# FUNNEL: Automatic Mining of Spatially Coevolving Epidemics
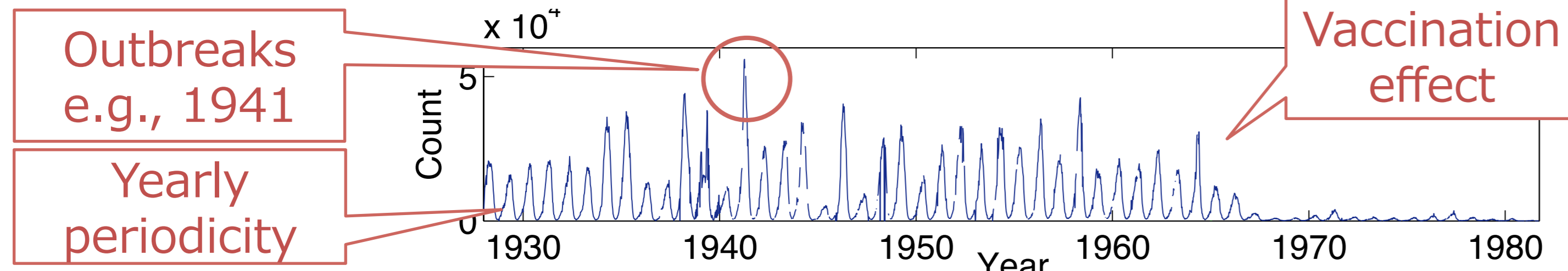
Kumamoto University

**Yasuko Matsubara**
Kumamoto University
yasuko@cs.kumamoto-u.ac.jp

**Yasushi Sakurai**
Kumamoto University
yasushi@cs.kumamoto-u.ac.jp

**Willem G. van Panhuis**
University of Pittsburgh
wav10@pitt.edu

**Christos Faloutsos**
Carnegie Mellon University
christos@cs.cmu.edu

## Motivation - Given: large set of epidemiological data

e.g., Measles cases from 1928 to 1982 (50 states)

Outbreaks e.g., 1941
Yearly periodicity
Vaccination effect

**Goal:** statistically summarize all the epidemic time-series
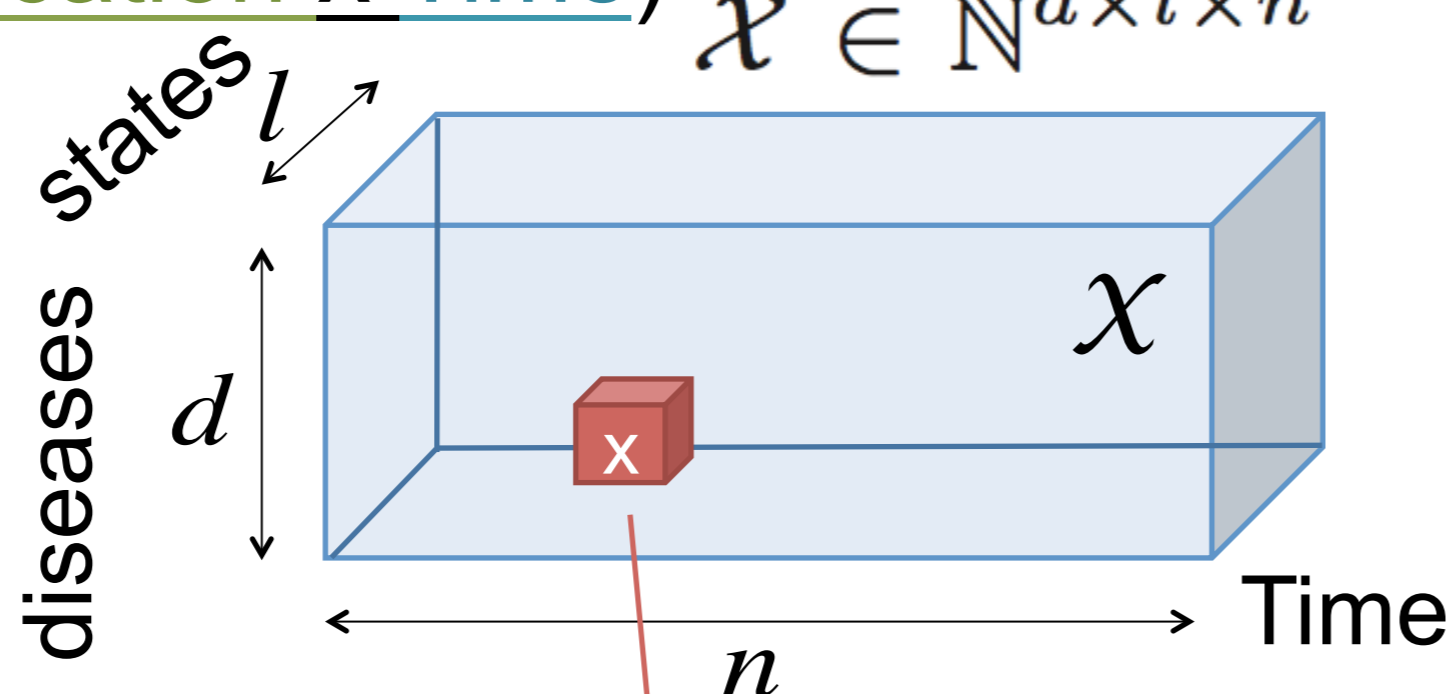
## Data description – Project Tycho

