# BIG DATA ANALYTICS FOR CYBER SECURITY

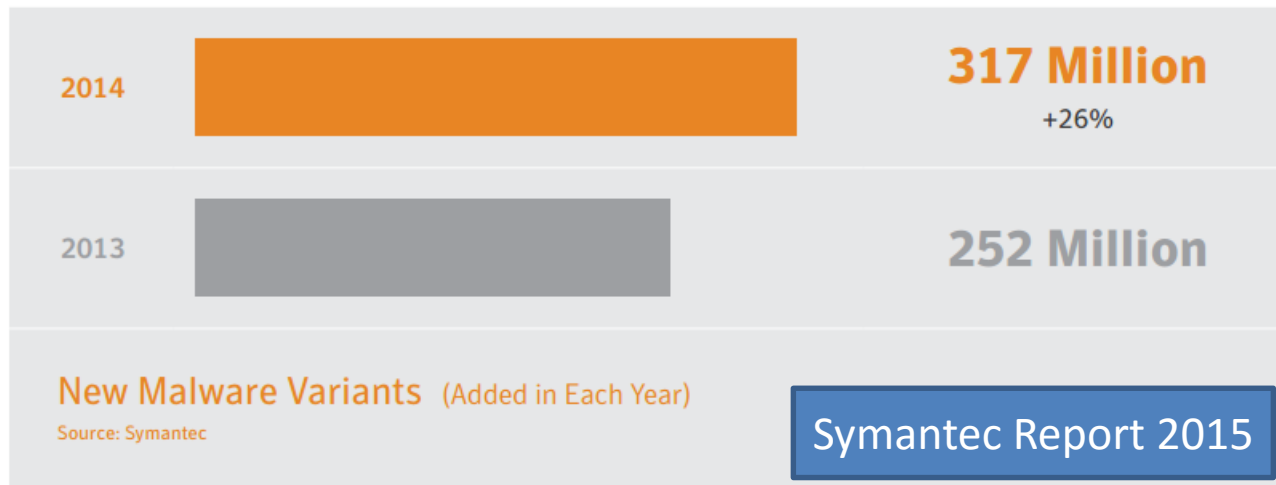**V.S. Subrahmanian**

University of Maryland, College Park

vs@cs.umd.edu, @vssubrah

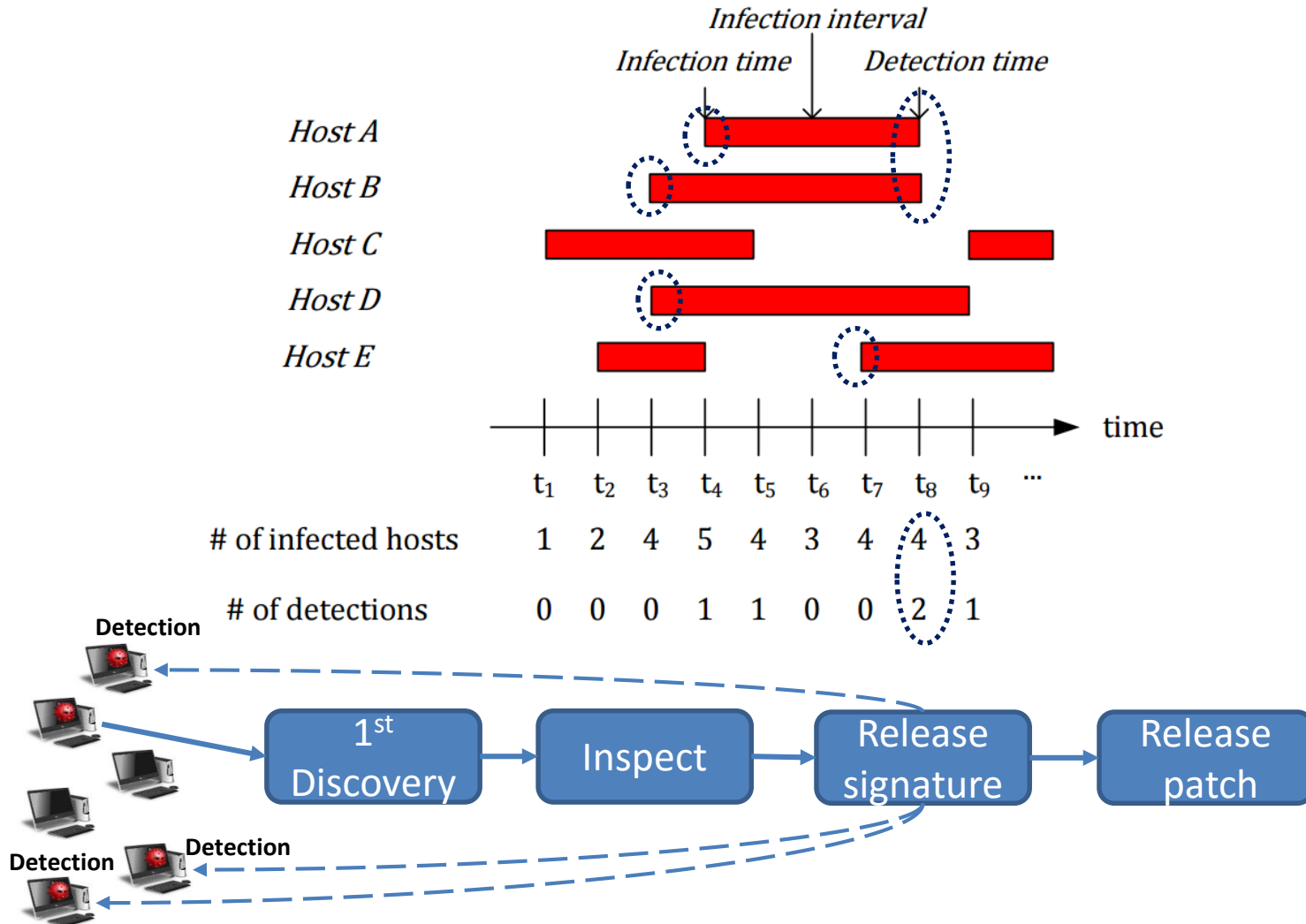Joint work with Chanhyun Kang, Noseong Park, B. Aditya Prakash, Edoardo Serra

