# FUNNEL: Automatic Mining of Spatially Coevolving Epidemics

Yasuko Matsubara, Yasushi Sakurai (Kumamoto University)
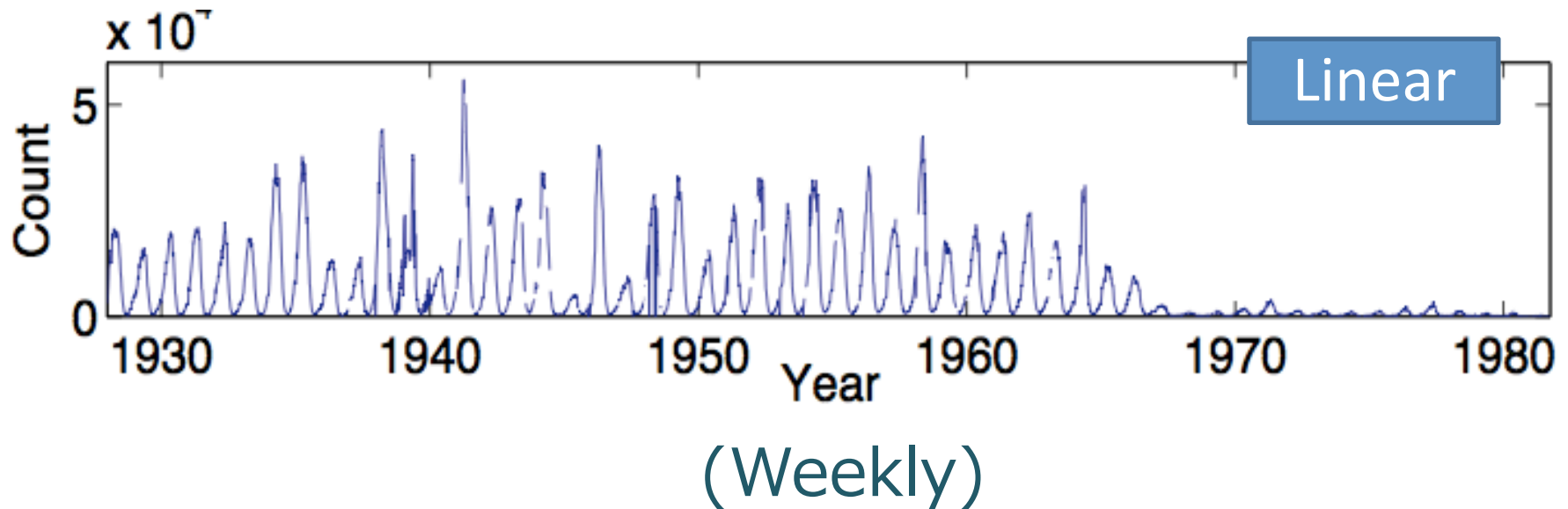
Willem G. van Panhuis (University of Pittsburgh)
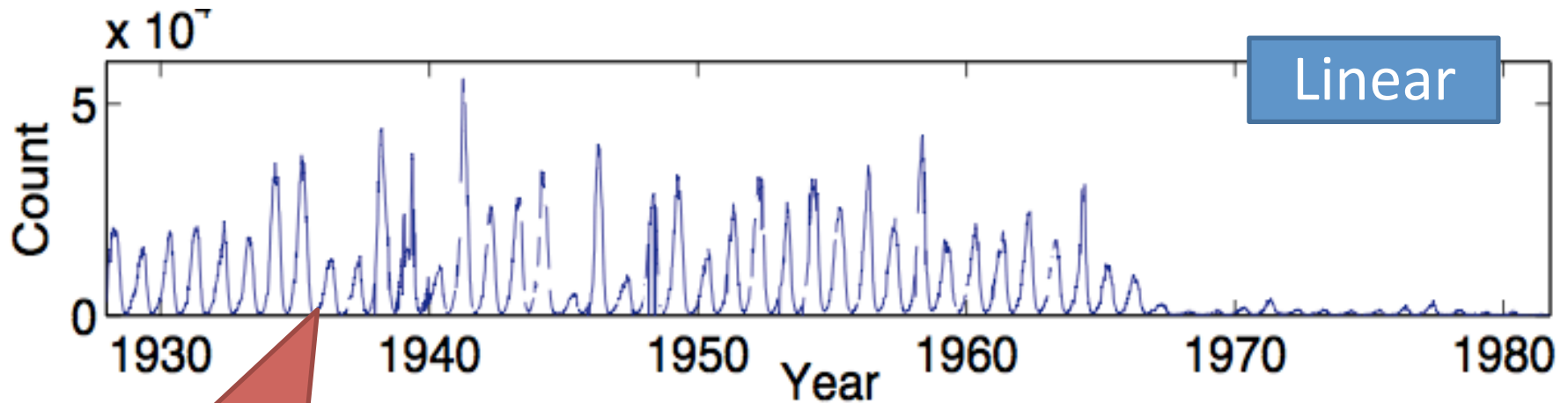
Christos Faloutsos (CMU)

# Motivation

Given: Large set of epidemiological data

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



Linear

(Weekly)

# Motivation

Given: Large set of epidemiological data

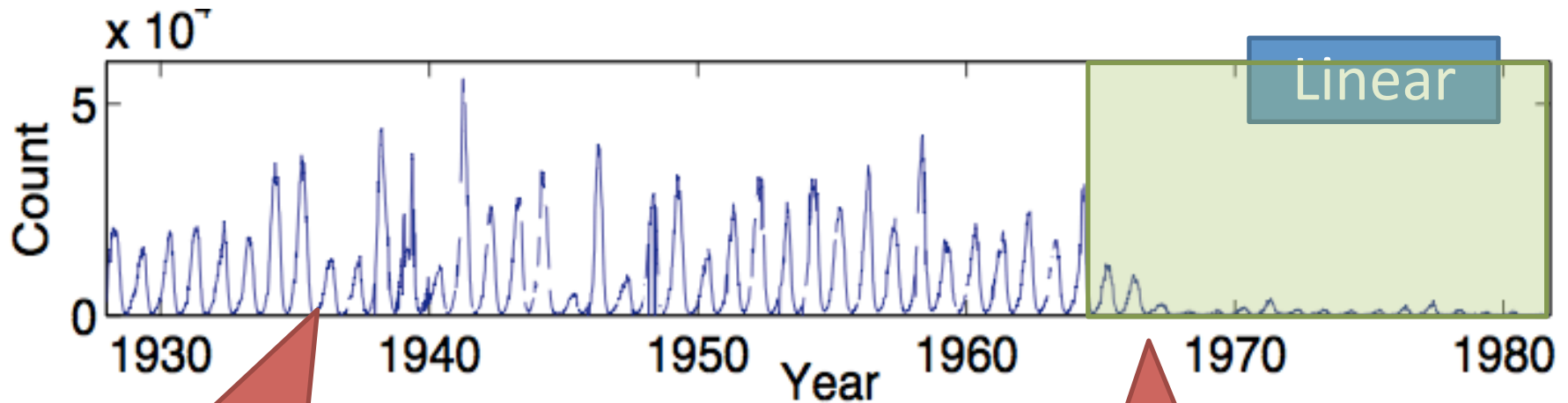e.g., Measles cases in the U.S.



Linear

Yearly periodicity

(Weekly)

# Motivation

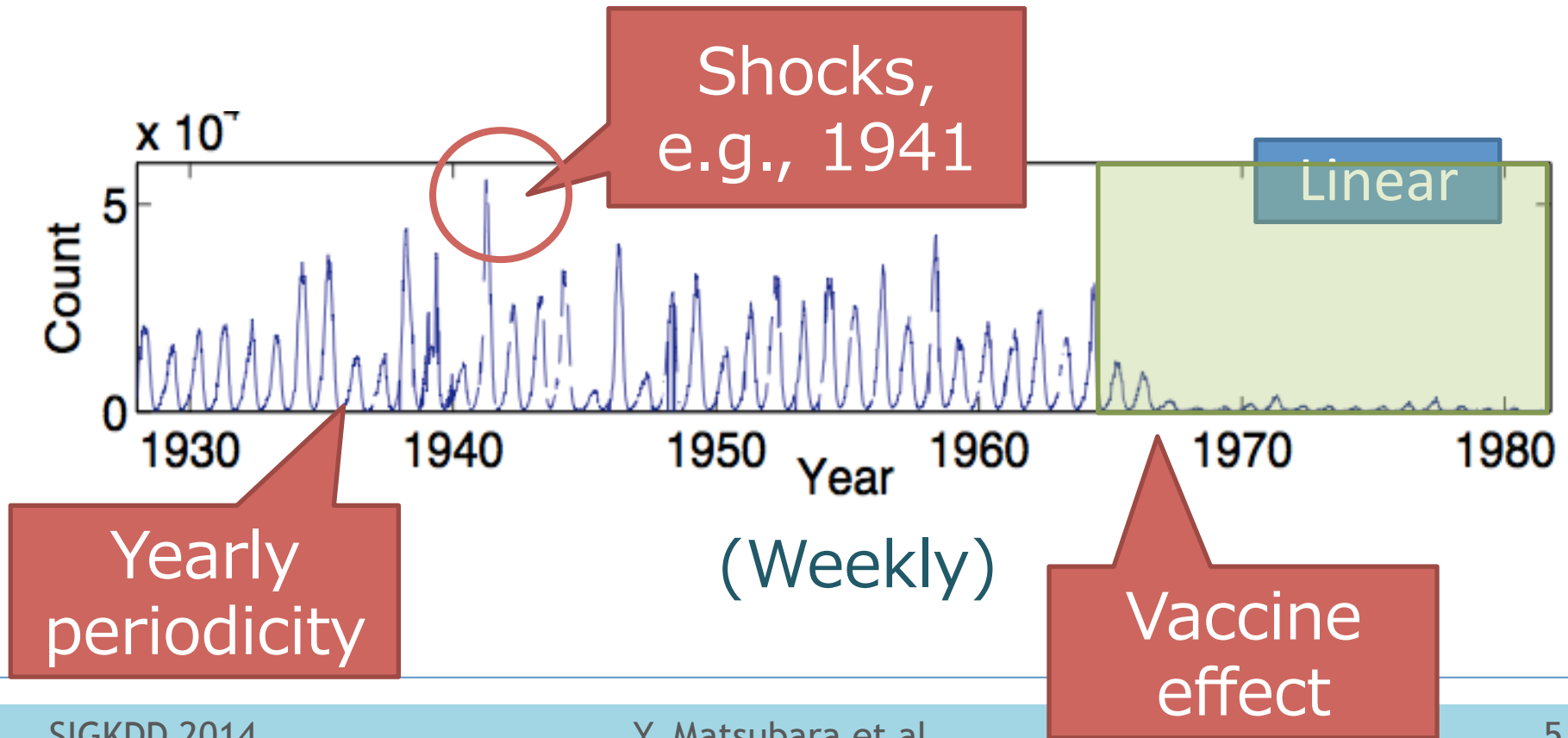Given: Large set of epidemiological data

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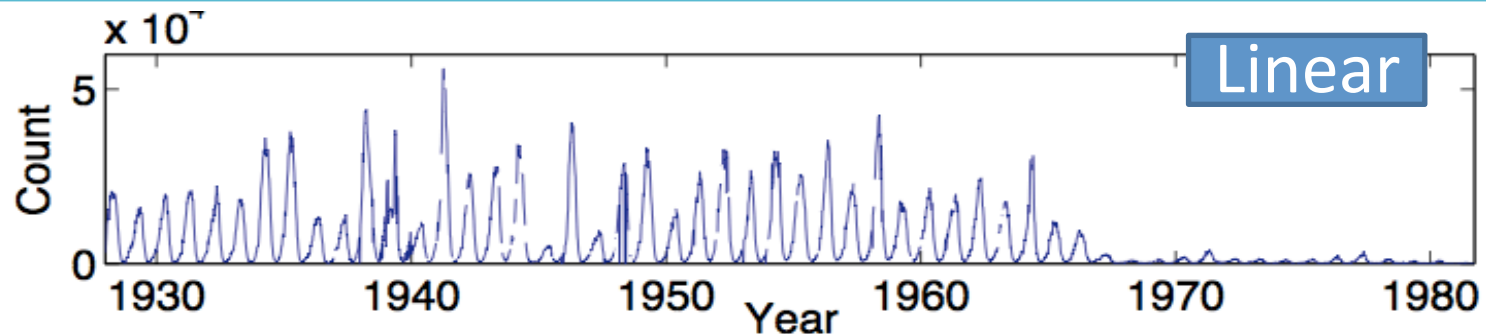


(Weekly)

Yearly periodicity

Vaccine effect

Linear

# Motivation

Given: Large set of epidemiological data

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



Shocks, e.g., 1941

Linear

Yearly periodicity

(Weekly)

Vaccine effect

# Motivation

Given: Large set of epidemiological data

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

Goal: summarize all the epidemic time-series, **"fully-automatically"**
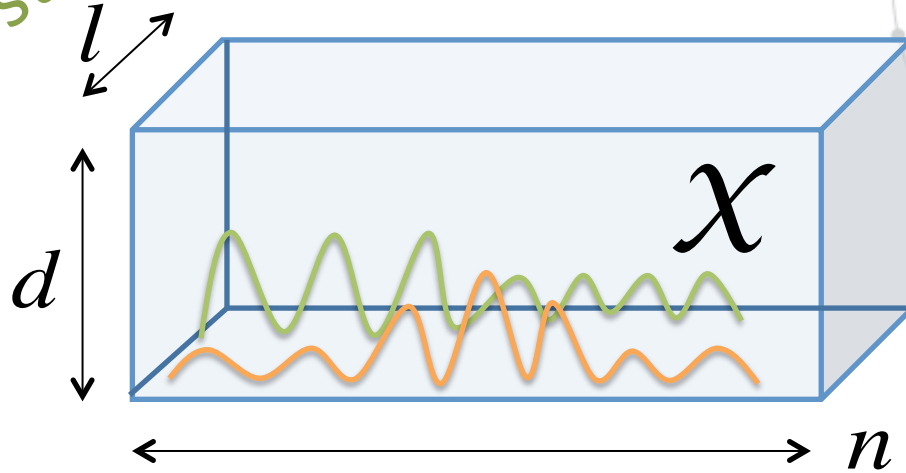
# Data description

Project Tycho: infectious diseases in the U.S.