PROJECT TYCHO DATA FOR HEALTH

- 56 contagious diseases for U.S. states
- from 1888 to the present (>125 years)

3rd order tensor (diseases x location x Time)

$\mathcal{X} \in \mathbb{N}^{d \times l \times n}$

# of cases in 1931, …

| Time | disease | loc | cases |
|------|---------|-----|-------|
| 04-01 | measles | PA | 4740 |
| 04-01 | measles | NY | 5310 |
| 04-01 | rubella | CA | 1923 |
| … | … | … | … |

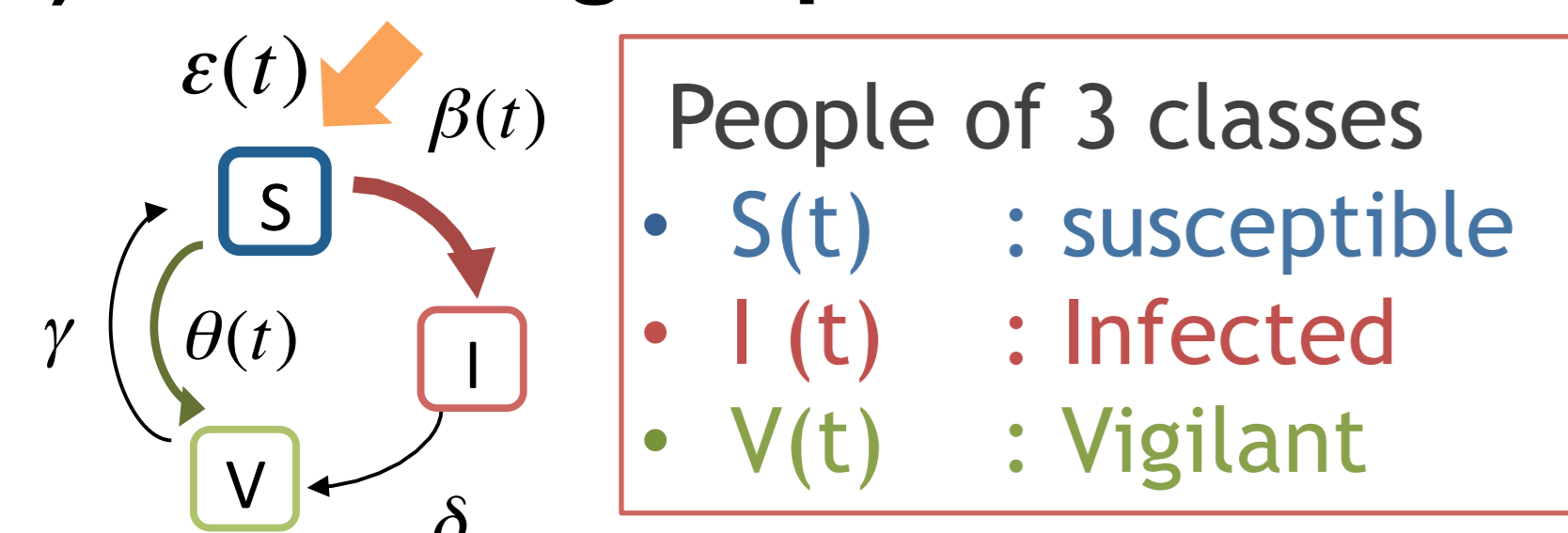states $l$, diseases $d$, $\mathcal{X}$, Time, $n$

