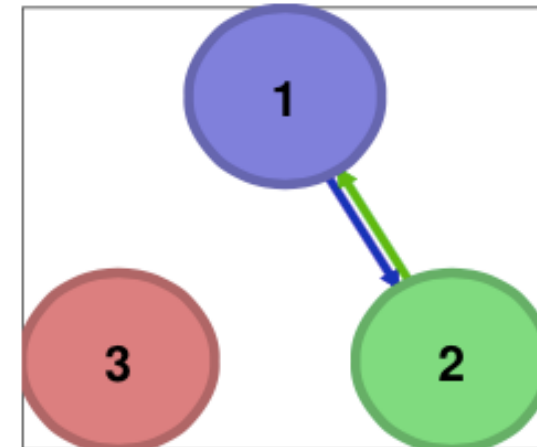


Q1. Effectiveness

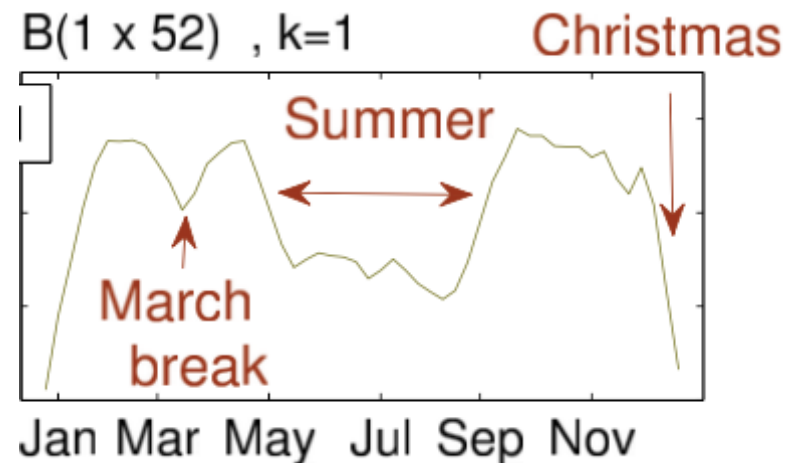
(#2) Programming language

C , **R** , **MATLAB**

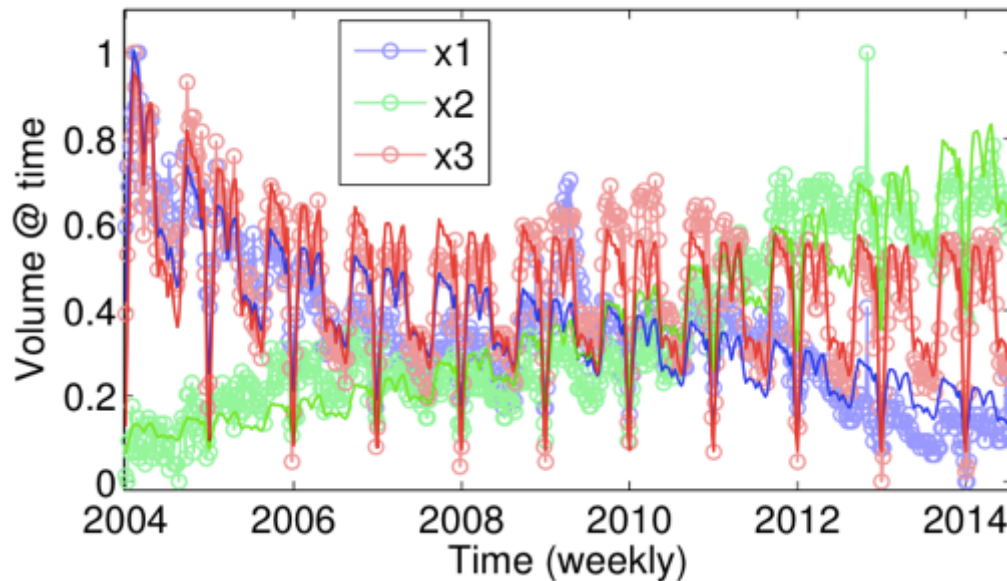
Interactions



Seasonality



Fitting result - RMSE=0.076417



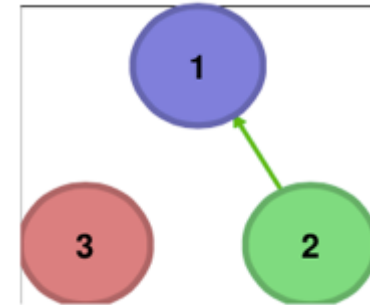


Q1. Effectiveness