50 states

$l$

$d$

56 diseases

$X$

$n$

1888

Time (weekly)   (> **125** years)

PROJECT TYCHO

DATA FOR HEALTH

# Data description

Project Tycho: infectious diseases in the U.S.



**50** states

**56** diseases

$l$

$d$

$x$

$x$

$n$

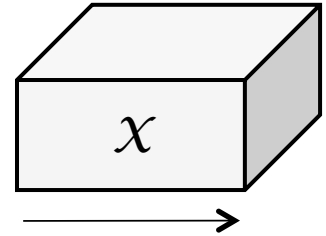1888    Time (weekly)    (> **125** years)

Element x : # of cases

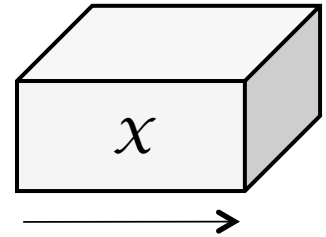e.g., 'measles', 'NY', 'April 1-7, 1931', '4000'

# Problem definition

Given:

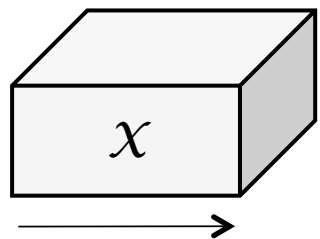Tensor $\mathcal{X}$ (disease x state x time)

# Problem definition

**Given**:

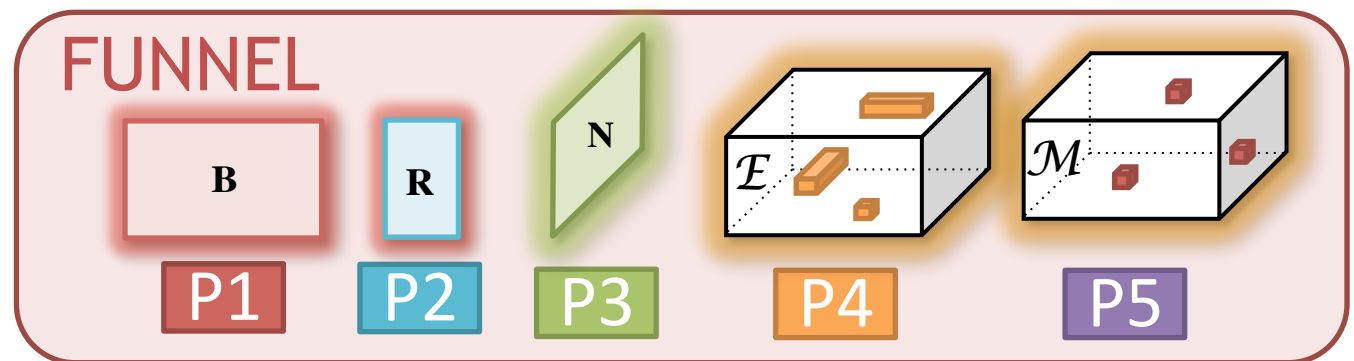Tensor $\mathcal{X}$ (disease x state x time)

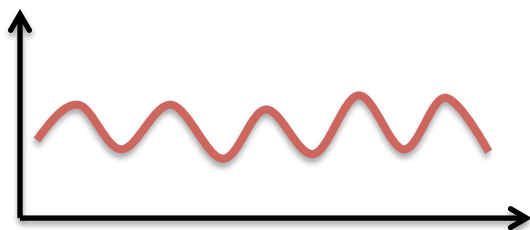**Find**:

Compact description of $\mathcal{X}$, *"automatically"*

FUNNEL

$\mathcal{X}$ =  B  R  N  $\mathcal{E}$  $\mathcal{M}$

P1  P2  P3  P4  P5

# Problem definition

Gi...

Te... ...state

Fi...

Compact des...tion of $\mathcal{X}$, "...tomatically"

**Seasonality**

**Discontinuities**

FUNNEL

$\mathcal{X}$ = | **B** | **R** | **N** | $\mathcal{E}$ | $\mathcal{M}$ |

P1  P2  P3  P4  P5

# Problem definition

Given:

Tensor $\mathcal{X}$

Find:

Compac... _...ically"_

**NO magic numbers !**

**Parameter-free!**

$$\mathcal{X} \quad = \quad \text{FUNNEL}$$

| | | | | |
|---|---|---|---|---|
| **B** | **R** | **N** | $\mathcal{E}$ | $\mathcal{M}$ |
| P1 | P2 | P3 | P4 | P5 |

# Roadmap

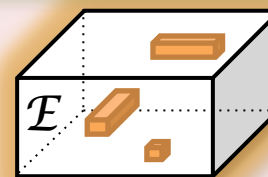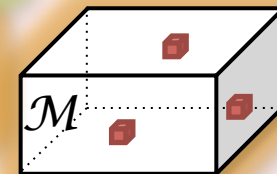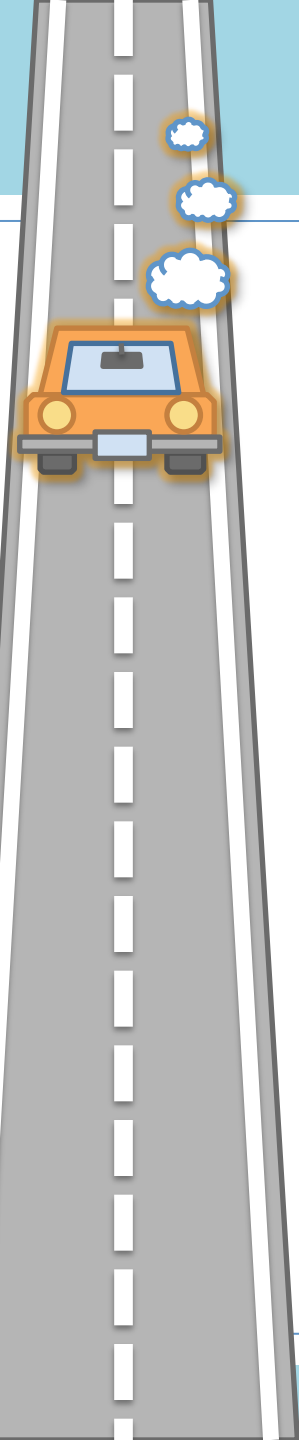✔ Motivation

– Modeling power of FUNNEL

– Overview – main ideas

– Proposed model – idea #1

– Algorithm – idea #2

– Experiments

– Discussion

– Conclusions

# Modeling power of FUNNEL
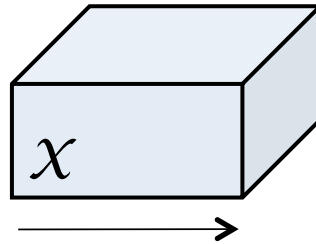
Questions about epidemics
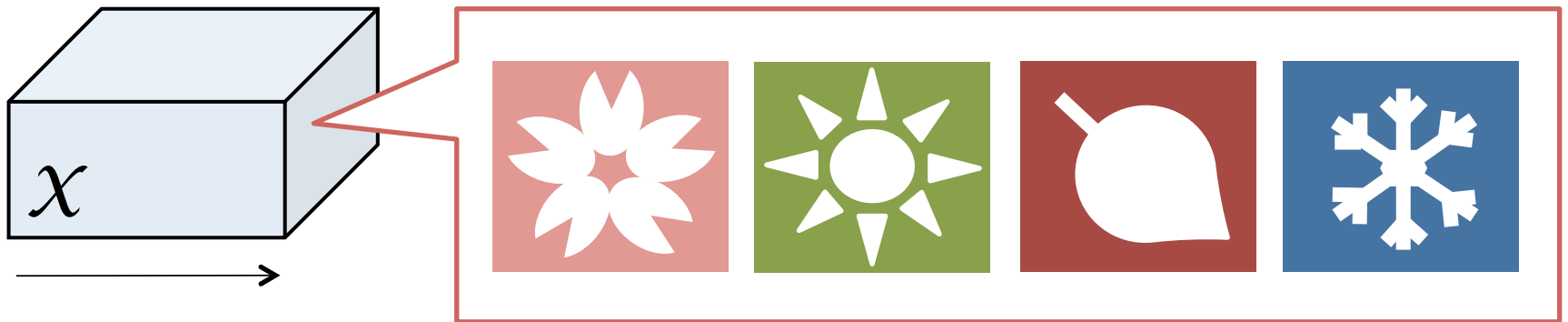


Q1  Q2  Q3  Q4  Q5

$x$

Q1

Are there any periodicities?

If yes, when is the peak season?

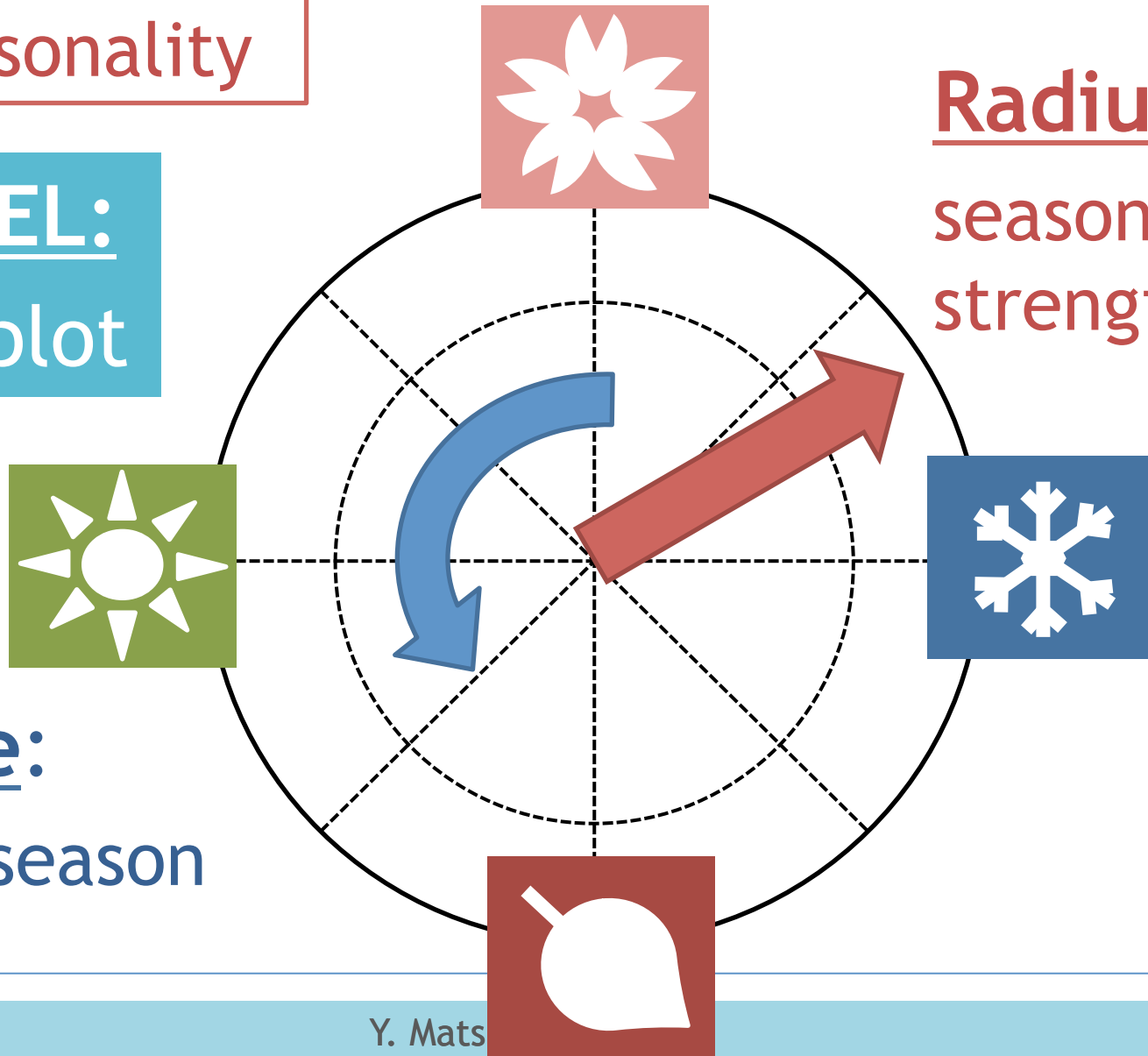# Answers

**P1** Seasonality

**FUNNEL:**
Polar plot

**Radius**:
seasonality strength

**Angle**:
peak season

P1 Seasonality

Questions

?