Element x: # of cases (weekly)
e.g., 'measles', 'PA', 'April 1-7, 1931', '4740'

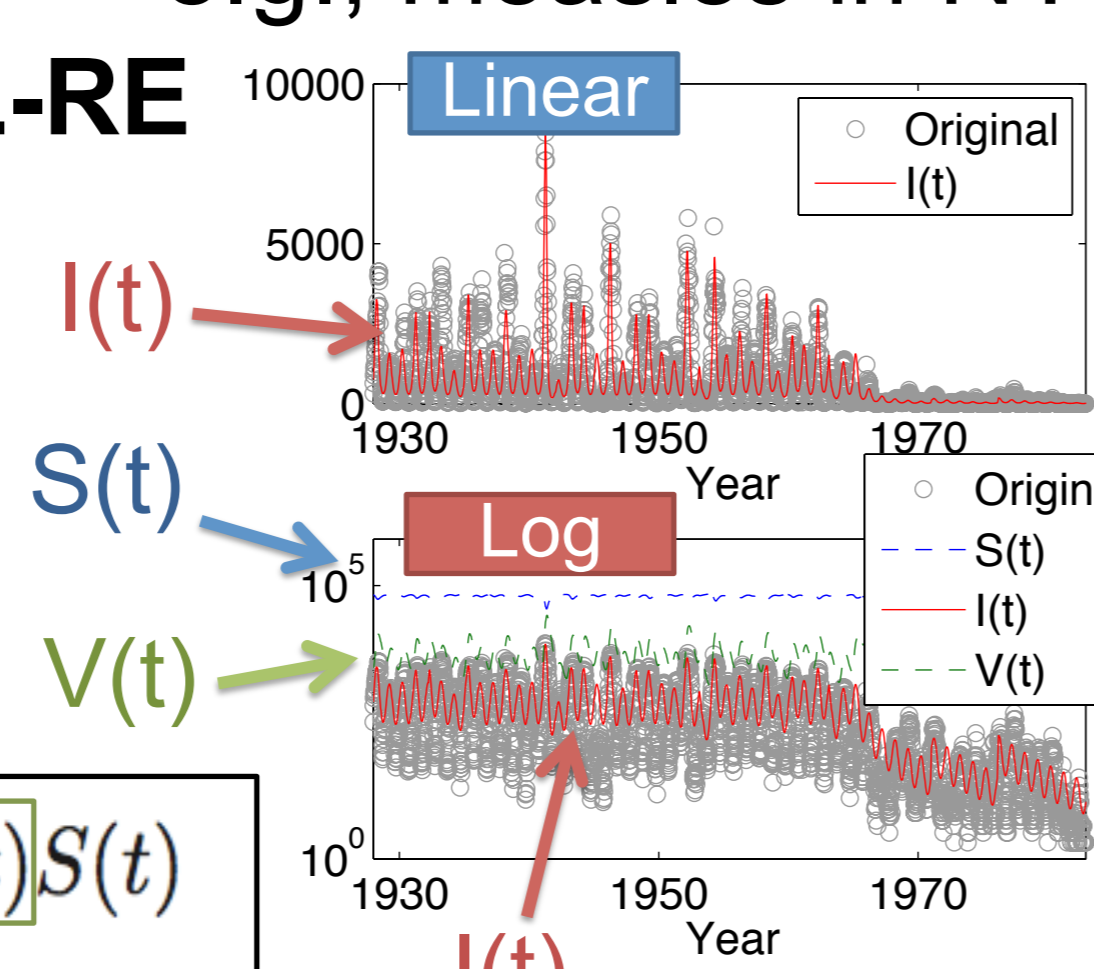## Observations - Properties of real epidemic data

- **P1** yearly periodicity (e.g., flu peaks in the winter)
- **P2** disease reduction effects (e.g., vaccination)
- **P3** area specificity and sensitivity (e.g., correlation)
- **P4** external shock events (e.g., outbreaks)
- **P5** mistakes, incorrect values (e.g., typos)