(#3) Social media

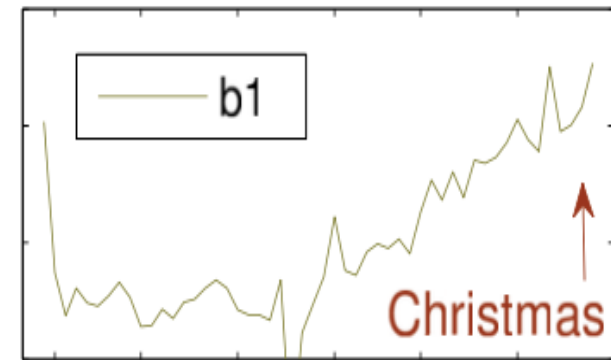
Tumblr , **Facebook** , **LinkedIn**

Interactions

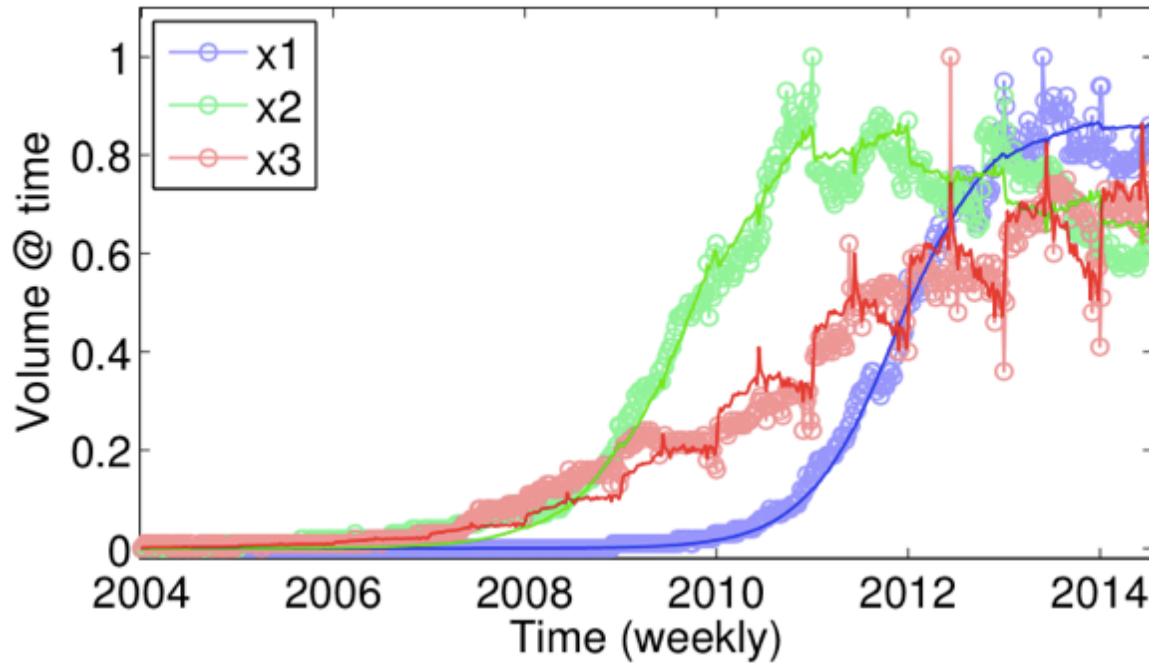


Seasonality

$B(1 \times 52)$, $k=1$



Fitting result - RMSE=0.039536





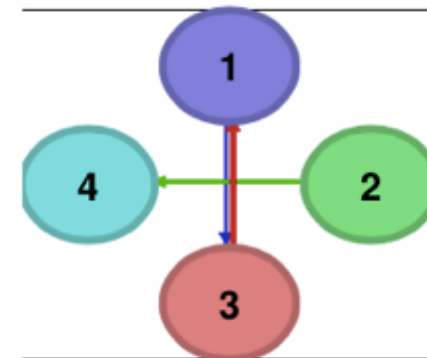
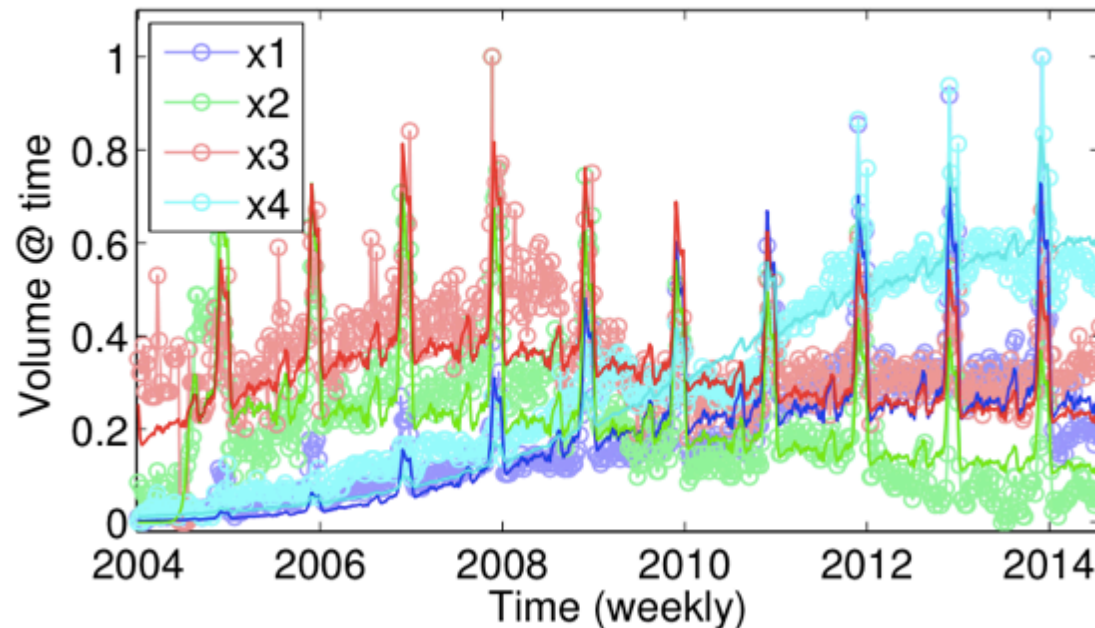
Q1. Effectiveness