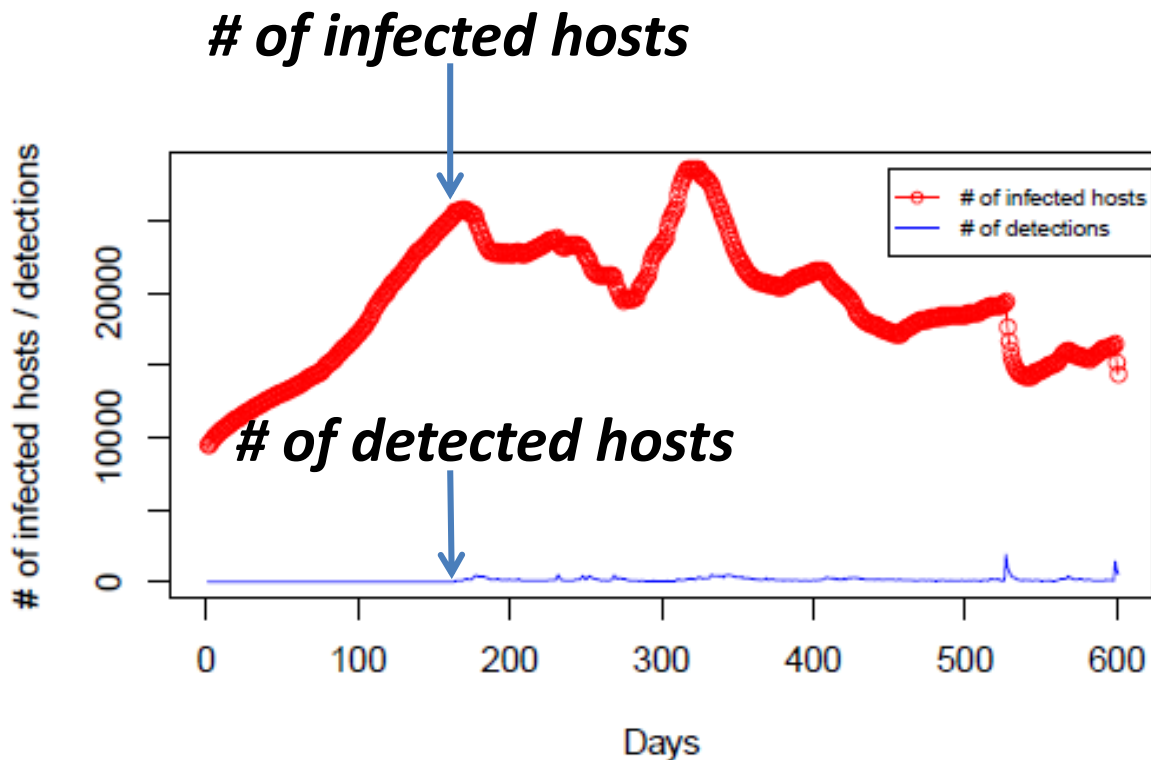UNIVERSITY OF
MARYLAND

# Malware Pandemic

| | | |
|---|---|---|
| **2014** | | **317 Million** +26% |
| 2013 | | **252 Million** |

**New Malware Variants** (Added in Each Year)
Source: Symantec

Symantec Report 2015

# Malware is hard to detect!

# Key Challenge

- Statistics from Symantec WINE Dataset
  - # of Detections <<< # of Infections

**# of infected hosts**

**# of detected hosts**

*Chart legend:*
- # of infected hosts
- # of detections

*Y-axis: # of infected hosts / detections*
*X-axis: Days*

# Problem Statement

# Our Approaches

- Feature based prediction method
  - Proposed a set of novel features

- Epidemic model inspired by SIR model

- Ensemble method that merges the previous two methods with other state-of-the-art techniques.

1st Method

# FEATURE BASED PREDICTION MODEL

# Feature Based Method

Each record= (Host, Malware, File name, Infection time, Detection time)

**Symantec Telemetry data**

**2. Compute host-level features**

'*Detection and Patch incompetence*' of each host
'*Detection and Patch ability*' of each host
'*Detection and Patch hardness*' of each malware
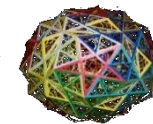
**3. Compute country-level features**

**4. Train a prediction model with the features**

**Prediction model**

*Features /
Infection Ratio*

*The expected number of
infections in future*

**Aggregate host level features, e.g. Average**

| | day | Feature #1 | Feature #2 | …… | Infected Host Ratio |
|---|---|---|---|---|---|
| 80% Training | d | | | | Ground Truth |
| | d+1 | | | | |
| | ⋮ | | | | |
| 20% Test | d+n | | | | Ground Truth vs. Predictions |
| | ⋮ | | | | |

~2 years data

# Detection/Patch Incompetence

- Each record= (Host $h$, Malware $m$, File name $f$, Infection time $i$, Detection time $t$)

- Detection time – Infection time <span style="color:red">(Detection Incompetence[1])</span>
  - How good/bad is a user $h$ at detecting malware $m$?
  - How easy/hard is it to detect malware $m$?

- Patch time – Infection time <span style="color:red">(Patch Incompetence[1])</span>
  - How good/bad is a user at patching a vulnerability/malware?
  - How easy/hard is it to patch a vulnerability/malware?

- Average these values for each host → host-level detection/patch incompetence

- Some other similar features, e.g., Detection time – Malware signature release time

1: These two are the most simplest features.

# Detection Ability/Hardness

- Each record= (Host $h$, Malware $m$, File name $f$, Infection time $i$, Detection time $t$)

- Detection Ability (ADA) of host $h$ is the weighted sum of Detection Hardness (ADH) of malware detected by $h$.

$$ADA(h) \quad = \sum_{(f,m,t) \in dH(h)} w_{12}(h, f, m, t) \cdot ADH(m)$$

A subset of WINE records, where Host = $h$

- Detection Hardness of malware $m$ is the weighted sum of Detection Ability of hosts that detected $m$.

$$ADH(m) \quad = \sum_{(f,h,t) \in dM(m)} w_{21}(m, f, h, t) \cdot ADA(h)$$

# BiFixpoint Algorithm
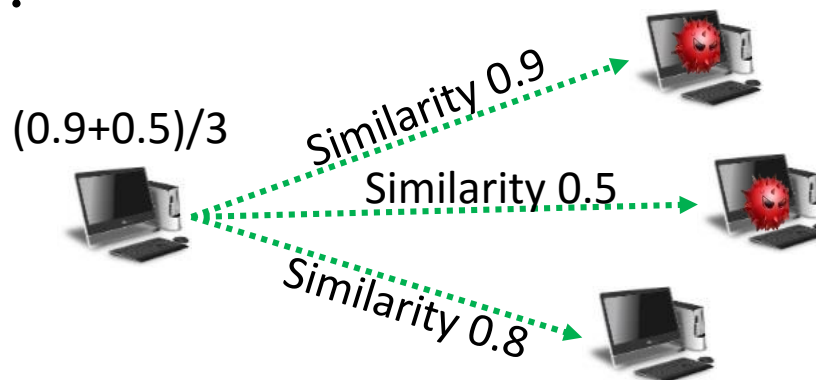
**Algorithm 1: BiFixpoint**

**Input** : $\mathcal{H}, \mathcal{M}, T$ (*$T$ is a training set *)
**Output**: $ADA, ADH$

1  **forall** $h \in \mathcal{H}$, $ADA(h) \leftarrow \frac{1}{|\mathcal{H}|}$ (* initialize *)
2  **forall** $m \in \mathcal{M}$, $ADH(m) \leftarrow \frac{1}{|\mathcal{M}|}$    ***Uniform initialization***
3  change $\leftarrow$ true;
4  **while** *change* **do**
5    $ADA'(h) \leftarrow \sum_{(f,m,t) \in dH(h)} w_{12}(h, f, m, t) * ADH(m)$
6    $ADH(m) \leftarrow \sum_{(f,h,t) \in dM(m)} w_{21}(m, f, h, t) * ADA(h)$
7    **if** $ADA' \sim ADA$ *and* $ADH' \equiv ADH$ **then**     ***Recursive calculation***
8      change $\leftarrow$ false
9    **else**
10     $ADA \leftarrow ADA'$ and $ADH \leftarrow ADH'$
11   **end**
12 **end**
13 **return** $ADA, ADH$

***We prove that convergence is always guaranteed!***

# Collaborative Features

- Given two **similar**[1] hosts $h_1$ and $h_2$
  - Suppose $h_1$ was infected by m.
  - $h_2$ is likely to be infected soon with prob ~ $sim(h_1, h_2)$.
- $cf(h,m)$ is the estimated prob. of host $h$ being infected by $m$ (considering similarity).
- $cf(C,m)$ is the sum of $cf(h,m)$, where $h$ is a host in country $C$.