## Proposed model: FUNNEL

e.g., measles in NY

### (a) With a single epidemic: FUNNEL-RE

$\varepsilon(t)$, $\beta(t)$, S, $\gamma$, $\theta(t)$, I, V, $\delta$

People of 3 classes
- $S(t)$ : susceptible
- $I(t)$ : Infected
- $V(t)$ : Vigilant

Linear — Original, I(t)
Log — Original, S(t), I(t), V(t)

$I(t)$, $S(t)$, $V(t)$

$$S(t+1) = S(t) - \beta(t)\epsilon(t)S(t)I(t) + \gamma V(t) - \theta(t)S(t)$$
$$I(t+1) = I(t) + \beta(t)\epsilon(t)S(t)I(t) - \delta I(t)$$
$$V(t+1) = V(t) + \delta I(t) - \gamma V(t) + \theta(t)S(t) \quad (3)$$

$\beta(t)$ : strength of infection (yearly cycle)
$\delta$ : healing rate   $\gamma$ : forgetting rate
$\theta(t)$ : disease reduction effect
$\varepsilon(t)$ : temporal susceptible rate

$$\beta(t) = \beta_0 \cdot \left(1 + P_a \cdot \cos\left(\frac{2\pi}{P_p}(t + P_s)\right)\right)$$
$$\theta(t) = \begin{cases} 0 & (t < t_\theta) \\ \theta_0 & (t \geq t_\theta) \end{cases} \quad P_p = 52$$

### (b) With multi-evolving epidemics: FUNNEL-full

$d$, $\mathcal{X}$, $l$, $n$

**Global**   **Local**   **Extra**

$N, \beta_0, \delta, \gamma, P_a, P_s$
$\theta_0, t_\theta$

$\mathbf{B}$ (6), $\mathbf{R}$ (2), $\mathbf{N}$, $\mathcal{E}$, $\mathcal{M}$

**P1** Base matrix
**P2** Disease reduction matrix
**P3** Geo-disease matrix
**P4** External shock tensor
**P5** Mistake tensor

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

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

Disease matrix   Time matrix   State matrix

## Optimization algorithm: FUNNEL-Fit

P4, P5, $\mathcal{E}$, $\mathcal{M}$

**Idea (1) Model description cost**

**Q. How can we find externals and mistakes??**

**A. Minimize coding cost!**

Dimensions of $\mathcal{X}$

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

Model description cost of $\mathcal{F}$

Coding cost of $\mathcal{X}$ given $\mathcal{F}$

**Idea (2) Multi-layer optimization**

FUNNEL $\mathcal{F}$

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

**Global**: $\mathbf{B}$ **P1**, $\mathbf{R}$ **P2**, $\mathbf{E}$ **P4**, $\mathbf{M}$ **P5**
**Local**: $\mathbf{N}$ **P3**, $\mathbf{E}$ **P4**, $\mathbf{M}$ **P5**

(step1) Global fitting $\mathcal{X}$
(step2) Local fitting $\mathcal{X}$

## Experiments - (a) Sense-making (5 properties)

**P1**
April (4), May (5), June (6), July (7), August (8), September (9), October (10), November (11), December (12), January (1), February (2), March (3)
Measles, Rubella, Mumps, Chickenpox, Influenza, Streptococcal sore throat, Gonorrhea, Typhoid fever, Cryptosporidiosis, Typhus fever, Lyme disease, Rocky mountain spotted fever, Smallpox

**P2** Measles — Vaccination — Original, I(t); Original, S(t), I(t), V(t)

**P3**

**P4** Smallpox epidemics in 1937–39

**P5** Mistake

### (b) Fitting accuracy (c) Scalability – (linear with data size)

RMSE: SIRS, SKIPS, Funnel-R, Funnel

Diseases — Wall clock time (s), Number of diseases (d)
States — Number of states (l)
Time — Number of time-ticks (n)

## Discussion - Application & Generality

(a) Forecasting (Influenza) — Funnel, AR(52), AR(26), AR(8)

(b) Epidemics on computer networks — Sircam, Badtrans, Netsky, Klez, Mytob

## Conclusions – FUNNEL has following advantages:

- **General & Sense-making**: it captures all essential aspects (P1-P5)
- **Fully-automatic**: it needs no training set
- **Scalable**: it scales linearly with the input size

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

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