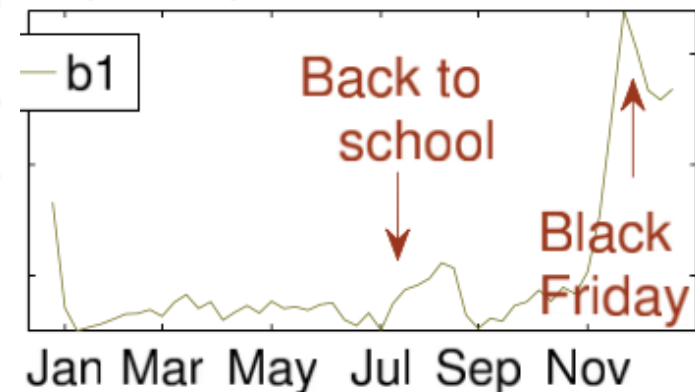
(#4) Apparel companies

Kohls , **JCPenny** , **Nordstrom** , **Forever21**

Fitting result - RMSE=0.074104



$B(1 \times 52)$, $k=1$



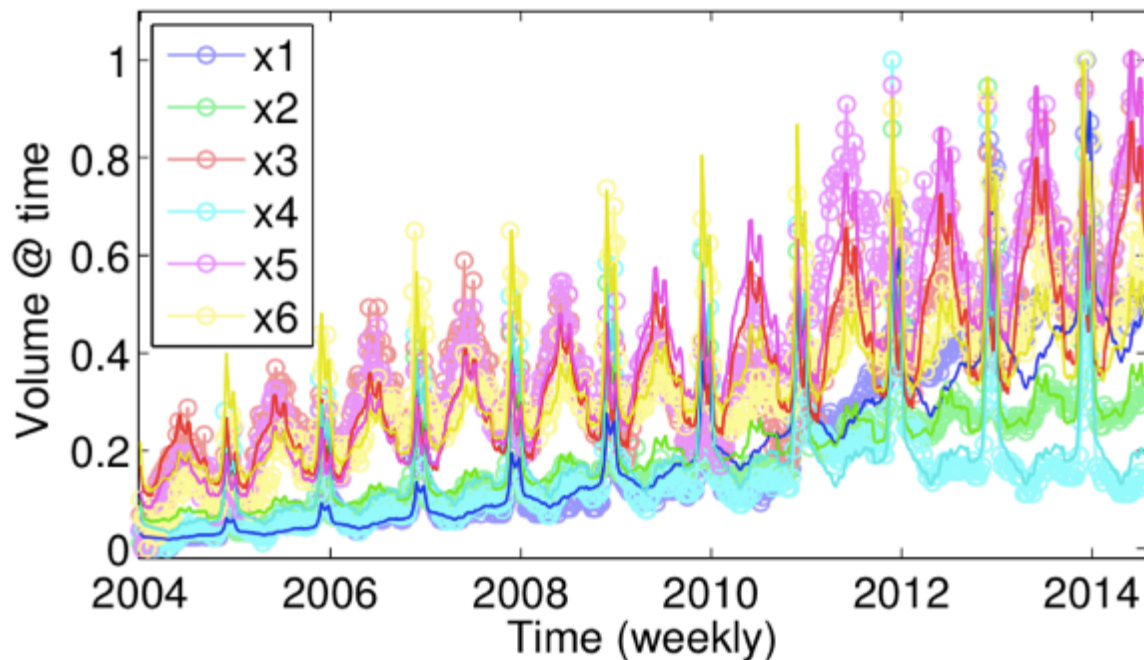


Q1. Effectiveness

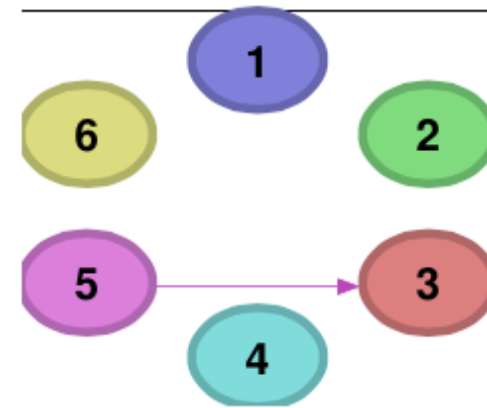
(#5) Retail companies

Amazon , Walmart , Home Depot ,
BestBuy , Lowes , Costco

Fitting result - RMSE=0.065173