(0.9+0.5)/3

Similarity 0.9

Similarity 0.5

Similarity 0.8

1: We defined various similarity measures based on calculated features.

# Time Lag Features

- Today's infection ratio depends on not only today's features but also past features.

- Very high dimensional feature space

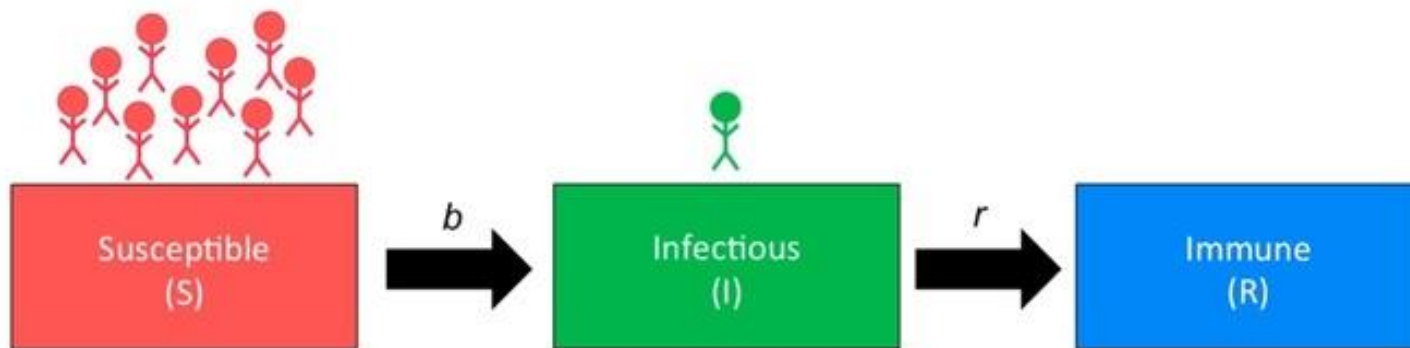| | day | Feature #1 | Feature #1 (-1 day) | Feature #1 (-7 day) | …… | Infected Host Ratio |
|---|---|---|---|---|---|---|
| 80% Training | d | | | | | Ground Truth |
| | d+1 | | | | | |
| | ⋮ | | | | | |
| 20% Test | d+n | | | | | Ground Truth vs. Predictions |
| | ⋮ | | | | | |

# Recap of Features

- Features from raw values
  - Detection time – Infection time (Detection Incompetence)
  - Patch time – Infection time (Patch Incompetence)
  - Some features calculated from raw data
- Features from BiFixpoint Algorithm
  - Detection ability, Patch ability for hosts
  - Detection hardness, Patch hardness for malware
- Collaborative Features
  - Infection numbers based on host similarity
- Country Human Development Index, …
- Time lag features
- Country level aggregation → Regression Problem

2<sup>nd</sup> Method
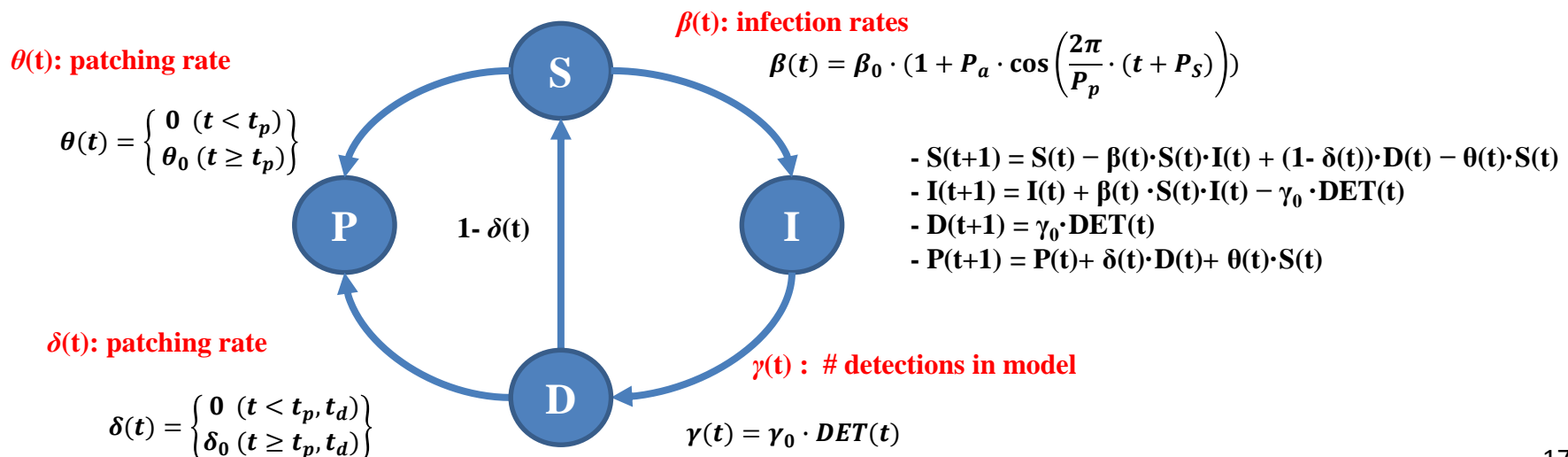
# EPIDEMIC PREDICTION MODEL

# Epidemic Model

- SIR Model models the the dynamics of infectious disease.
- Sometimes used for social rumor diffusion.
- Does not fit the spread of malware.
  - **Recovered** doesn't precisely capture the dynamics of malware spread.
  - Transition rate is not designed for malware.
  - Network data may not always be available.



$b$ = the rate at which susceptible people become infectious
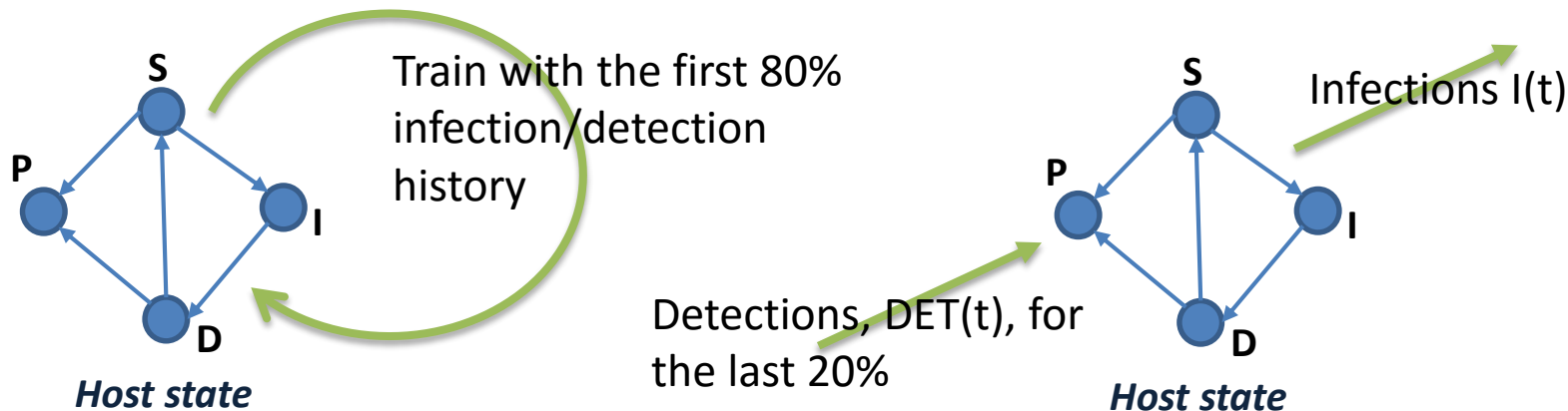$r$ = the rate at which infectious people recover/develop immunity