P1 Seasonality

Questions

Q: Does Influenza have seasonality? If yes, when?

# Answers

**P1** Seasonality

**P1** Seasonality

Influenza in Feb.
Detected by FUNNEL
(strong seasonality)

**Detected!**

# Answers

P1 Seasonality

Questions

?

# Answers

P1  Seasonality

Questions

Q: How about measles ?

# Answers

P1 Seasonality

Questions

?

Y. Mats

# Answers

**P1** Seasonality

Measles (children's) in spring

Detected!

May (5)

Measles 0.4

Rubella 0.3

Chickenpox

February (2)

0.1

Influenza

Smallpox

Rocky mountain spotted fever

(1)

Lymedisease

Streptococo

Gonorrhea

December (12)

August (8)

Typhus fever

Typhoidfever

Cryp

September (9)

November (11)

# Answers

P1 Seasonality

Questions

?

# Answers

P1 Seasonality

Questions

Q: Which disease peaks in summer?

# Answers

**P1** Seasonality



**Questions**

**?**

# Answers

P1 Seasonality



Detected!

Lyme-disease (tick-borne) in summer

May (5) · March (3) · February (2) · (1) · December (12) · November (11) · ptember (9) · September (9)

Measles · 0.4
Rubella · 0.3 · Mumps
Smallpox · 0.2 · Chickenpox
0.1 · Influenza
Rocky m · Streptococc
Gonorrhea
medisease
Typhus fever · Typhoidfever
Cryp

P1 Seasonality

Questions

?

# Answers

P1 Seasonality

Questions

Q: Which disease has no periodicity?

Y. Mats

# Answers

P1 Seasonality

**Questions**

**?**

P1 Seasonality



Detected!

May (5)

Measles 0.4

Rubella 0.3

June (6)

Smallpox 0.2 Ch

0.1

Rocky m...ed fever

Influenza

Streptococo

Gonorrhea

Lymedisease

Gonorrhea
(STD)
no periodicity

...us fever

Typhoidfever

December (12)

Cryp

(9)

November (11)

## Q2

# Can we see any discontinuities?

# Answers

**P2** Disease reduction effect

Measles

Detected!

**1965:** Detected by FUNNEL

**1963:** Vaccine licensure

## Q3

# What's the difference between measles in NY and in FL?

P3  area sensitivity

FUNNEL's guess of susceptibles (measles)



CA

TX

**Detected!**

**NY, PA (more children)**

**FL (fewer children)**

## Q4

## Are there any external shock events, like wars?

P4  external shock events

Funnel can detect external shocks "**fully-automatically**" !

Scarlet fever

Detected by FUNNEL

Detected!

World war II

## Q5

How can we remove <u>mistakes</u> and <u>incorrect values</u>?

P5 mistakes

It can also detect typos, "**automatically**" !!

Typhoid fever cases

Mistake

Missing values

Detected!

# Modeling power of FUNNEL

Our model can capture 5 properties

| P1 | Seasonality |
| P2 | Disease reductions |
| P3 | Area sensitivity |
| P4 | External events |
| P5 | Mistakes |

# Roadmap

✔ Motivation

✔ Modeling power of FUNNEL

– Overview – main ideas

– Proposed model – idea #1

– Algorithm – idea #2

– Experiments

– Discussion

– Conclusions

# Problem definition

Given:

Tensor $\mathcal{X}$ (disease x state x time)

$x$

Find:

Compact description of $\mathcal{X}$, *"automatically"*

$\mathcal{X}$ = FUNNEL

| B | R | N | $\mathcal{E}$ | $\mathcal{M}$ |
| P1 | P2 | P3 | P4 | P5 |

# Two main ideas

## Idea #1: Grey-box model



## Idea #2: MDL for fitting

**NO magic numbers !**
**(parameter-free)**

# Two main ideas

Idea #1: <u>Grey-box model</u> - domain knowledge



FUNNEL

**B** P1  **R** P2  **N** P3  $\mathcal{E}$ P4  $\mathcal{M}$ P5

$\mathcal{X}$ =

(SIRS+) : 6 parameters

$$S(t+1) = S(t) - \beta(t)\epsilon(t)S(t)I(t)$$
$$I(t+1) = I(t) + \beta(t)\epsilon(t)S(t)I(t)$$
$$V(t+1) = V(t) + \delta I(t) - \gamma V(t) +$$

Vaccine

Shocks

# Two main ideas

Idea #2: <u>Fitting with MDL</u> -> parameter free!

FUNNEL

$\mathcal{X}$ =

**B** P1   **R** P2   **N** P3   $\mathcal{E}$ P4   $\mathcal{M}$ P5

$$Cost_T(\mathcal{X}; \mathcal{F}) = \log^*(d) + \log^*(l) + \log^*(n)$$
$$+Cost_M(\mathbf{B}) + Cost_M(\mathbf{R}) + Cost_M(\mathbf{N})$$
$$+Cost_M(\mathcal{E}) + Cost_M(\mathcal{M}) + Cost_C(\mathcal{X}|\mathcal{F})$$

**NO magic numbers**

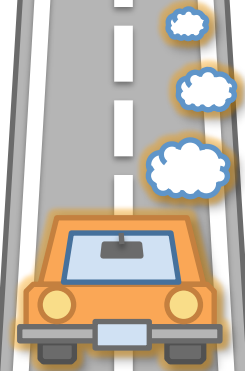**Parameter-free!**

# Roadmap
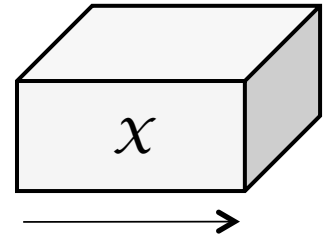
✔ Motivation

✔ Modeling power of FUNNEL

✔ Overview – main ideas

– Proposed model – idea #1

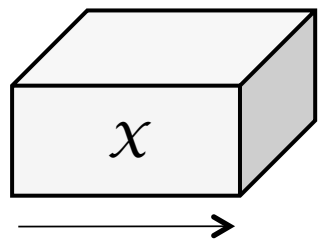– Algorithm – idea #2

– Experiments

– Discussion

– Conclusions

Y. Matsubara et al.

# Proposed model: FUNNEL



diseases $d$

states $l$

$\mathcal{X}$

Time $n$

single epidemic

Multi-evolving epidemics

(a) FUNNEL-single

(b) FUNNEL-full

# Proposed model: FUNNEL



states $l$

diseases $d$

$\mathcal{X}$

$n$ — Time

single epidemic

Multi-evolving epidemics

(a) FUNNEL-single

(b) FUNNEL-full

# FUNNEL – with a single epidemic

**Given:**

"single" epidemic sequence

e.g., measles in NY



**Find:**

nonlinear equation, model parameters

FUNNEL

# FUNNEL – with a single epidemic

## With a single epidemic: Funnel-RE

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



Linear

I(t)

S(t)

Log

V(t)

I(t)

# FUNNEL – with a single epidemic

With a single epidemic: Funnel-RE

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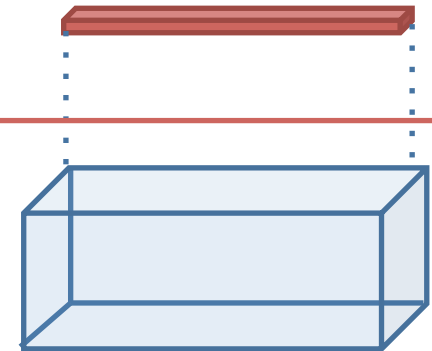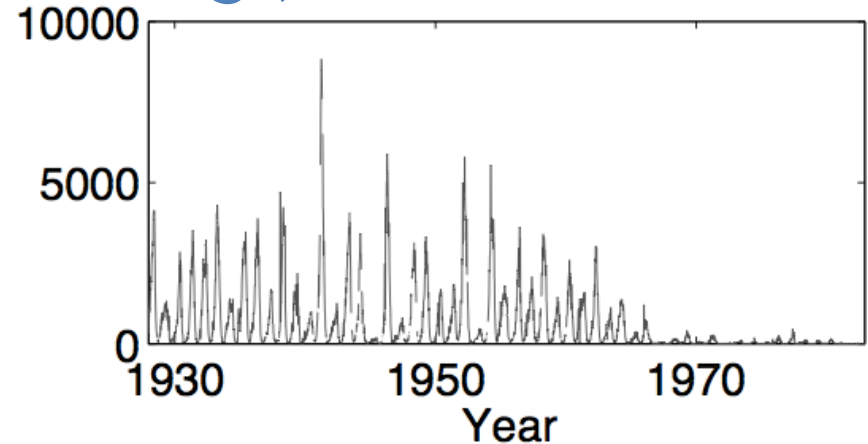I(t+1) &= I(t) + \beta(t)\epsilon(t)S(t)I(t) - \delta I(t) \\
V(t+1) &= V(t) + \delta I(t) - \gamma V(t) + \theta(t)S(t)
\end{aligned} \quad (3)
$$

**S(t)** : susceptible
**I (t)** : Infected
**V(t)** : Vigilant
/Vaccinated

**Details**

## With a single epidemic: Funnel-RE

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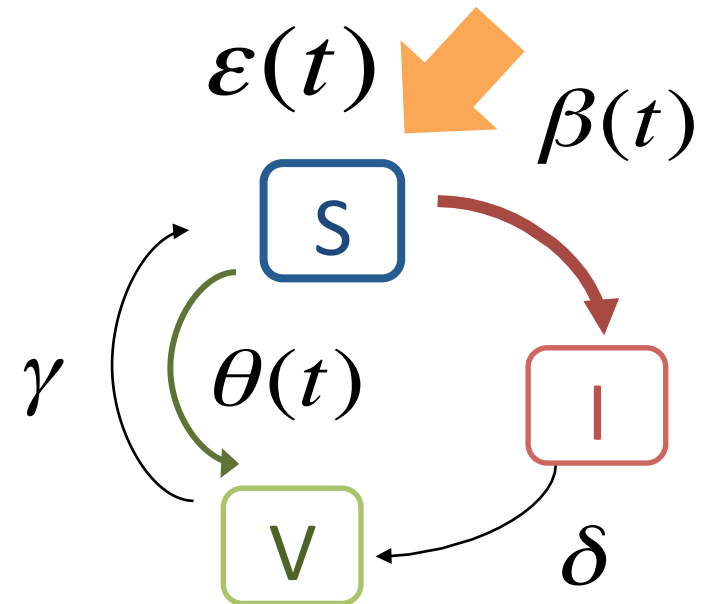I(t+1) &= I(t) + \boxed{\beta(t)}\epsilon(t)S(t)I(t) - \delta I(t) \\
V(t+1) &= V(t) + \delta I(t) - \gamma V(t) + \theta(t)S(t)
\end{aligned} \qquad (3)$$

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

$$\beta(t) = \beta_0 \cdot \left(1 + P_a \cdot cos\left(\tfrac{2\pi}{P_p}(t + P_s)\right)\right)$$
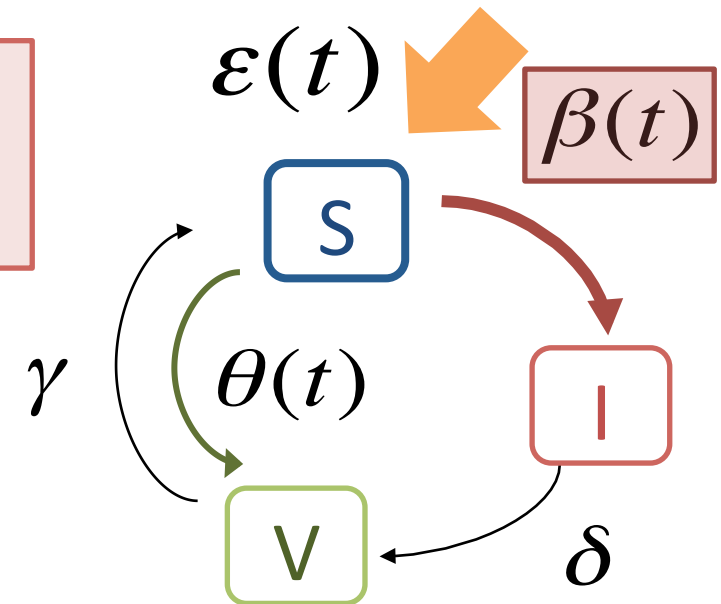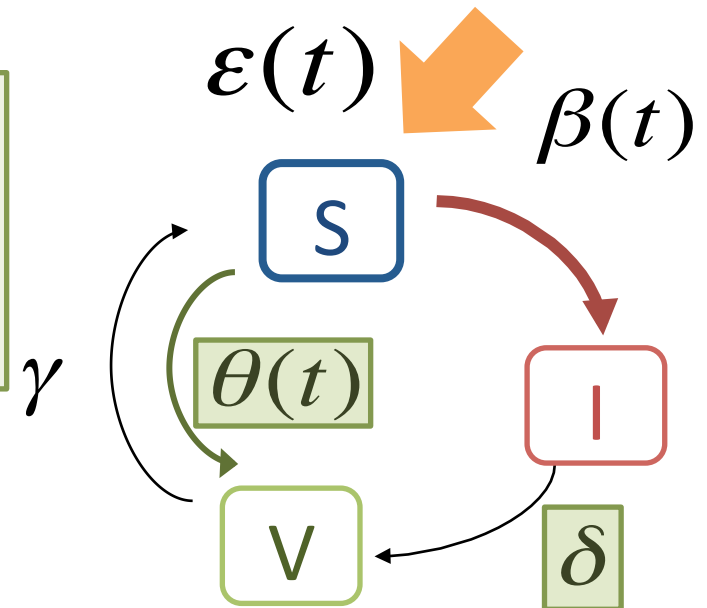
$$P_p = 52$$