Interaction



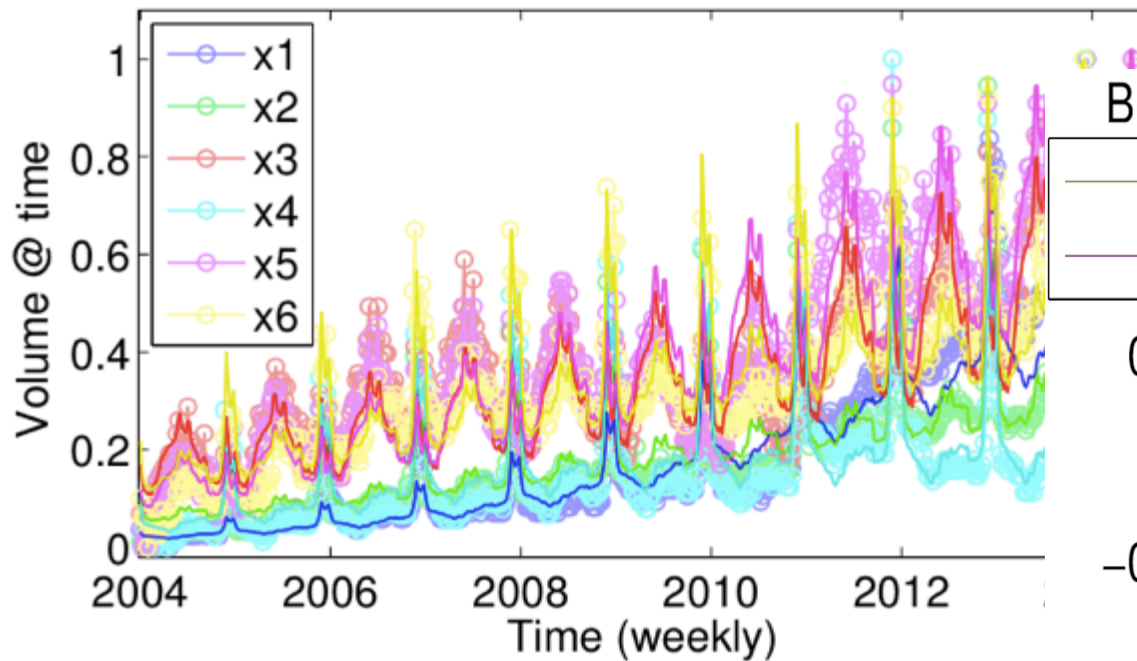


Q1. Effectiveness

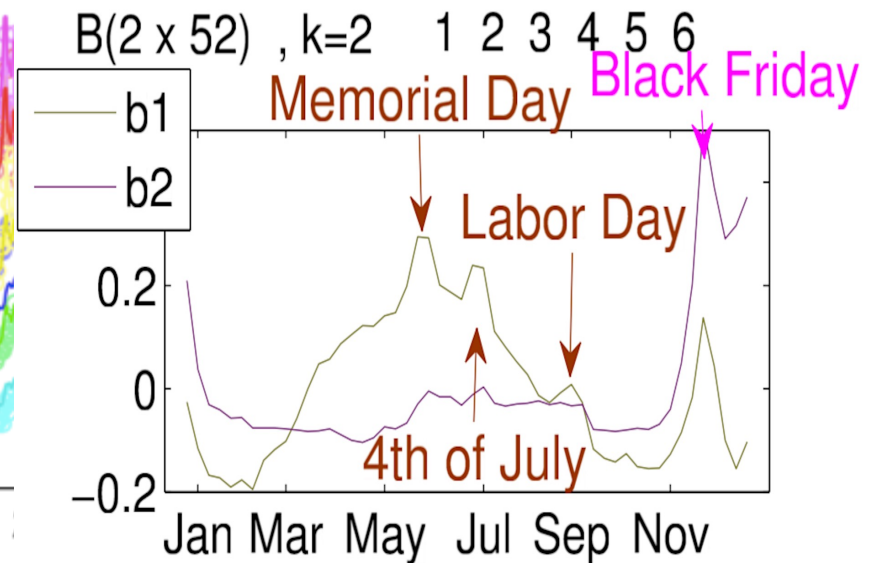
(#5) Retail companies

Amazon , Walmart , Home Depot ,
BestBuy , Lowes , Costco

Fitting result - RMSE=0.065173



Seasonality

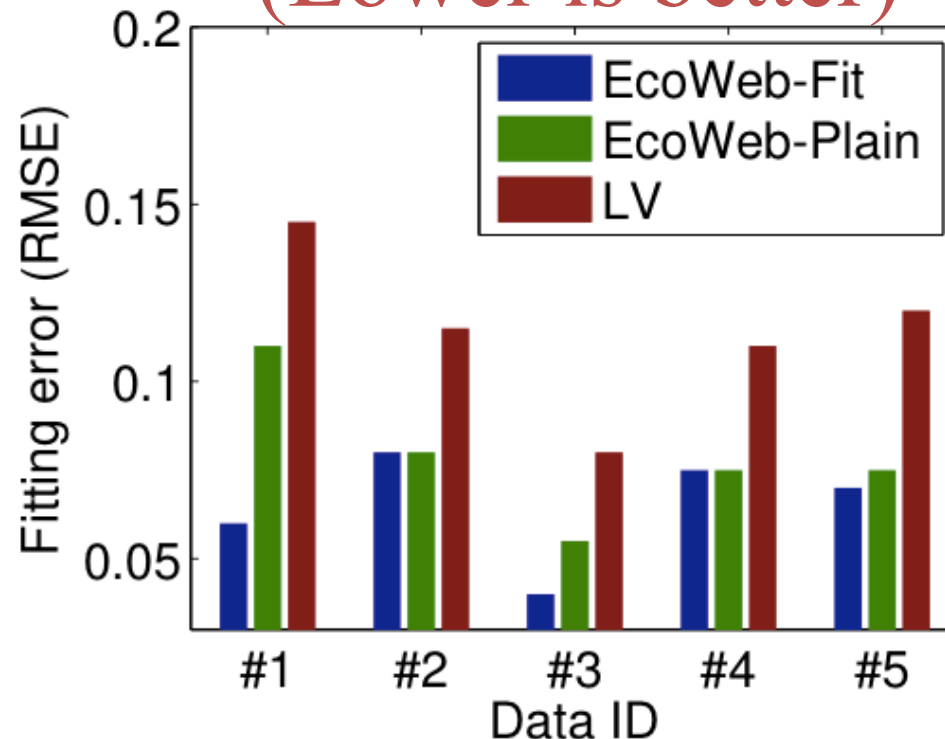




Q2. Accuracy

RMSE between original and fitted volume

(Lower is better)