# DIPS Epidemic Model

- "Recovered" → "Detected" and "Patched"

- Carefully designed transition rates

- S(t), I(t), D(t) and P(t) are the number of susceptible, infected, detected and patched hosts at time t

- S(t), I(t), D(t) and P(t) are recursively defined.

$\theta(t)$: patching rate

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

$\beta(t)$: infection rates

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

S

P  1- $\delta(t)$  I

D

- S(t+1) = S(t) − β(t)·S(t)·I(t) + (1- δ(t))·D(t) − θ(t)·S(t)
- I(t+1) = I(t) + β(t)·S(t)·I(t) − γ₀·DET(t)
- D(t+1) = γ₀·DET(t)
- P(t+1) = P(t)+ δ(t)·D(t)+ θ(t)·S(t)

$\delta(t)$: patching rate

$$\delta(t) = \begin{cases} 0 \ (t < t_p, t_d) \\ \delta_0 \ (t \geq t_p, t_d) \end{cases}$$

$\gamma(t)$ : # detections in model

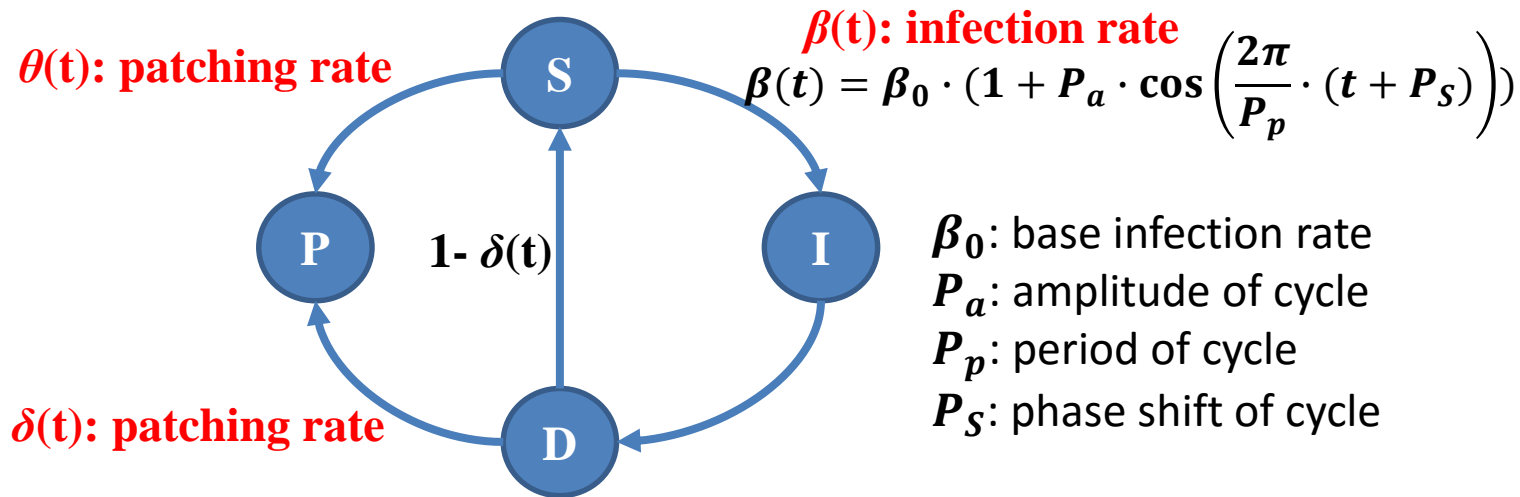$$\gamma(t) = \gamma_0 \cdot DET(t)$$

# How to predict with DIPS

- Find the optimal set of parameters with Least Square Method to minimize the sum of (true-prediction)$^2$

- Train with the target country-malware pair.
  - Initialization → local optimal → not stable learning

- Learning algorithm (two phases)
  - First, train the parameters with all countries and malware
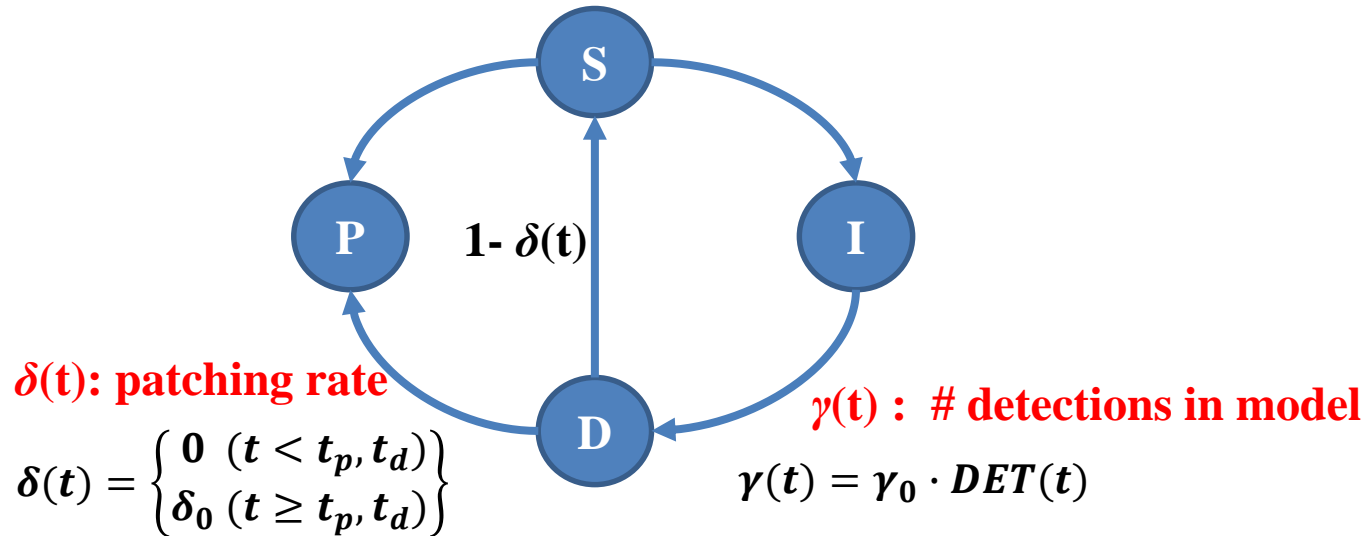  - Second, train again only for the target country-malware