## With a single epidemic: Funnel-RE

$$S(t+1) = S(t) - \beta(t)\epsilon(t)S(t)I(t) + \gamma V(t) - \boxed{\theta(t)}S(t)$$

$$I(t+1) = I(t) + \beta(t)\epsilon(t)S(t)I(t) - \boxed{\delta}I(t)$$

$$V(t+1) = V(t) + \boxed{\delta}I(t) - \gamma V(t) + \boxed{\theta(t)}S(t) \tag{3}$$

$\delta$ : healing rate

$\theta(t)$ : disease reduction effect

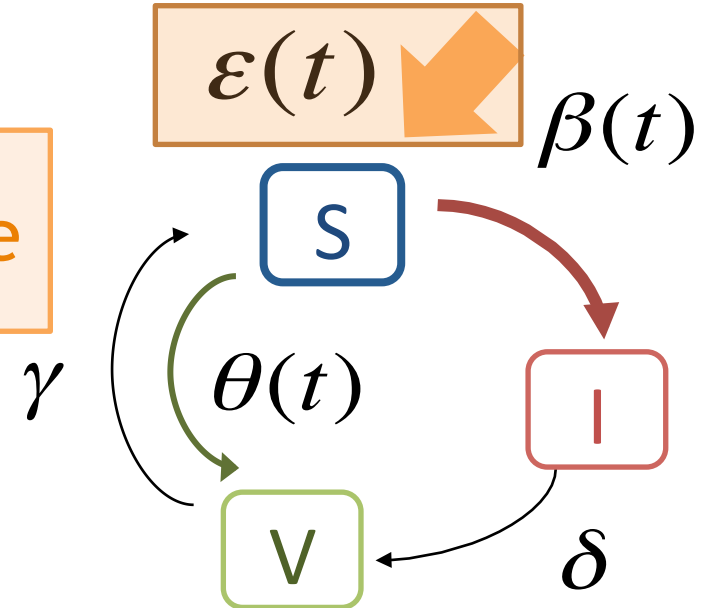$$\theta(t) = \begin{cases} 0 & (t < t_\theta) \\ \theta_0 & (t \geq t_\theta) \end{cases}$$

$\varepsilon(t)$
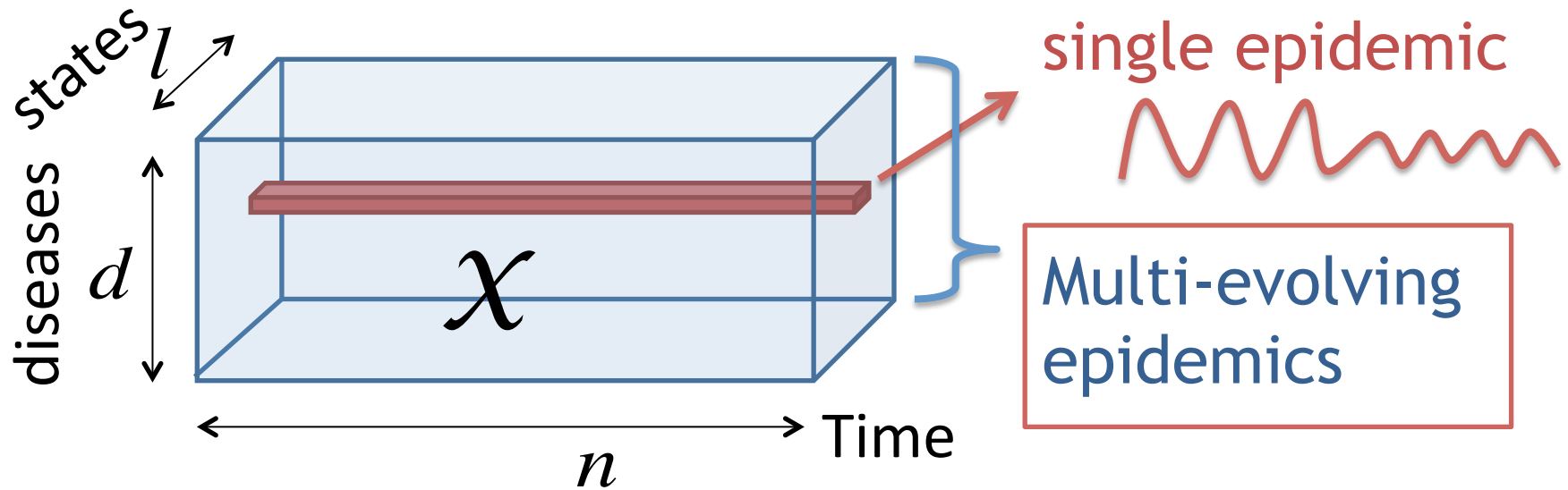
$\beta(t)$

S

$\theta(t)$

$\gamma$

I

V

$\delta$

Details

## With a single epidemic: Funnel-RE

$$S(t+1) = S(t) - \beta(t)\epsilon(t)S(t)I(t) + \gamma V(t) - \theta(t)S(t)$$

$$I(t+1) = I(t) + \beta(t)\epsilon(t)S(t)I(t) - \delta I(t)$$

$$V(t+1) = V(t) + \delta I(t) - \gamma V(t) + \theta(t)S(t) \tag{3}$$

$\varepsilon(t)$

$\beta(t)$

$\varepsilon(t)$ : temporal susceptible rate

S

$\gamma$    $\theta(t)$

I

V

$\delta$

# Proposed model: FUNNEL



states $l$

diseases $d$

$\mathcal{X}$

$n$ → Time

single epidemic

Multi-evolving epidemics

(a) FUNNEL-single

(b) FUNNEL-full

# Proposed model: FUNNEL-full



P1 P2 global/country   P3 local/state

$$\mathcal{X} = B, R, N,$$

with $\mathcal{X}$ having dimensions: states $l$, diseases $d$, time $n$.

$B$ parameters: $N, \beta_0, \delta, \gamma, P_a, P_s$ (size $d \times 6$)

$R$ parameters: $\theta_0, t_\theta$ (size $d \times 2$)

$N$ (size $d \times l$)

P4 P5 extra - $\mathcal{E}$: shocks & $\mathcal{M}$: mistakes

$$\mathcal{E}, \mathcal{M}$$

# Proposed model: FUNNEL-full

P1 P2 global/country

$$x = B, R,$$

where the base matrix $\mathbf{B}$ ($d \times 6$) contains $N, \beta_0, \delta, \gamma, P_a, P_s$ and the disease reduction matrix $\mathbf{R}$ ($d \times 2$) contains $\theta_0, t_\theta$.

- states $l$
- diseases $d$
- time $n$

**Global**

**P1** Base matrix $\mathbf{B}$ (d x 6)

**P2** Disease reduction matrix $\mathbf{R}$ (d x 2)

# Proposed model: FUNNEL-full

**Details**

states

$l$

diseases

$d$

$x$

$n$

time

$=$

$l$
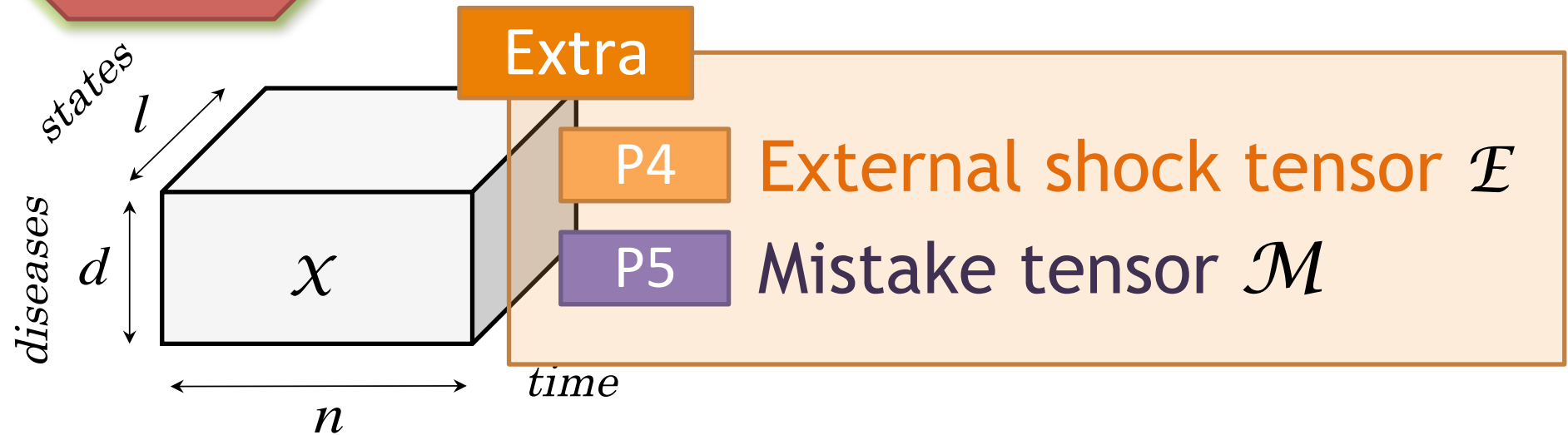
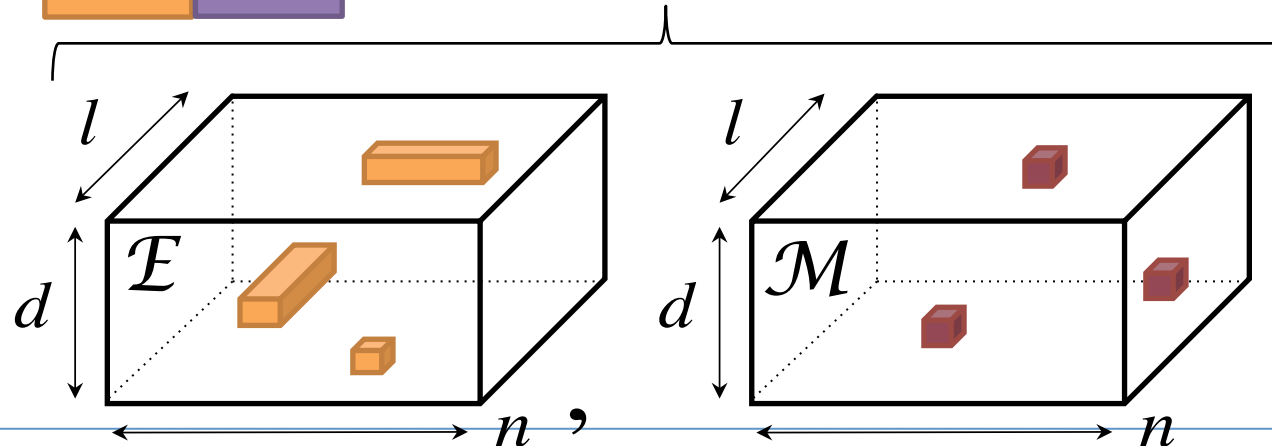$\mathbf{N}$

$d$

**Local**

P3  Geo-disease matrix $\mathbf{N}$ (d x l)

$$\mathbf{N} = \{N_{ij}\}_{i,j=1}^{d,l}$$ : potential population of disease i in state j
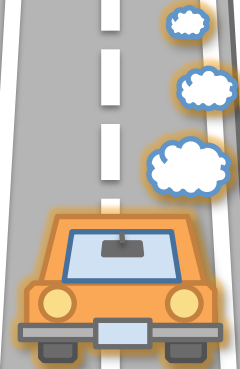
# Proposed model: FUNNEL-full

Extra

P4 External shock tensor $\mathcal{E}$

P5 Mistake tensor $\mathcal{M}$

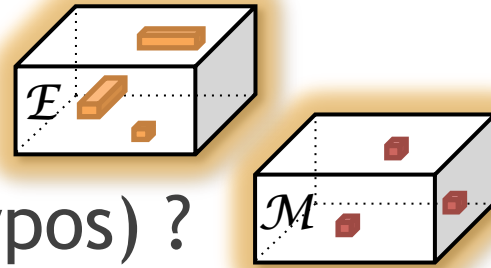P4 P5 extra - $\mathcal{E}$: shocks & $\mathcal{M}$: mistakes