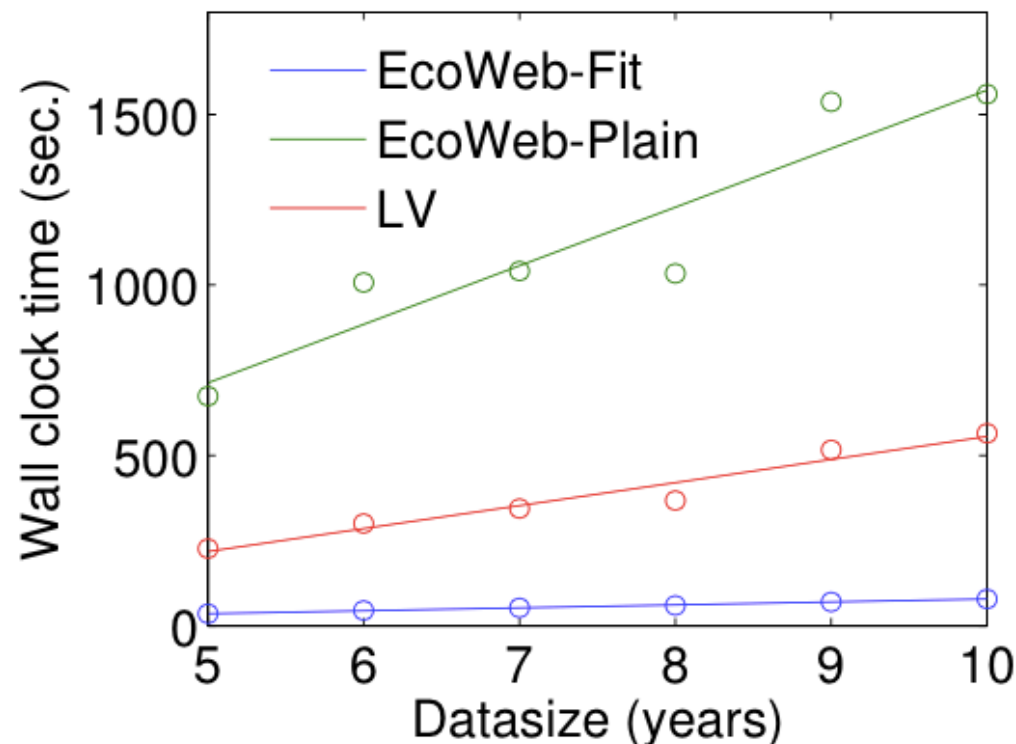
EcoWeb consistently wins!



Q3. Scalability

Wall clock time vs. dataset size (years)