S

P

I

D

Train with the first 80% infection/detection history

*Host state*

S

P

I

D

Infections I(t)

Detections, DET(t), for the last 20%

*Host state*

26

# DIPS - Susceptible



$\theta(t)$: **patching rate**

$\beta(t)$: **infection rate**

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

**1- $\delta(t)$**

$\beta_0$: base infection rate
$P_a$: amplitude of cycle
$P_p$: period of cycle
$P_S$: phase shift of cycle

$\delta(t)$: **patching rate**

- S→I in between $t$ and $t+1$: $\beta(t) \cdot S(t) \cdot I(t)$[1]
- D→S: $(1- \delta(t)) \cdot D(t)$
- S→P: $\theta(t) \cdot S(t)$
- $S(t+1) = S(t) - \beta(t) \cdot S(t) \cdot I(t) - \theta(t) \cdot S(t) + (1- \delta(t)) \cdot D(t)$

1: This is from SIR model.
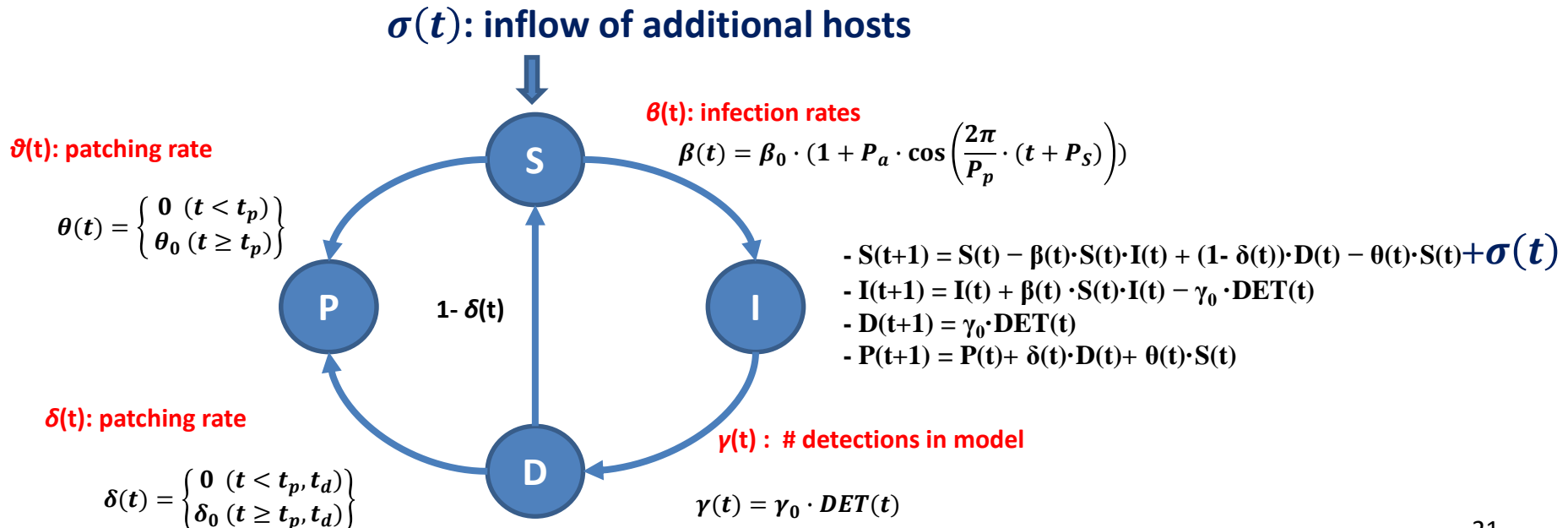
# DIPS - Detected



$\delta(t)$: patching rate

$$\delta(t) = \begin{cases} 0 & (t < t_p, t_d) \\ \delta_0 & (t \geq t_p, t_d) \end{cases}$$

$\gamma(t)$ : # detections in model

$$\gamma(t) = \gamma_0 \cdot DET(t)$$

- I→D: $\gamma_0 \cdot DET(t)$, where $DET(t)$ is the true detection numbers at time $t$
- D→S: $(1 - \delta(t)) \cdot D(t)$
- D→P: $\delta(t) \cdot D(t)$

- $D(t) = \gamma_0 \cdot DET(t)$

# DIPS-exp Epidemic Model

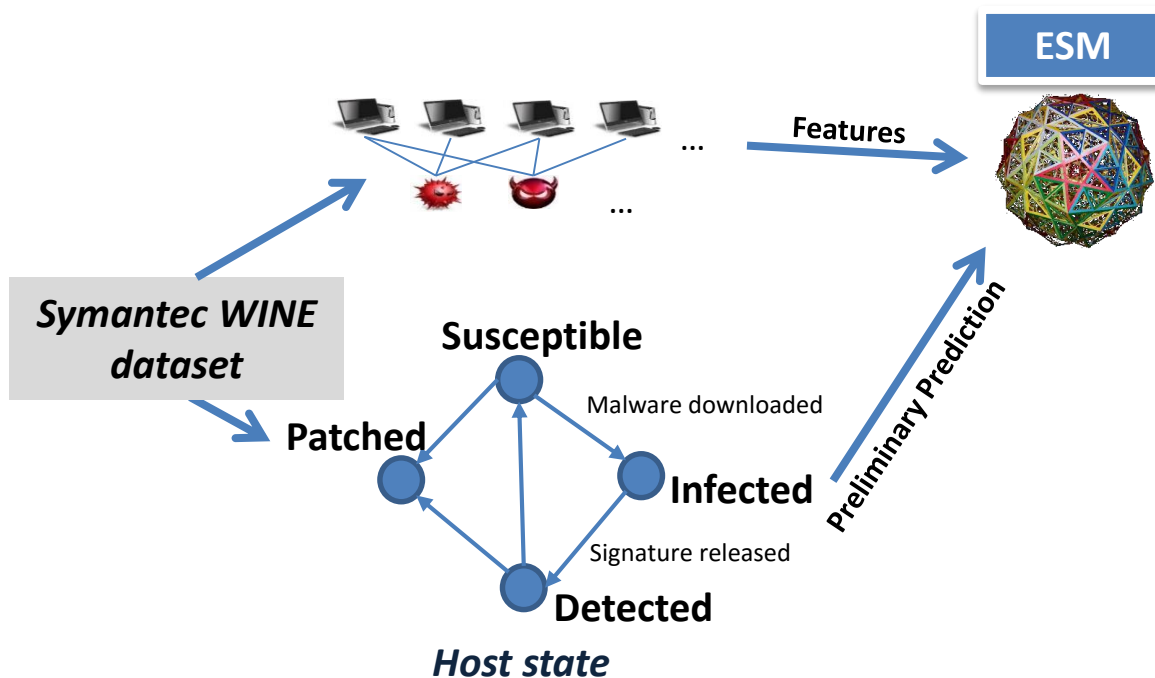- Modeling of "Birth" of the SIR model
- $\sigma(t)$ is added.