# Roadmap

✔ Motivation

✔ Modeling power of FUNNEL

✔ Overview – main ideas

✔ Proposed model – idea #1

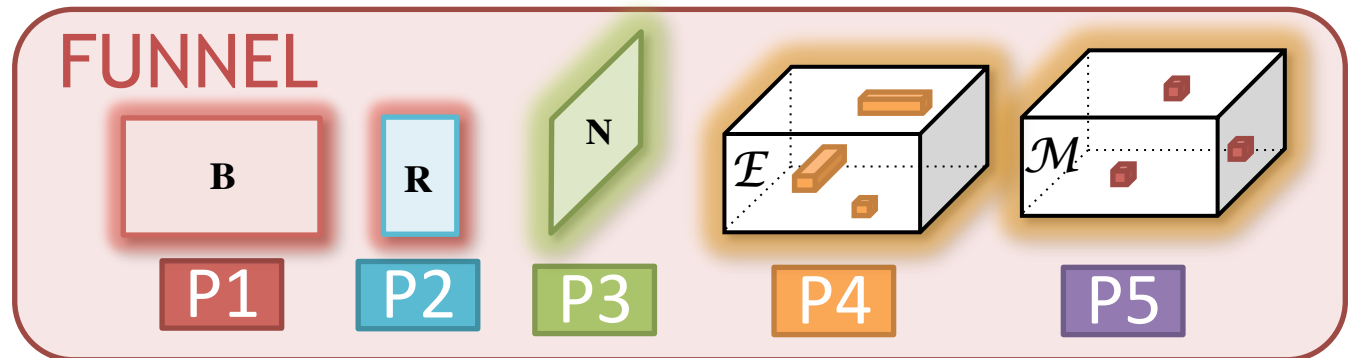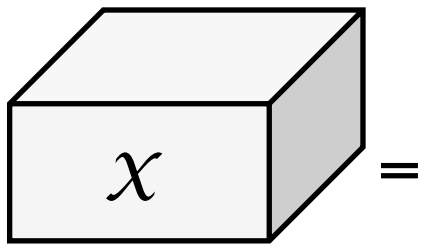– Algorithm – idea #2

– Experiments

– Discussion

– Conclusions

# Challenges

**Q1.** How to automatically
- find "external shocks" ?
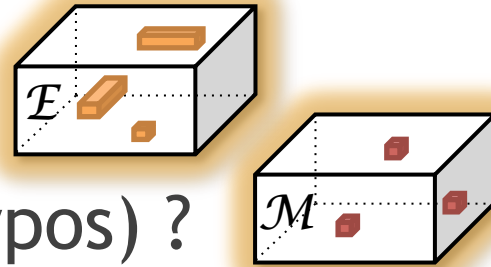- ignore "mistakes" (i.e., typos) ?

$\mathcal{E}$

$\mathcal{M}$

**Q2.** How to efficiently estimate model parameters ?

$$\mathcal{X} \quad = \quad \text{FUNNEL}$$

| B | R | N | $\mathcal{E}$ | $\mathcal{M}$ |
|---|---|---|---|---|
| P1 | P2 | P3 | P4 | P5 |

# Challenges

**Q1.** How to automatically
- find "external shocks" ?
- ignore "mistakes" (i.e., typos) ?

$\mathcal{E}$

$\mathcal{M}$

> Idea (1) : <u>Model description cost</u>

**Q2.** How to efficiently estimate model parameters ?

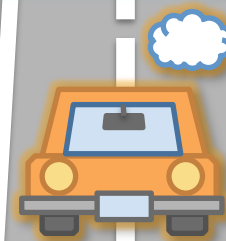$\mathcal{X}$ =

FUNNEL

**B**   **R**   **N**   $\mathcal{E}$   $\mathcal{M}$

P1   P2   P3   P4   P5

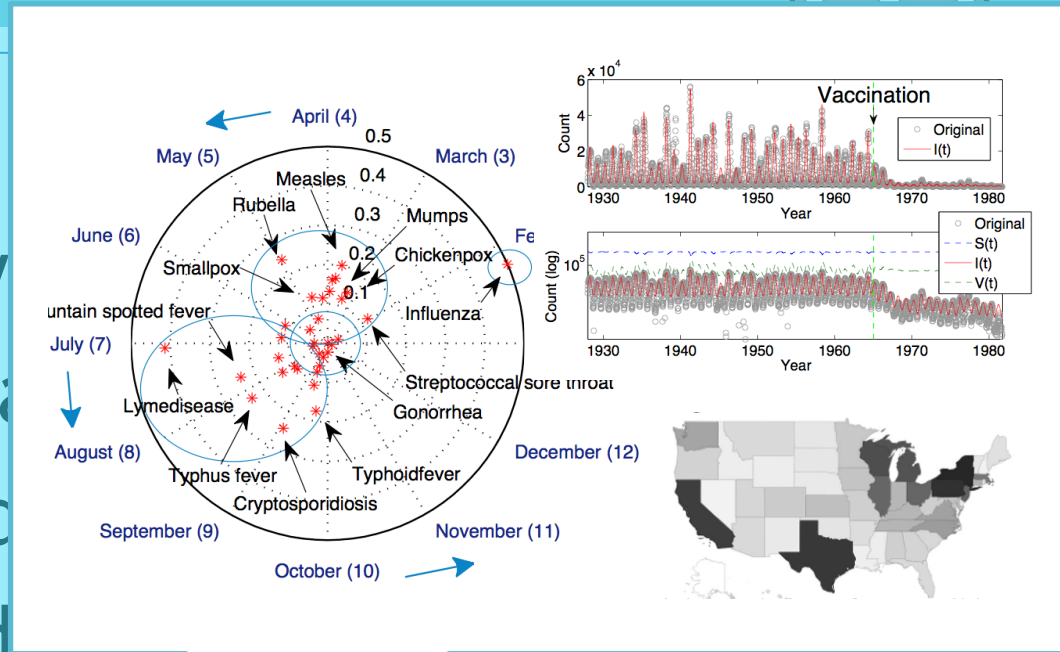> Idea (2): <u>Multi-layer optimization (linear)</u>

# Roadmap

- ✔ Motivation
- ✔ Modeling power of FUNNEL
- ✔ Overview – main ideas
- ✔ Proposed model – idea #1
- ✔ Algorithm – idea #2
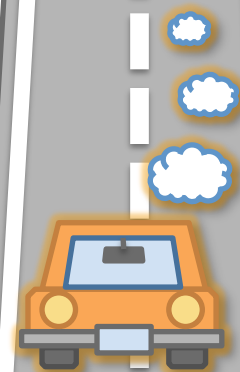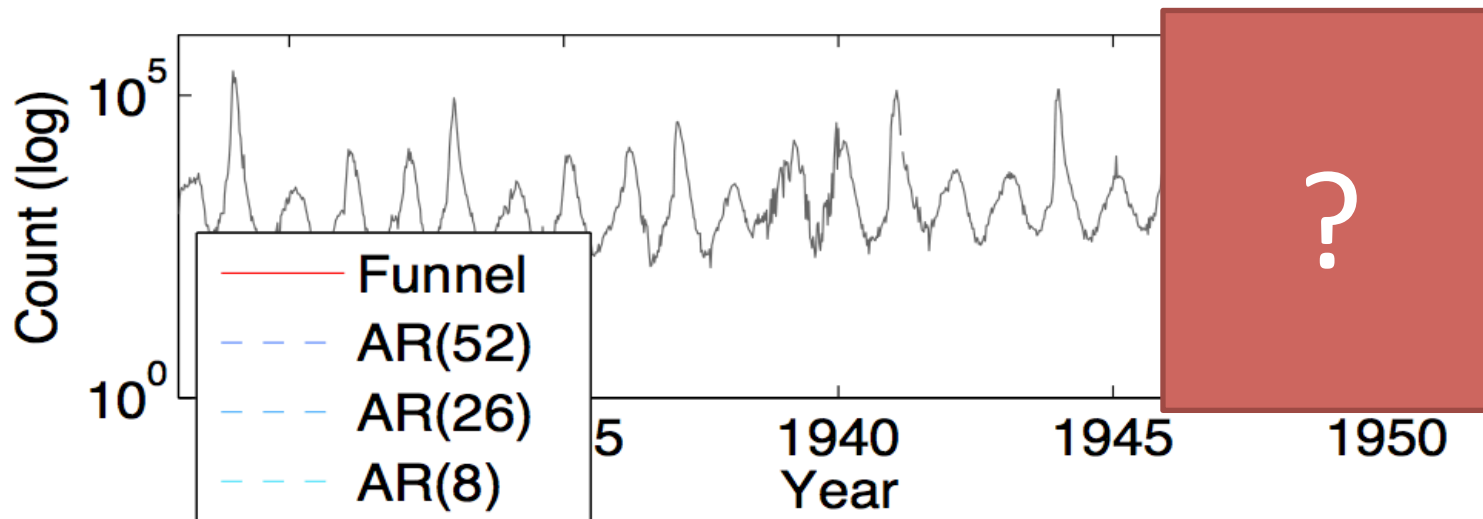- – Experiments
- – Discussion
- – Conclusions

# Roadmap

- ✔ Motivation
- ✔ Modeling pow
- ✔ Overview – ma
- ✔ Proposed mo
- ✔ Algorithm – idea
- ✔ Experiments
- Discussion
- Conclusions

# Roadmap

✔ Motivation

✔ Modeling power of FUNNEL

✔ Overview – main ideas

✔ Proposed model – idea #1

✔ Algorithm – idea #2

✔ Experiments

– Discussion

– Conclusions

## Forecasting future epidemics

| Train: **2/3** sequences | Forecast: **1/3** following years |



(a) Influenza

## Forecasting future epidemics

| Train: 2/3 sequences | Forecast: 1/3 following years |



(a) Influenza

**Funnel** can capture  future epidemics (AR: fail)

## Forecasting future epidemics

| Train: **2/3** sequences | Forecast: **1/3** following years |



(c) Typhoid fever

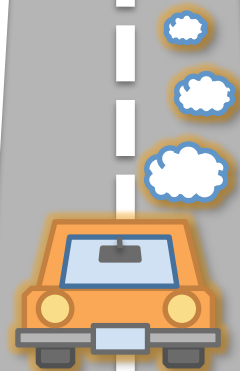**Funnel** can capture  future epidemics (AR: fail)

## Epidemics on computer networks

Spread via email attachment

Spread through corporate networks



(10 years)

**Funnel** is general: it fits computer virus very well!

# Roadmap

✔ Motivation

✔ Modeling power of FUNNEL

✔ Overview – main ideas

✔ Proposed model – idea #1

✔ Algorithm – idea #2

✔ Experiments

✔ Discussion

– Conclusions

# Conclusions

FUNNEL has the following advantages
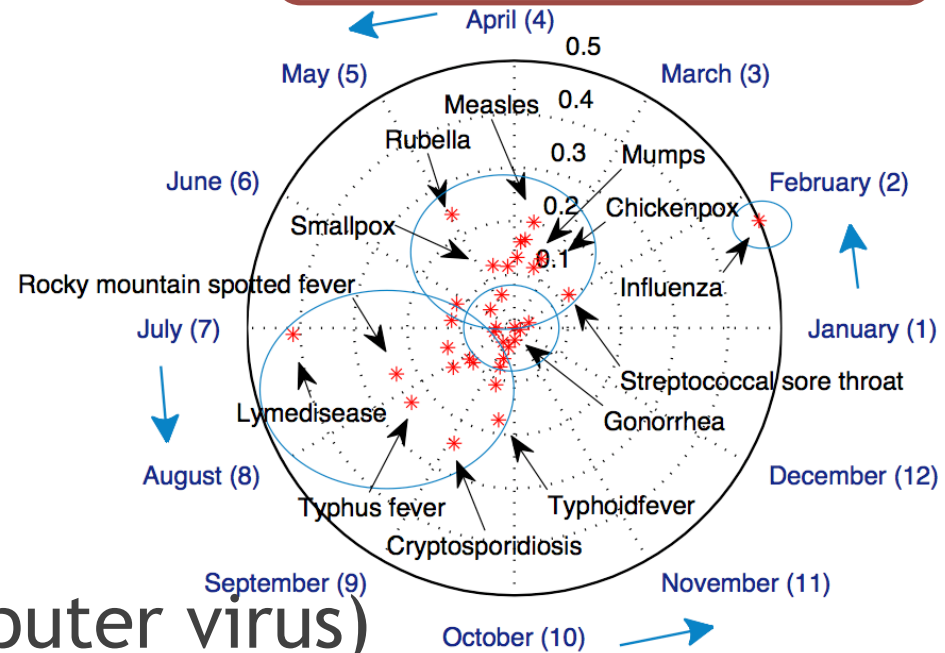
✔ **Sense-making**