EcoWeb-Fit scales linearly, i.e., $O(n)$



7x faster than **LV**, 20x faster than **EcoWeb-Plain**

EcoWeb at work - forecasting

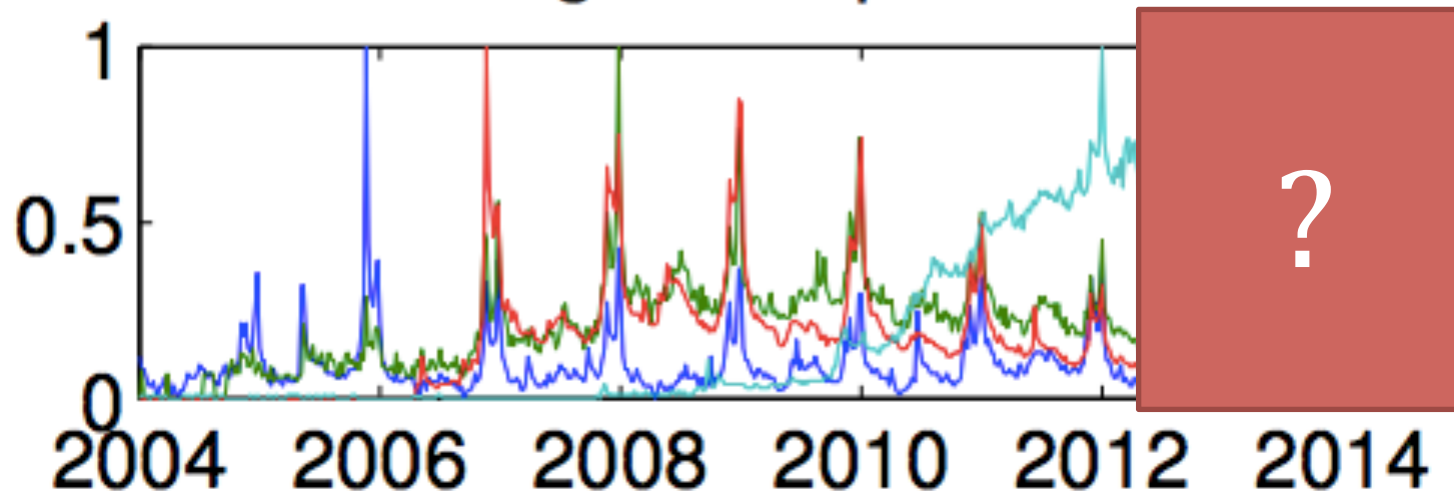


Forecasting future activities

Train:
2/3 sequences

Forecast:
1/3 following years

Original sequences



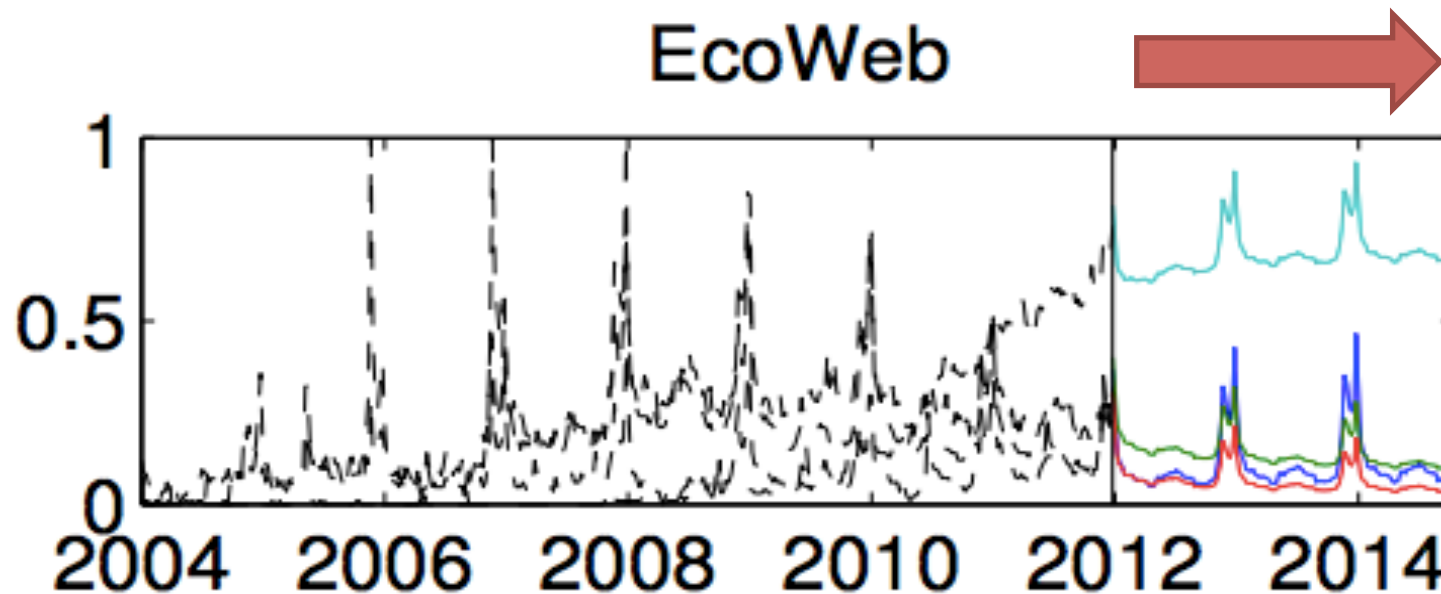
EcoWeb at work - forecasting



Forecasting future activities

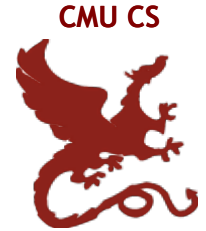
Train:
2/3 sequences

Forecast:
1/3 following years

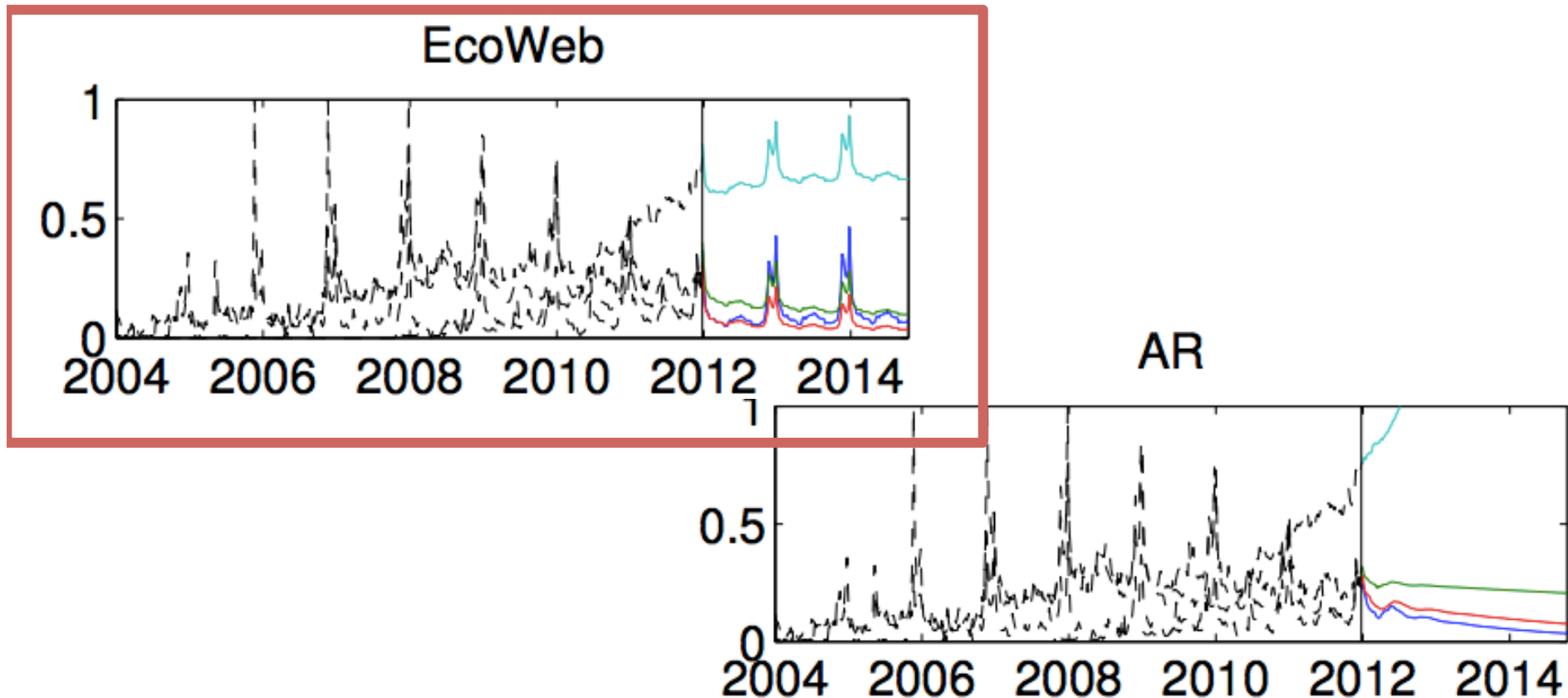


EcoWeb can capture future patterns

EcoWeb at work - forecasting