$\sigma(t)$: **inflow of additional hosts**

$\vartheta(t)$: patching rate

$\beta(t)$: infection rates

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

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

$1 - \delta(t)$

- $S(t+1) = S(t) - \beta(t) \cdot S(t) \cdot I(t) + (1 - \delta(t)) \cdot D(t) - \theta(t) \cdot S(t) + \sigma(t)$
- $I(t+1) = I(t) + \beta(t) \cdot S(t) \cdot I(t) - \gamma_0 \cdot DET(t)$
- $D(t+1) = \gamma_0 \cdot DET(t)$
- $P(t+1) = P(t) + \delta(t) \cdot D(t) + \theta(t) \cdot S(t)$

$\delta(t)$: patching rate

$\gamma(t)$ : # detections in model

$$\delta(t) = \begin{cases} 0 & (t < t_p, t_d) \\ \delta_0 & (t \geq t_p, t_d) \end{cases}$$

$$\gamma(t) = \gamma_0 \cdot DET(t)$$
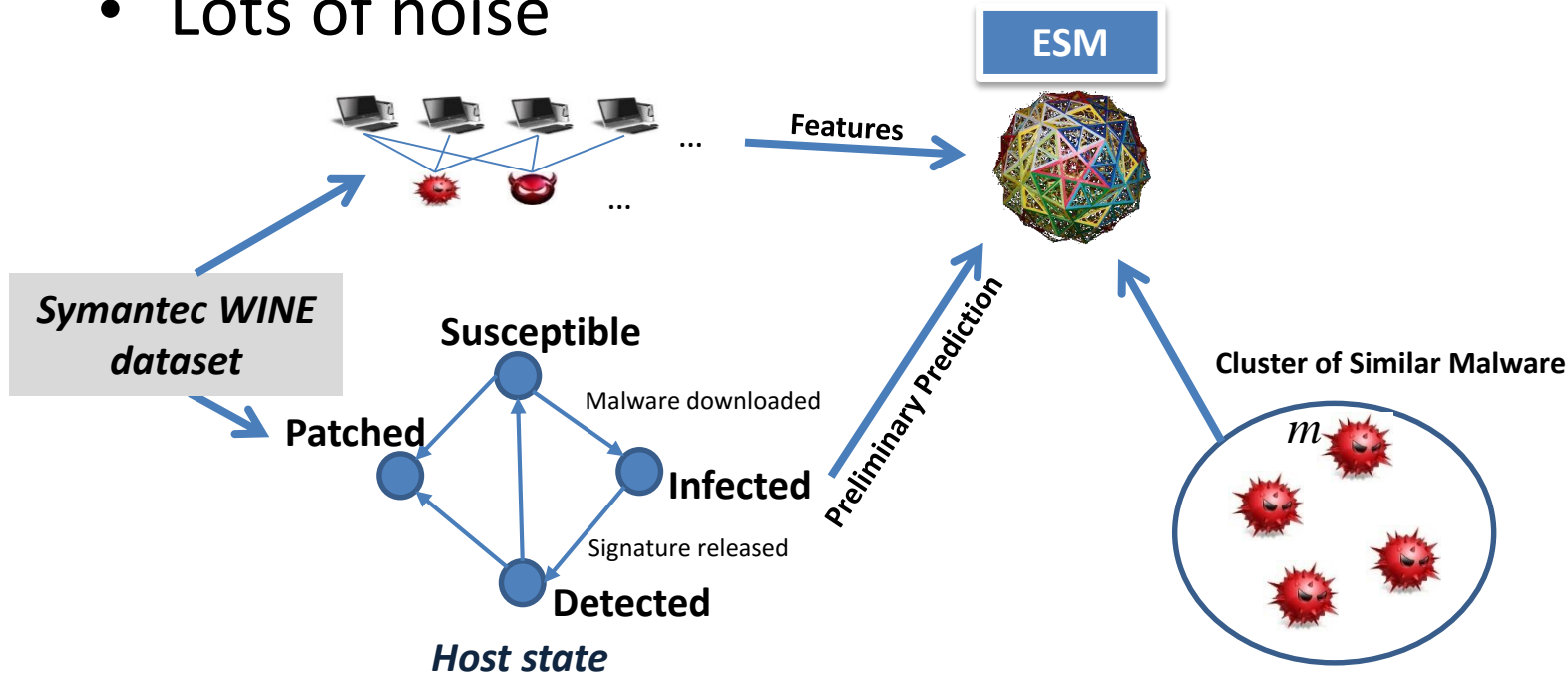
3rd Method

# ENSEMBLE PREDICTION MODEL

# Combine Prediction Models

- Combine Feature Method and DIPS.

- Use DIPS prediction results as additional features.

# Not Enough Training Data

- To predict number of hosts infected by malware $m$, train jointly with similar malware

- Discover similar malware with Dynamic Time Warping to calculate time-series similarity

- Lots of noise

**ESM**

**Features**

**Symantec WINE dataset**

**Susceptible**

Malware downloaded

**Patched**

**Infected**

**Preliminary Prediction**

Signature released

**Detected**

*Host state*

**Cluster of Similar Malware**

$m$

# Robust Regression

- Need a robust regression

- Gaussian Process Regression
  - Very strong Bayesian regression method
  - Less parametric (Parameters are calculated from data with maximum likelihood.)

| Linear Regression | $\hat{y}(w, x) = w_0 + w_1 x_1 + ... + w_p x_p$ |
|---|---|
| Ridge Regression | $\min_{w} \|Xw - y\|_2^2 + \alpha \|w\|_2^2$ |
| Lasso Regression | $\min_{w} \frac{1}{2n_{samples}} \|Xw - y\|_2^2 + \alpha \|w\|_1$ |

Linear combination of weighted features + regularization term

# ESM Model

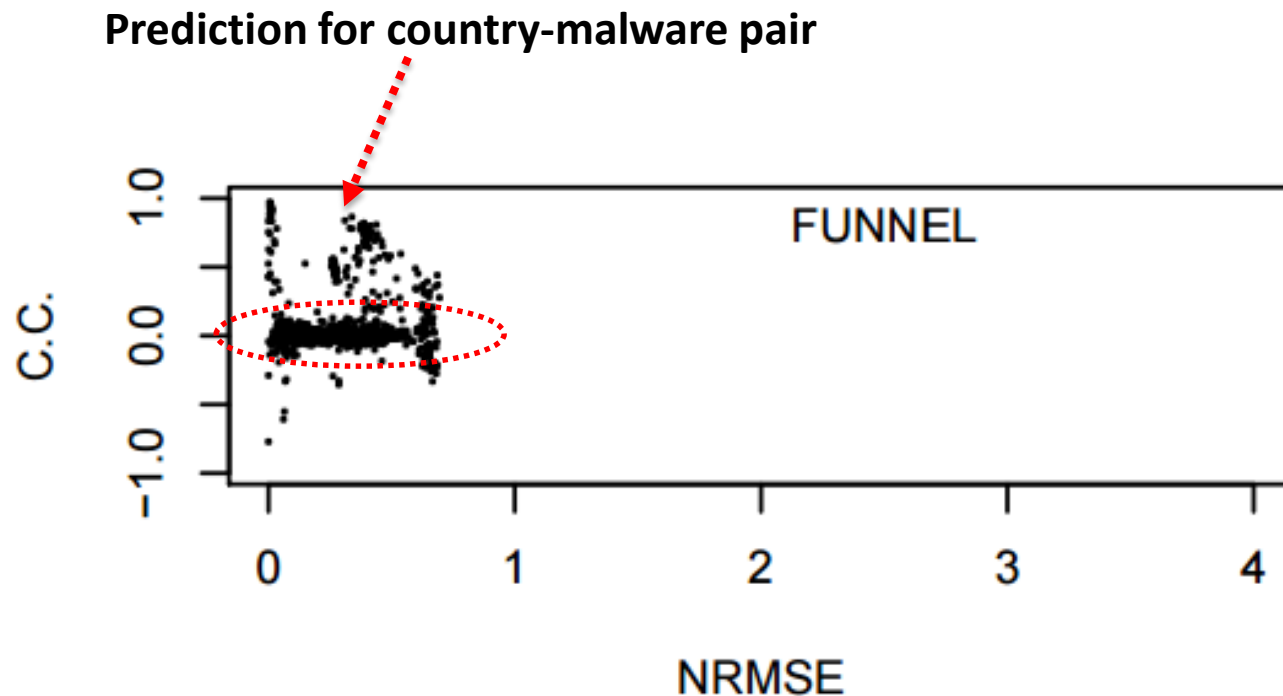| | Feature #1 | Feature #2 | DIPS output | ...... | Infected Host Ratio |
|---|---|---|---|---|---|
| 80% Training (m) | | | | | |
| 80% Training (m1) | | | | | |
| ⋮ | | | | | |
| 20% Test (m) | | | | | |