Captures all essential aspects:

| P1 | P2 | P3 | P4 | P5 |

✔ **Fully-automatic**

No training set

✔ **Scalable**

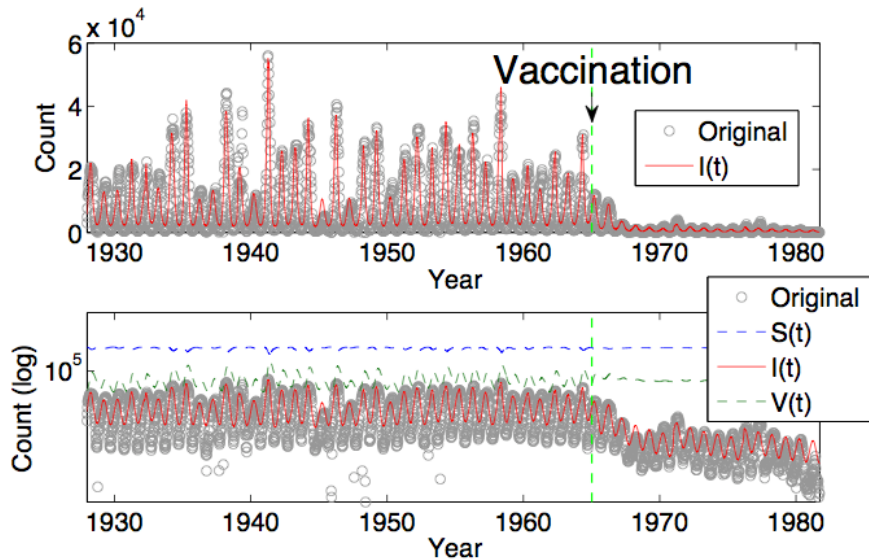It scales linearly

✔ **General**

Real epidemics (+ computer virus)

# Thank you!



Data: http://www.tycho.pitt.edu/

Code: http://www.cs.kumamoto-u.ac.jp/~yasuko/software.html

# FUNNEL: Automatic Mining of Spatially Coevolving Epidemics

Yasuko Matsubara, Yasushi Sakurai (Kumamoto University)

Willem G. van Panhuis (University of Pittsburgh)
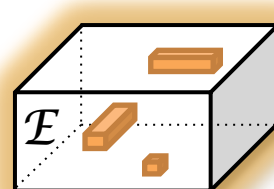
Christos Faloutsos (CMU)
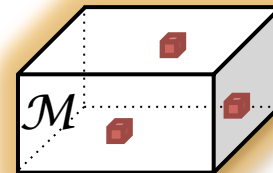
# Proposed model: FUNNEL-full
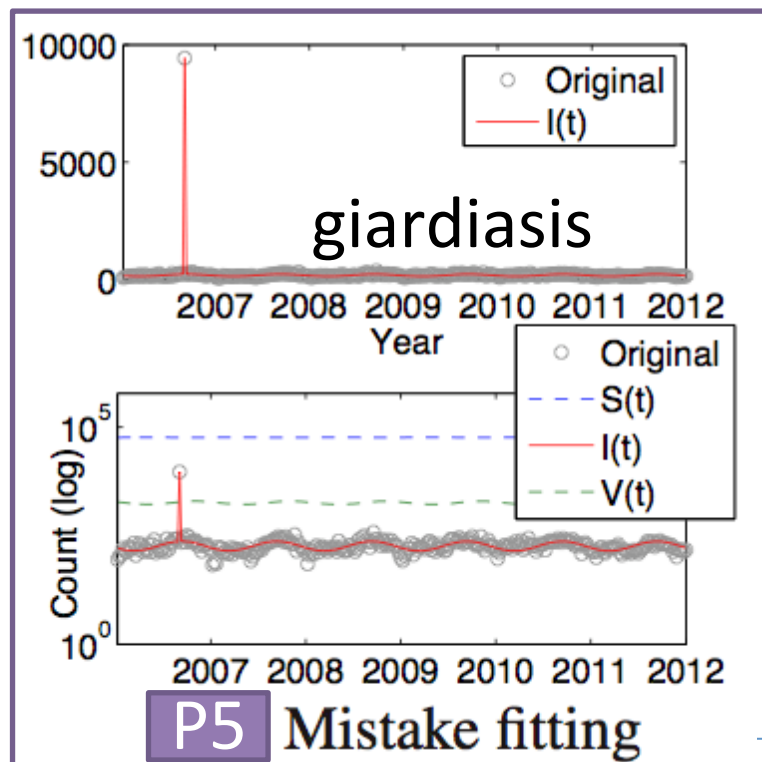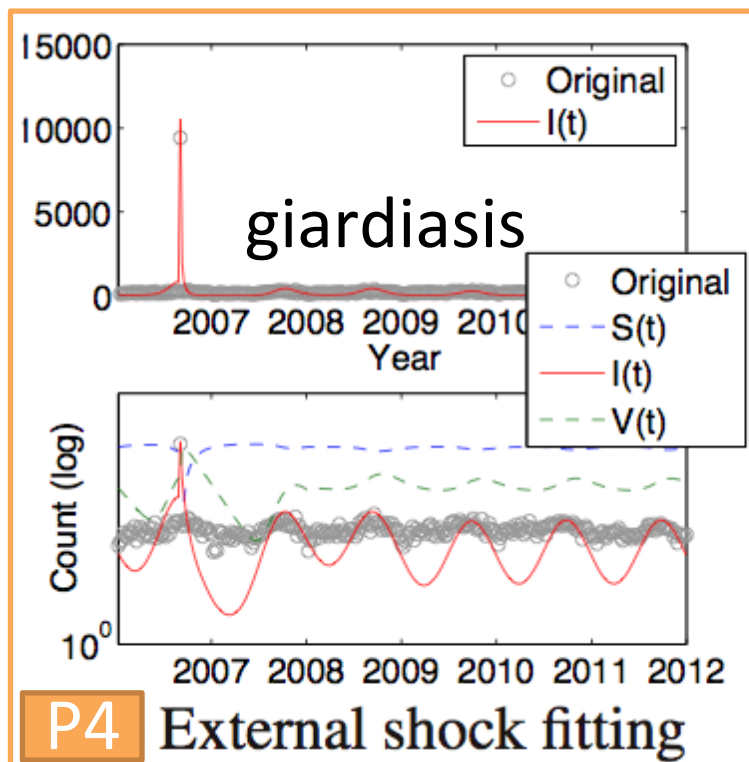
**Details**

## What's the difference??

**External shock** vs. **Mistake**    $\mathcal{E}$    vs.    $\mathcal{M}$    P4    P5

P4 **External shock fitting**
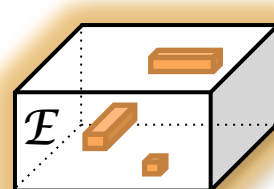
giardiasis

P5 **Mistake fitting**

giardiasis

# Proposed model: FUNNEL-full

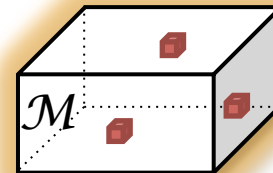**Details**

What's the difference??

External shock vs. Mistake

$\mathcal{E}$ vs. $\mathcal{M}$

P4     P5

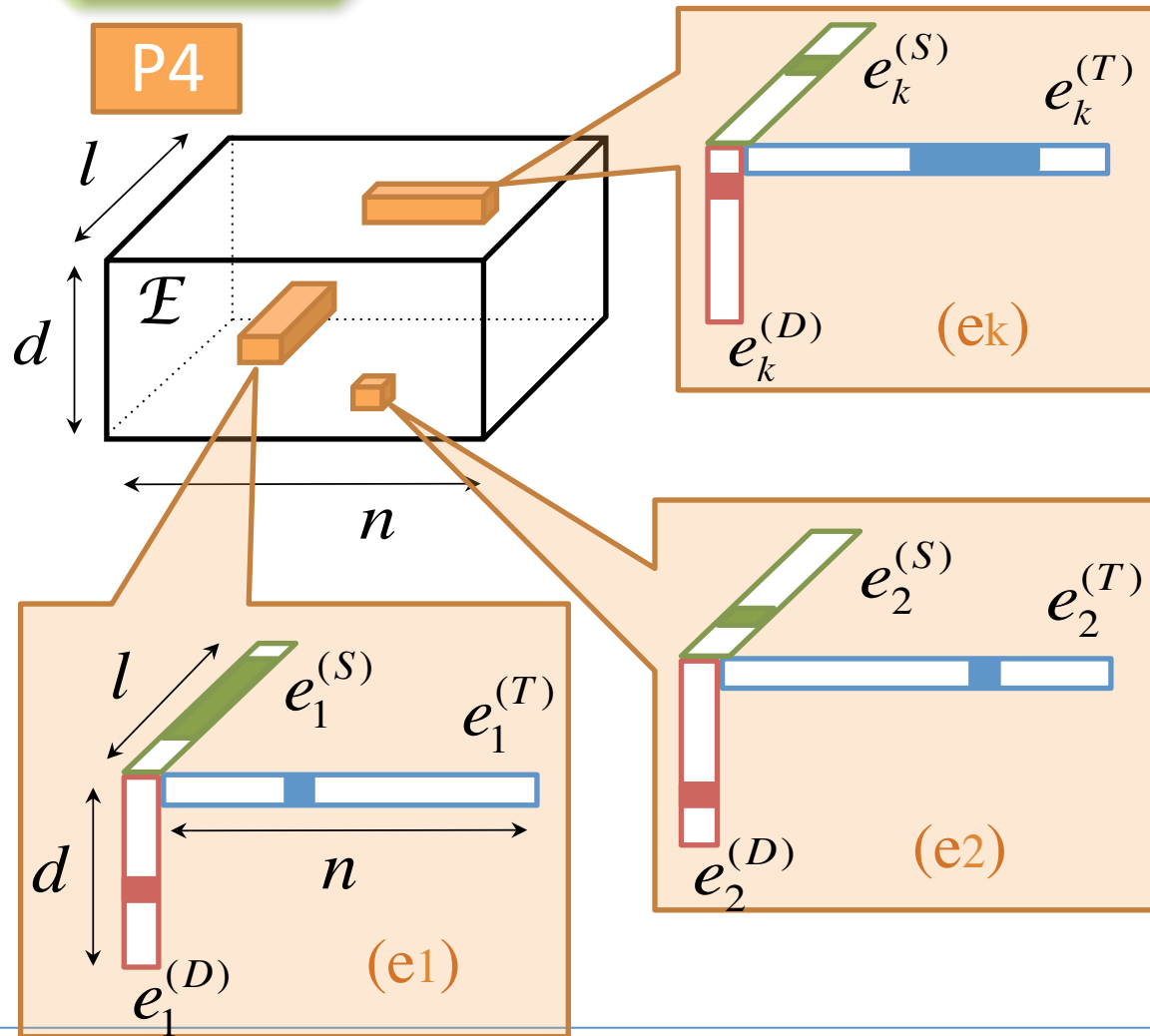influenced

independent

P4 **External shock fitting**

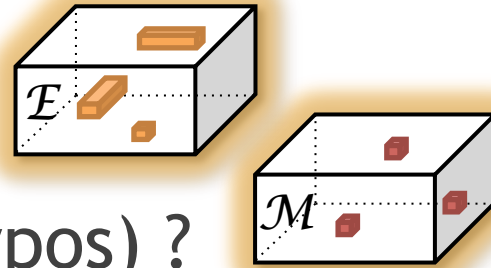P5 **Mistake fitting**

# Proposed model: FUNNEL-full

$$\mathcal{E} = \{\mathbf{E}^{(D)}, \mathbf{E}^{(T)}, \mathbf{E}^{(S)}\}$$

Disease matrix   Time matrix   State matrix

# Idea (1): Model description cost

**Q1**. How should we
- find "external shocks" ?
- ignore "mistakes" (i.e., typos) ?



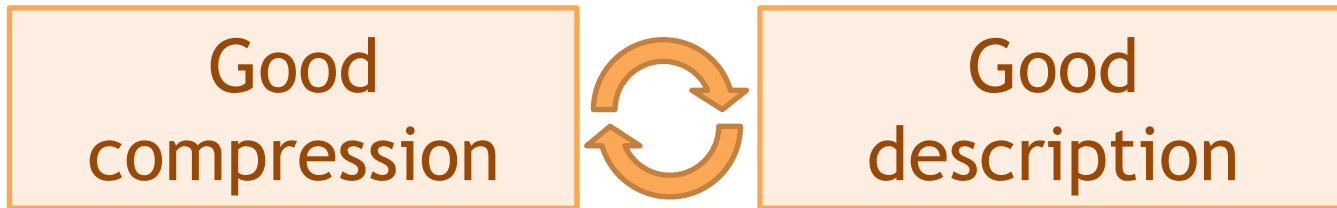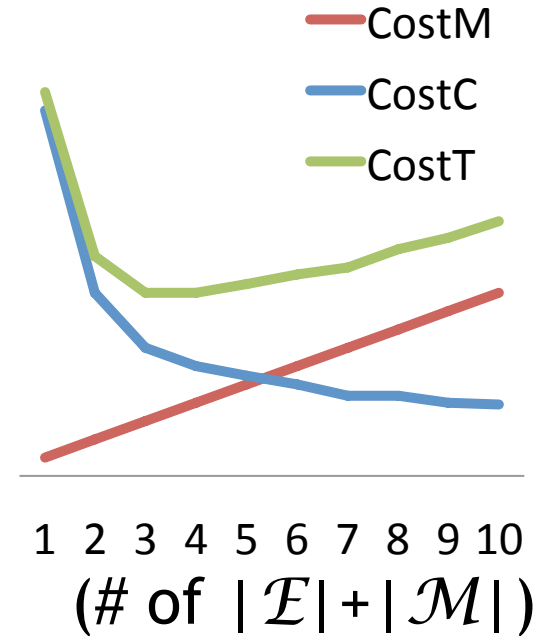Idea (1) : <u>Model description cost</u>