Forecasting future activities

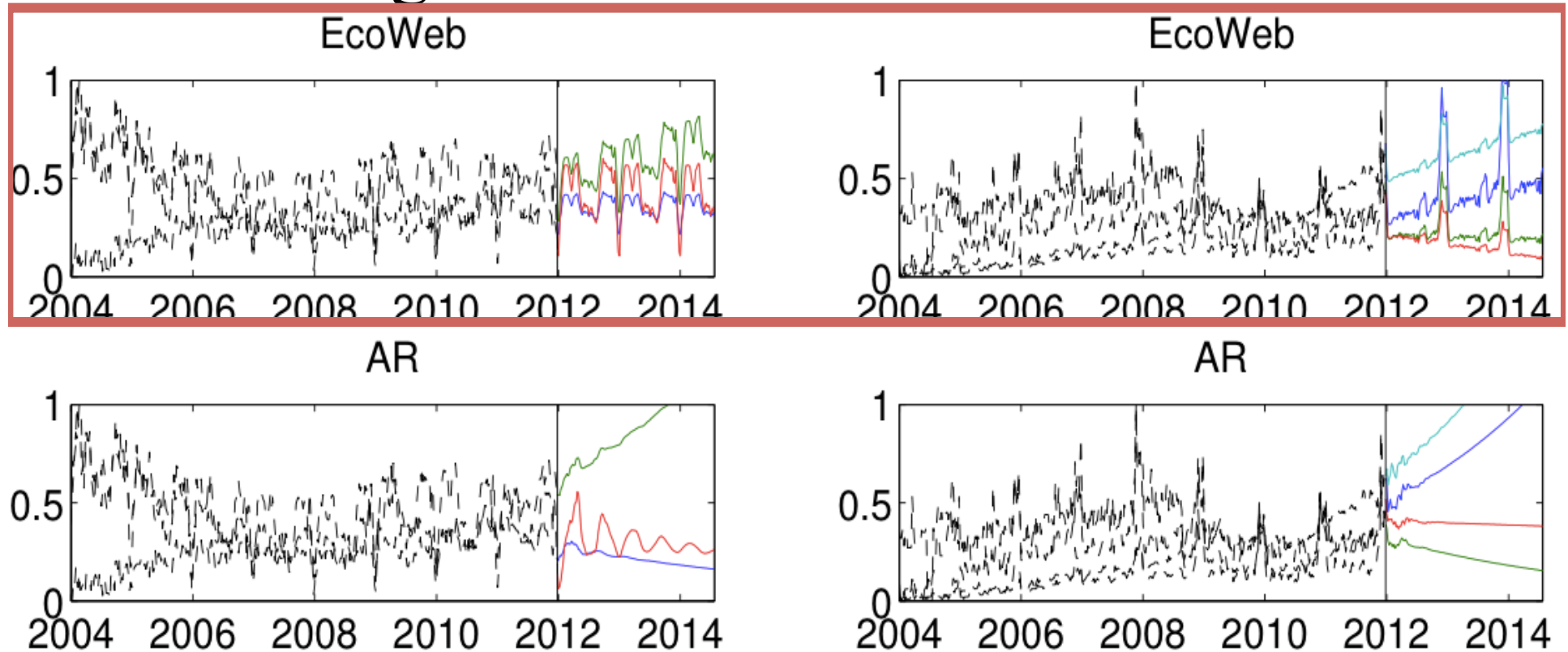


EcoWeb can capture future patterns!

EcoWeb at work - forecasting



Forecasting future activities



(b) Programming languages (#2)

(c) Apparel companies (#4)

EcoWeb can capture future patterns!



Part 2

Roadmap



Problem

- ✓ Why: “non-linear” modeling

Fundamentals

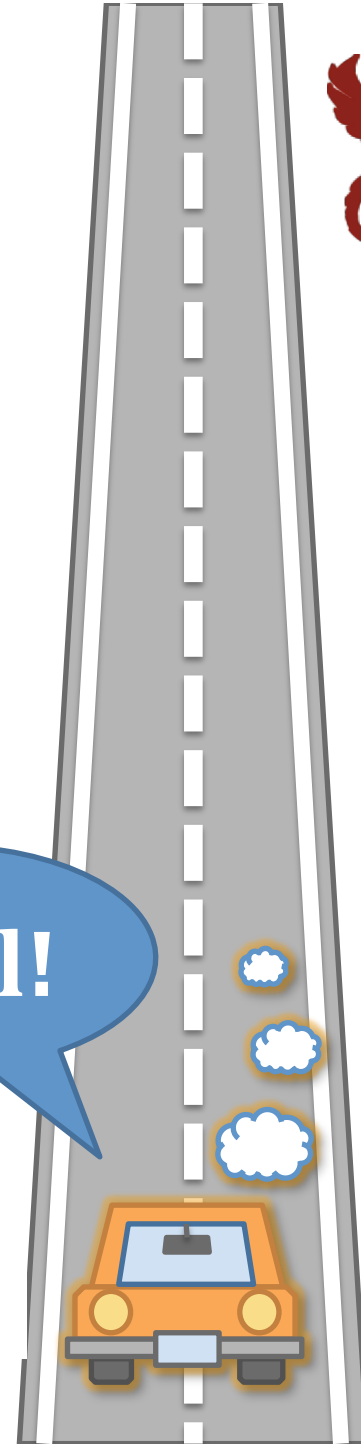
- ✓ Non-linear (grey-box) models

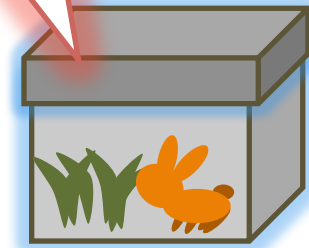
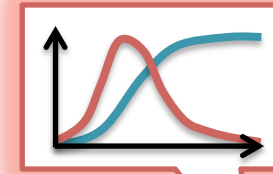
Applications

- ✓ Epidemics
- ✓ Information diffusion
- ✓ Online competition



vs.