ESM

Features

GPR

*The expected number of infections in future*

Symantec WINE dataset

Susceptible

Malware downloaded

Patched

Infected

Signature released

Detected

Preliminary Prediction

Cluster of Similar Malware

$m$

*Host state*

26

# Experiment Environment

- Top 50 Most Infectious Malware, Top 40 Country in GDP per capita → 2000 Predictions
- 1.45M unique hosts, 2.99M records
- FBP
- DIPS, DIPS-exp
- FUNNEL: state-of-the-art epidemic model
- ESM0 (FBP + DIPS + DIPS-exp +Similar Malware)
- ESM1 (ESM0 + FUNNEL)

# Measurements

- MAE*=|true infections - predicted infections|

- MSE = (true infection ratio - predicted infection ratio)$^2$

- RMSE = sqrt(MSE)

- NRMSE

- Pearson Correlation Coefficient

# FUNNEL (prior art)

- State of the art epidemic model for human disease

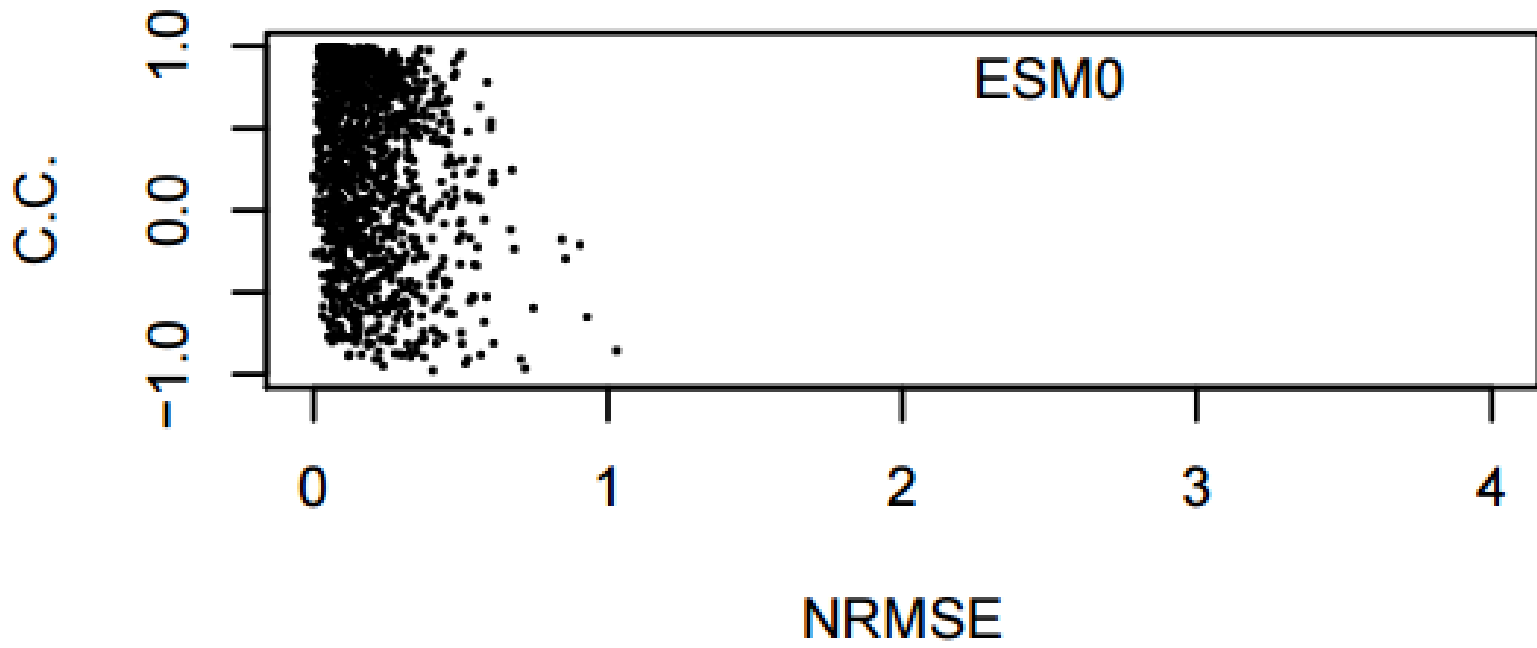- C.C. between truths and predictions are very bad.