- Minimize coding cost
- find "optimal" # of externals/mistakes
- "automatically"

# Idea (1): Model description cost

Idea: Minimize encoding cost!



$$\min \left( \boxed{\text{Cost}_M(\mathcal{F})} + \boxed{\text{Cost}_C(\mathcal{X}|\mathcal{F})} \right)$$

Model cost          Coding cost

Good compression ⟳ Good description

# Idea (1): Model description cost

Total cost of tensor $\mathcal{X}$, given $\mathcal{F}$

$$\mathcal{F} = \{\mathbf{B}, \mathbf{R}, \mathbf{N}, \mathcal{E}, \mathcal{M}\}.$$
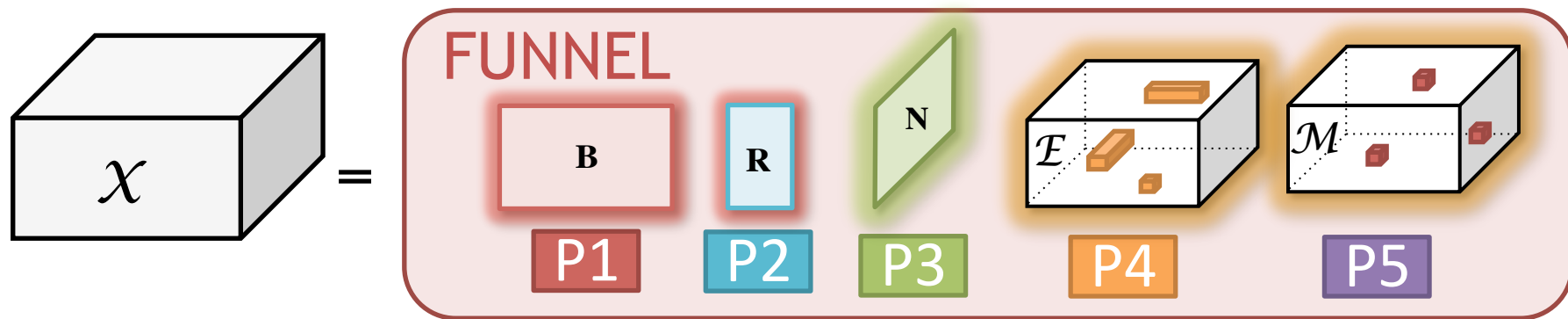
$$Cost_T(\mathcal{X}; \mathcal{F}) = \log^*(d) + \log^*(l) + \log^*(n)$$
$$+ Cost_M(\mathbf{B}) + Cost_M(\mathbf{R}) + Cost_M(\mathbf{N})$$
$$+ Cost_M(\mathcal{E}) + Cost_M(\mathcal{M}) + Cost_C(\mathcal{X}|\mathcal{F})$$

# Idea (1): Model description cost

Total cost of tensor $\mathcal{X}$, given $\mathcal{F}$

$$\mathcal{F} = \{\mathbf{B}, \mathbf{R}, \mathbf{N}, \mathcal{E}, \mathcal{M}\}.$$

Dimensions of X

$$Cost_T(\mathcal{X}; \mathcal{F}) = \log^*(d) + \log^*(l) + \log^*(n)$$
$$+ Cost_M(\mathbf{B}) + Cost_M(\mathbf{R}) + Cost_M(\mathbf{N})$$
$$+ Cost_M(\mathcal{E}) + Cost_M(\mathcal{M}) + Cost_C(\mathcal{X}|\mathcal{F})$$

Model description cost of F

Coding cost of X given F

# Idea (2): Multi-layer optimization
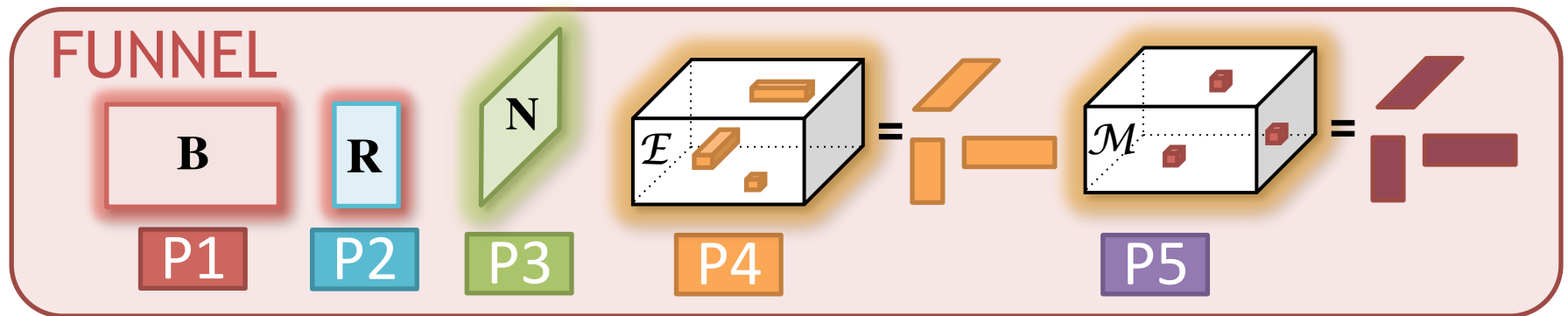
**Q2.** How to efficiently estimate model parameters ?



Idea (2): Multi-layer optimization
- Find "optimal" solution w.r.t.
  - Global level parameters
  - Local level parameters

Find "optimal" solution w.r.t. **Global** **Local**
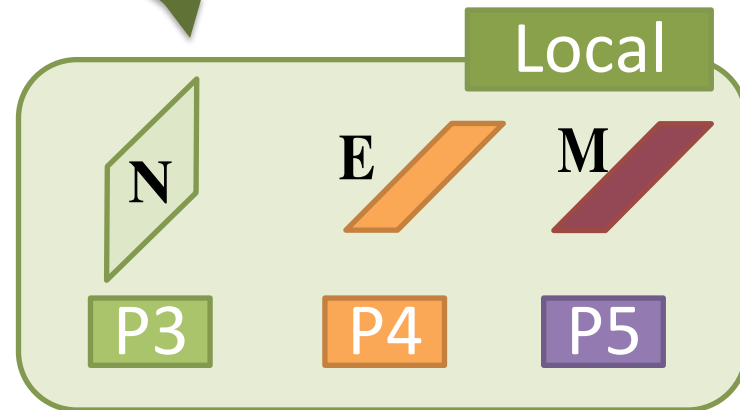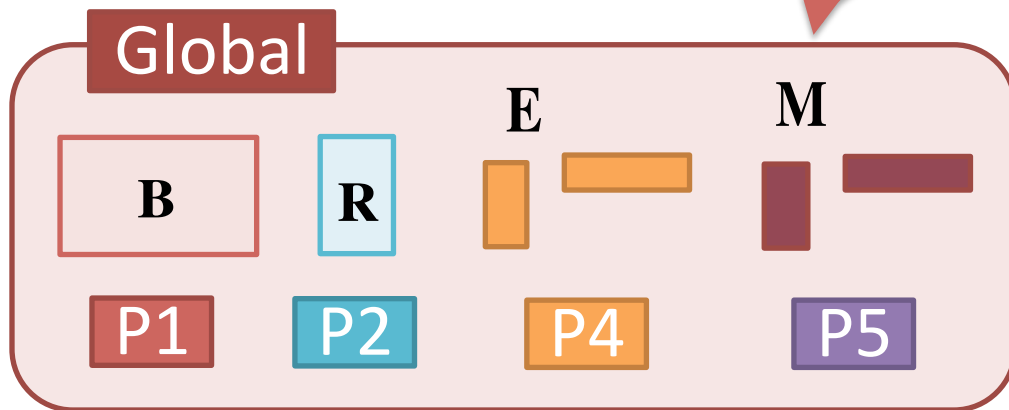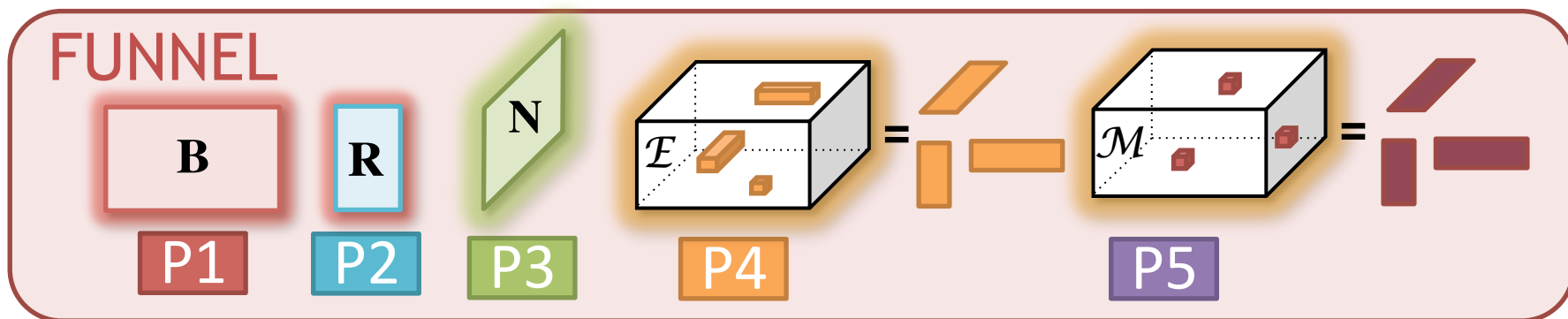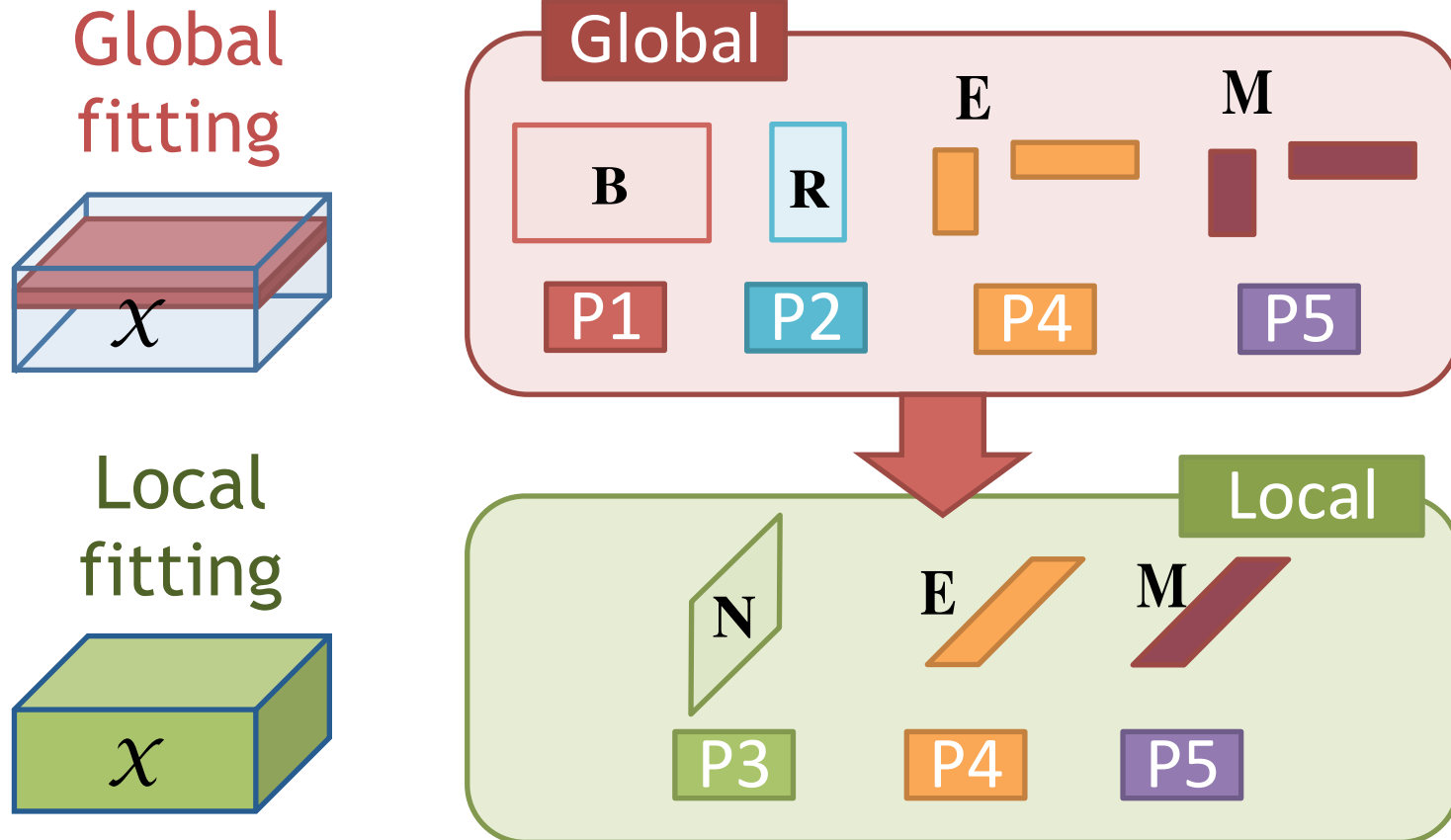
FUNNEL

**P1** **P2** **P3** **P4** **P5**

**Global** **Local**

# Idea (2): Multi-layer optimization

Find "optimal" solution w.r.t. Global Local

# Idea (2): Multi-layer optimization

## Multi-layer fitting algorithm

Global fitting

$x$

Local fitting

$x$

**Global**

**E**  **M**

**B**  **R**

P1  P2  P4  P5

**Local**

**N**  **E**  **M**

P3  P4  P5

# Experiments

We answer the following questions...