✓ Why: “non-linear” modeling

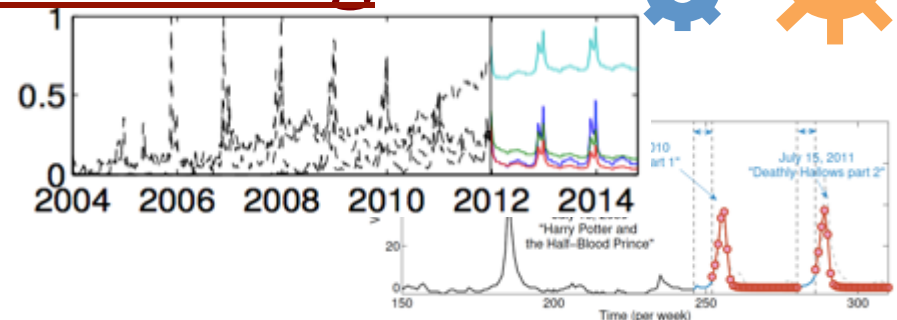
- Black box: lag plots (k-NN search)
- Grey-box: given a model

✓ Fundamentals: popular non-linear models

- Logistic function, Lotka-Volterra, Competition, ...
- Epidemics (SI, SIR, SEIR, etc.), ...

✓ Applications: non-linear mining

- Epidemics
- Information diffusion
- Online competition





References (1)



Fundamentals

- Non-linear forecasting
 - D. Chakrabarti and C. Faloutsos *F4: Large-Scale Automated Forecasting using Fractals* CIKM 2002, Washington DC, Nov. 2002.
 - Sauer, T. (1994). *Time series prediction using delay coordinate embedding*. (in book by Weigend and Gershenfeld, below) Addison-Wesley.
 - Takens, F. (1981). *Detecting strange attractors in fluid turbulence*. Dynamical Systems and Turbulence. Berlin: Springer-Verlag.
- Non-linear equations and modeling
 - F. Brauer and C. Castillo-Chavez. *Mathematical models in population biology and epidemiology*, volume 40. Springer Verlag, New York, 2001.
 - R. M. Anderson and R. M. May. *Infectious Diseases of Humans Dynamics and Control*. Oxford University Press, 1992.
 - F. M. Bass. A new product growth for model consumer durables. *Management Science*, 15(5):215–227, 1969.
 - D. Easley and J. Kleinberg. *Networks, Crowds, and Markets: Reasoning About a Highly Connected World*. Cambridge University Press, 2010.
 - R. M. Anderson and R. M. May. *Infectious Diseases of Humans*. Oxford University Press, 1991.
 - R. M. May. Qualitative stability in model ecosystems. *Ecology*, 54(3):638–641, 1973.
 - M. Nowak. *Evolutionary Dynamics*. Harvard University Press, 2006.
 - Schuster, H. G. and Wagner, P. A model for neuronal oscillations. *Biol. Cybern.*, 1990.
- Others
 - A. G. Hawkes and D. Oakes. A cluster representation of a self-exciting process. *J. Appl. Prob.*, 11:493–503, 1974.



References (2)

Applications

- Epidemics

- Rohani, P., Earn, D. J. D., Finkenstadt, B. F. & Grenfell, B. T. Population dynamic interference among childhood diseases. *Proc. R. Soc. Lond. B* 265, 2033–2041 (1998).
- Rohani, P., Green, C.J., Mantilla-Beniers, N.B. & Grenfell, B.T. Ecological Interference Among Fatal Infections. *Nature* 422: 885–888 (2003).
- L. Stone, R. Olinky, and A. Huppert. Seasonal dynamics of recurrent epidemics. *Nature*, 446:533–536, March 2007.
- Y. Matsubara, Y. Sakurai, W. G. van Panhuis, and C. Faloutsos. FUNNEL: automatic mining of spatially coevolving epidemics. In *KDD*, pages 105–114, 2014.

- Information diffusion

- J. Leskovec, L. Backstrom, and J. M. Kleinberg. Meme-tracking and the dynamics of the news cycle. In *KDD*, pages 497–506, 2009.
- J. Yang and J. Leskovec. Patterns of temporal variation in online media. In *WSDM*, pages 177–186, 2011.
- J. Yang and J. Leskovec. Modeling information diffusion in implicit networks. In *ICDM*, pages 599–608, 2010.
- R. Crane and D. Sornette. Robust dynamic classes revealed by measuring the response function of a social system. In *PNAS*, 2008.
- F. Figueiredo, J. M. Almeida, Y. Matsubara, B. Ribeiro, and C. Faloutsos. Revisit behavior in social media: The phoenix-r model and discoveries. In *PKDD*, pages 386–401, 2014.
- Y. Matsubara, Y. Sakurai, B. A. Prakash, L. Li, and
- C. Faloutsos. Rise and fall patterns of information diffusion:
- model and implications. In *KDD*, pages 6–14, 2012.

- Online activities and competition

- B. A. Prakash, A. Beutel, R. Rosenfeld, and C. Faloutsos. Winner takes all: competing viruses or ideas on fair-play networks. In *WWW*, pages 1037–1046, 2012.
- A. Beutel, B. A. Prakash, R. Rosenfeld, and C. Faloutsos. Interacting viruses in networks: can both survive? In *KDD*, pages 426–434, 2012.
- Y. Matsubara, Y. Sakurai, and C. Faloutsos. The web as a jungle: Non-linear dynamical systems for co-evolving online activities. In *WWW*, 2015.

Part 2



Non-linear mining and forecasting

Yasushi Sakurai (Kumamoto University)

Yasuko Matsubara (Kumamoto University)

Christos Faloutsos (Carnegie Mellon University)