**Prediction for country-malware pair**
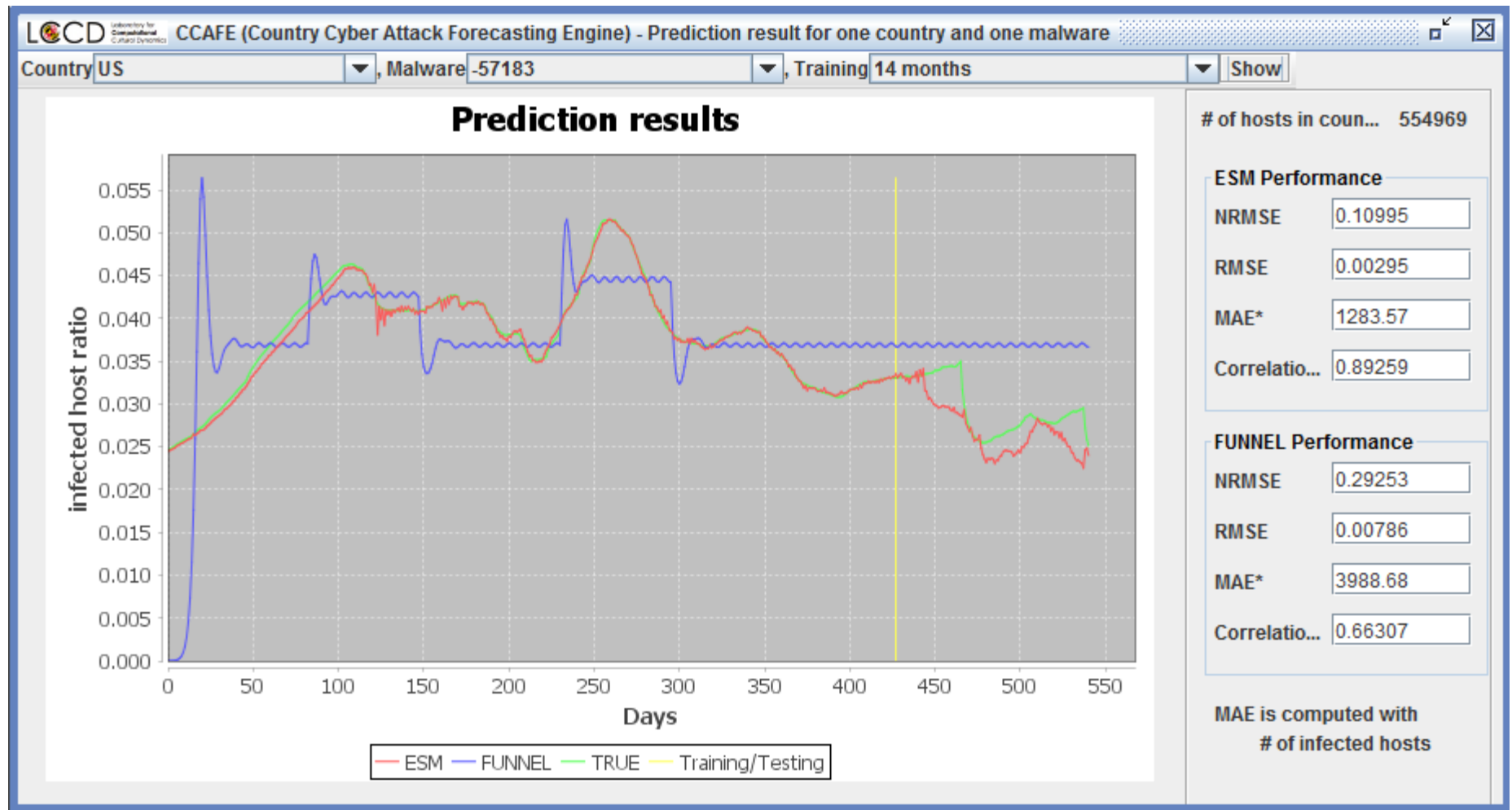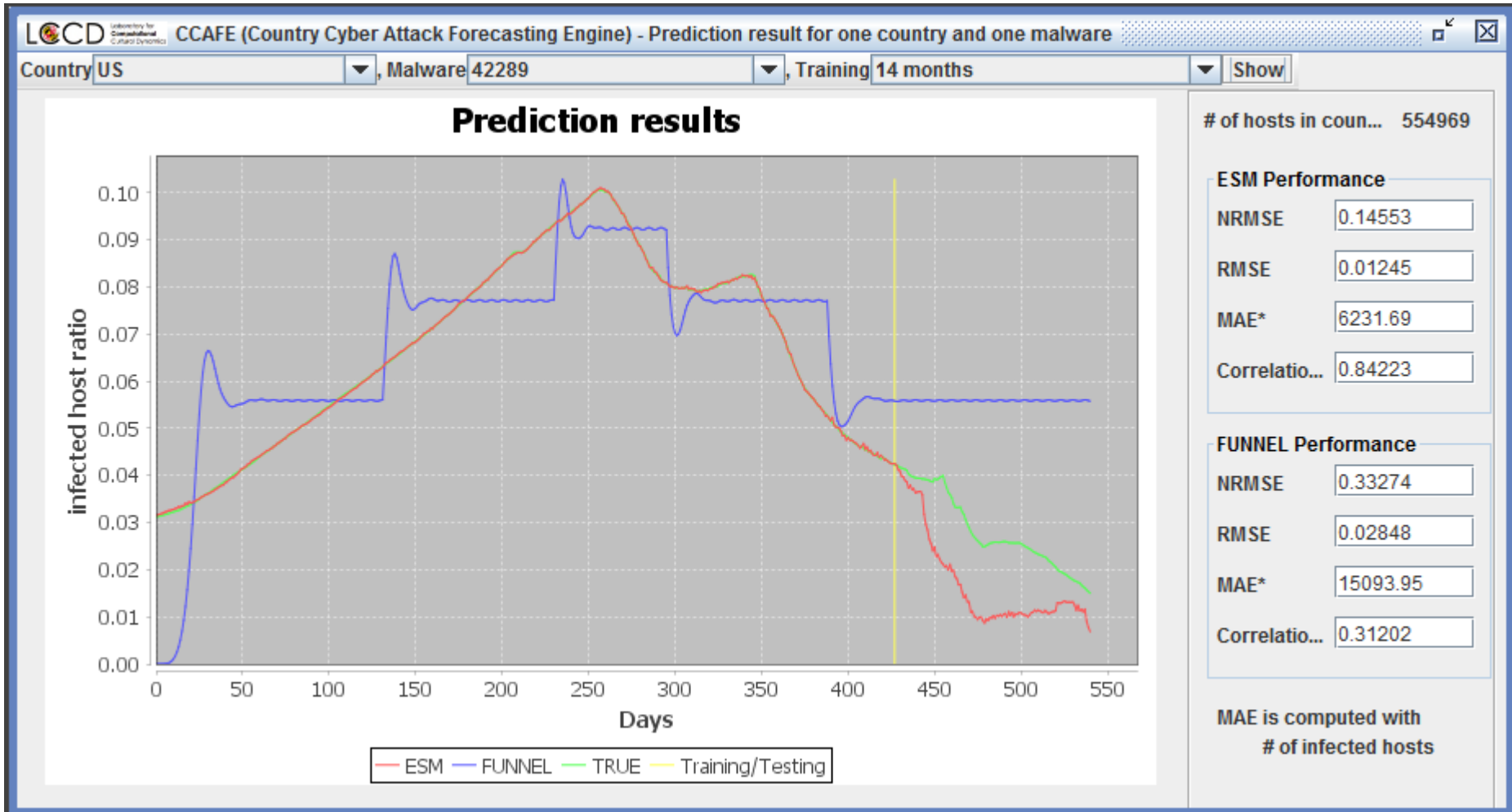
# Feature Based Prediction

# DIPS

# ESM0

# Error Values

| Model | MAE* | RMSE | NRMSE |
|---|---|---|---|
| FBP | 73.74 | 0.00170 | 0.179 |
| FUNNEL | 127.83 | 0.00269 | 0.226 |
| DIPS | 32.36 | 0.00083 | 0.165 |
| DIPS-EXP | 36.56 | 0.00096 | 0.223 |
| $ESM_0$ | 39.41 | 0.00115 | 0.150 |
| $ESM_1$ | 41.84 | 0.00118 | 0.151 |
| $FBP^+_{Funnel}$ | 79.01 | 0.00189 | 0.179 |

# References

C. Kang, N. Park, B.A. Prakash, E. Serra, V.S. Subrahmanian. Ensemble Models for Data-Driven Prediction of Malware Infections, *Proc. 2016 ACM Web Science & Data Mining Conference (WSDM)*, Feb 2016.

# Contact Information

V.S. Subrahmanian

Dept. of Computer Science & UMIACS

University of Maryland

College Park, MD 20742.

vs@cs.umd.edu

@vssubrah

www.cs.umd.edu/~vs/