Q1. Sense-making

Can it help us understand the given epidemics?

Q2. Accuracy

How well does it match the data?

Q3. Scalability

How does it scale in terms of computational time?

# Q1. Sense-making

Our preliminary observations:

- **P1** yearly periodicity
- **P2** disease reduction effects
- **P3** area specificity and sensitivity
- **P4** external shock events
- **P5** mistakes, incorrect values

# Q1. Sense-making

**P1** Disease seasonality
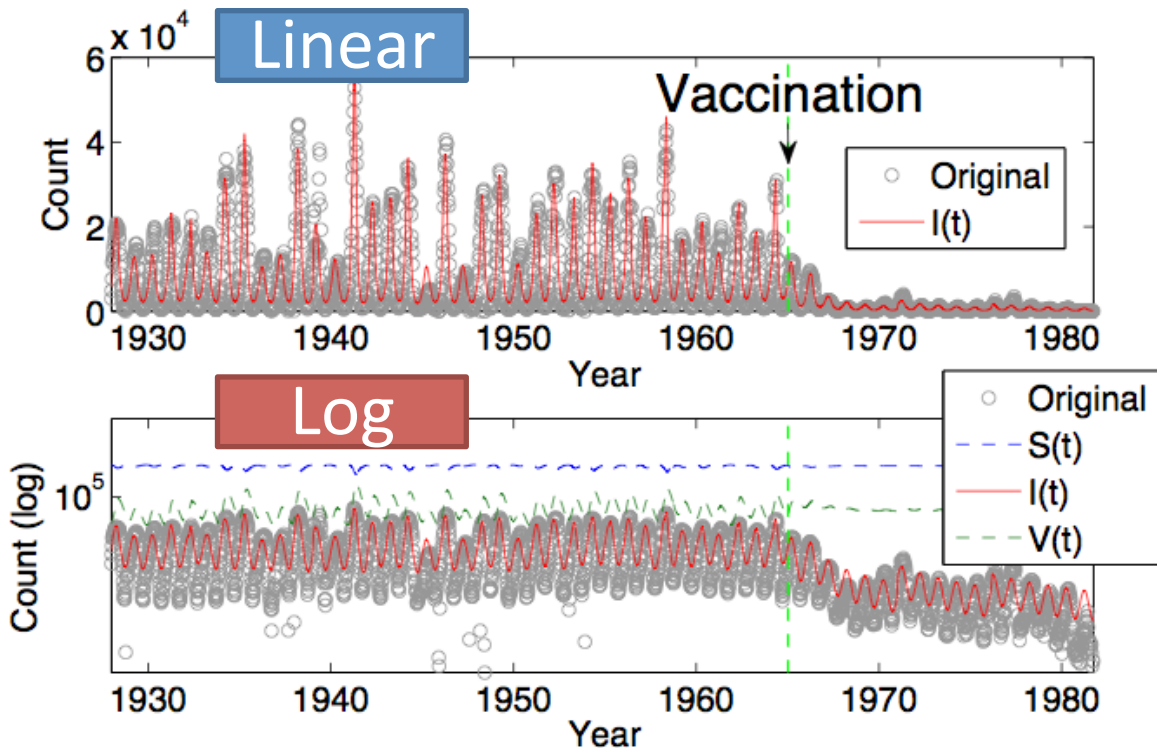


Radius:
seasonality strength
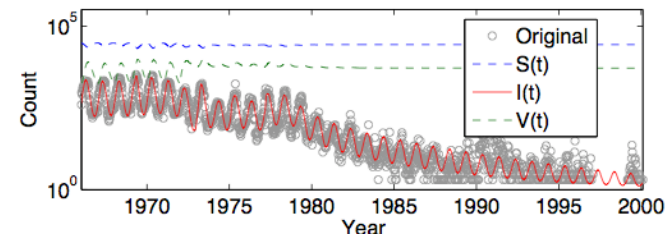Angle:
peak season

# Q1. Sense-making

P1  Disease seasonality



Influenza in Feb.

Children's in spring

Tick-borne in summer

Gonorrhea no periodicity

April (4)
0.5
May (5)
March
Measles 0.4
Rubella
0.3 Mumps
June (6)
Chickenpox
Smallpox 0.2
February (2)
0.1
Rocky mountain spotted fever
Influenza
July (7)
January (1)
Lyme disease
Streptococcal sore throat
Gonorrhea
August (8)
December (12)
Typhus fever
Typhoid
Cryptosporidiosis
September (9)
October (10)

# Q1. Sense-making

**P2** Disease reduction effect

Measles (vaccine licensure: 1963)



Linear

Log

1967

(c) Mumps (**P1**), (**P2**), (**P4**)

1969

(d) Rubella (**P1**), (**P2**), (**P4**)

# Q1. Sense-making

Potential population of susceptibles (measles)



CA

TX

NY, PA

**FL (less children)**

# Q1. Sense-making

P3  area specificity and sensitivity

Measles in NY and PA



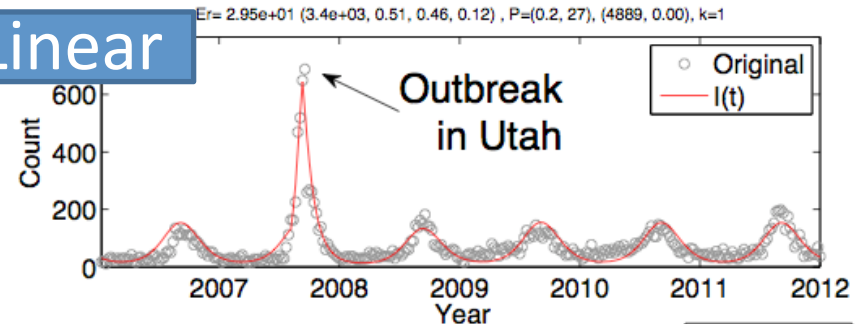Linear

Log

(a) New York State (NY)
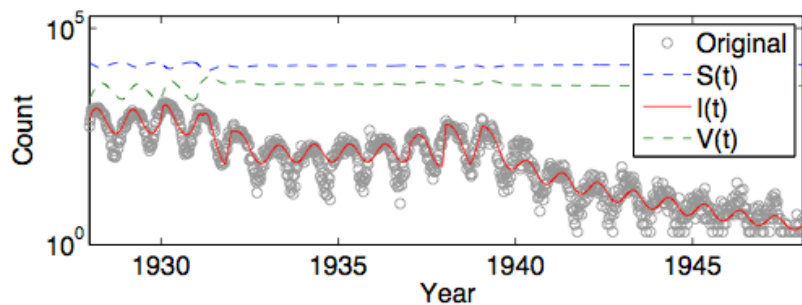
(b) Pennsylvania (PA)

P4 | external shock events

We can detect external shocks "**automatically**" !!



(h) Smallpox (**P1**), (**P2**), (**P4**)

(j) Cryptosporidiosis (**P1**), (**P3**), (**P4**)

P5 mistakes, incorrect values

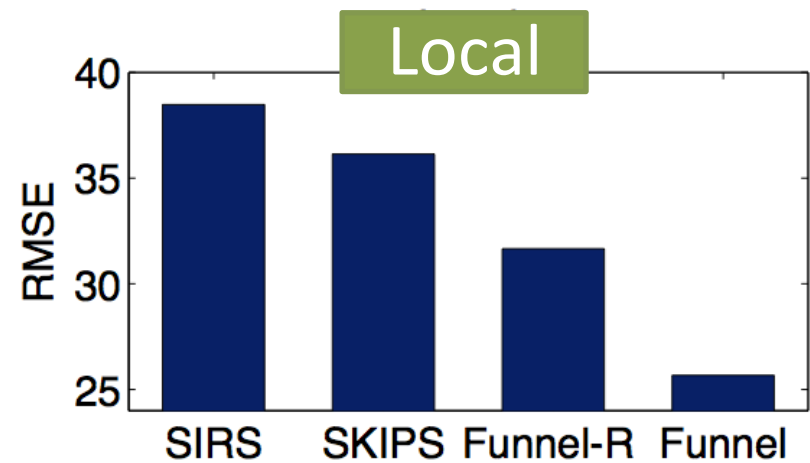We can also detect typos "**automatically**" !!
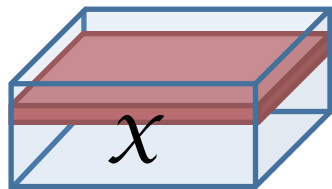


Linear

Missing values

Mistake

Log

# Q2. Accuracy

Fitting accuracy for Global Local sequences
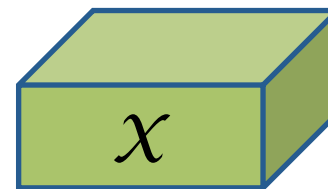(lower is better)



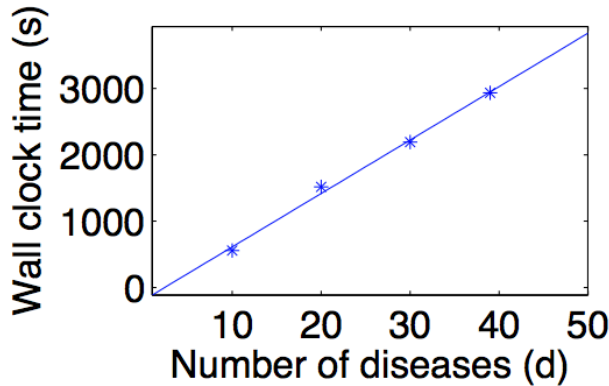(a) Global fitting      (b) Local fitting

$$\{\bar{x}_i(t)\}_{i,t}^{d,n}$$

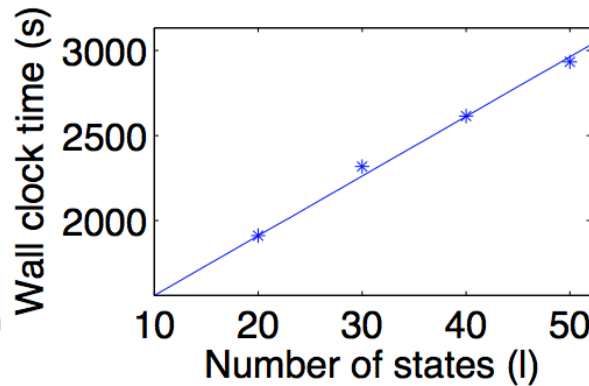$$\bar{x}_i(t) = \sum_{j=1}^{l} x_{ij}(t)$$
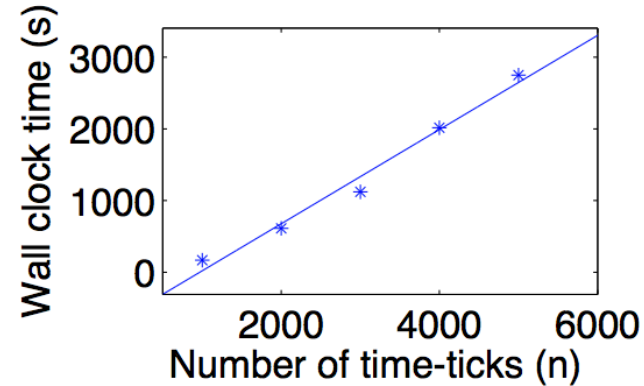
$$\{x_{ij}(t)\}_{i,j,t}^{d,l,n}$$

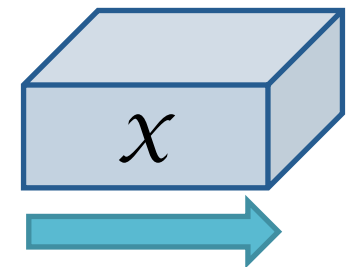# Q3. Scalability

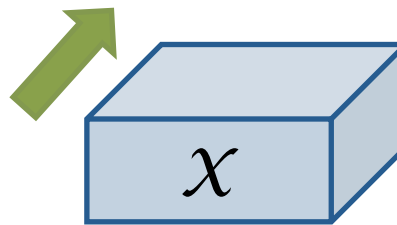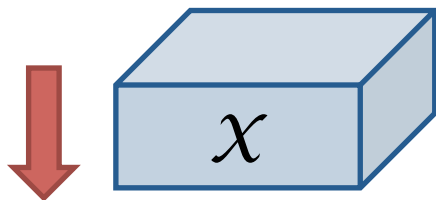Wall clock time vs. diseases , states , Time



(a) Diseases ($d$)    (b) States ($l$)    (c) Duration ($n$)

FunnelFit is linear w.r.t. data size : O(dln)