**Orna Agmon Ben-Yehuda**

**Presents**

# Communication-Efficient Online Detection of Network-Wide Anomalies

Ling Huang*      XuanLong Nguyen*

Minos Garofalakis[§]    Joe Hellerstein*

Michael Jordan*      Anthony Joseph*      Nina Taft[§]

*UC Berkeley      [§]Intel Research

**Coming  on Spring 2011 to a**

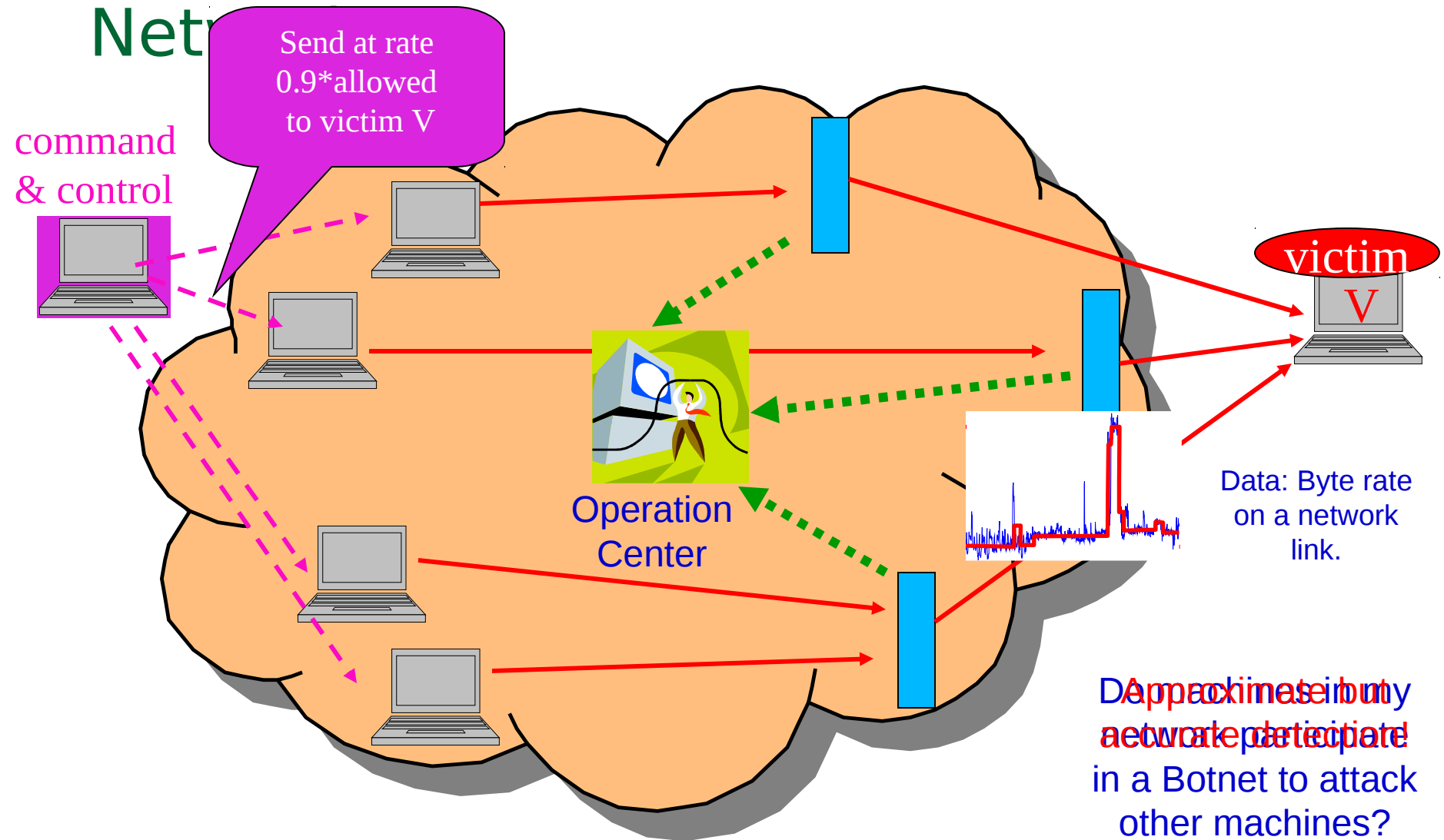**Seminar 236803 on Processing and Mining Distributed Data**

**Near you**

# Network-Wide Anomalies

- Are bad:
  - Router mis-configurations
  - Border Gateway Protocol (BGP) policy modifications
  - Device failures
- Or even malicious:
  - DDOS attacks
  - Viruses, spam sending
  - Port scanning
- But also just unpredictable
  - Flash Crowds (mob) supercomputing

# Detection Problems in Enterprise Network



Send at rate 0.9*allowed to victim V

command & control

Operation Center

victim V

Data: Byte rate on a network link.

Approximate but accurate detection

Data: machine in network participate in a Botnet to attack other machines?

For efficient and scalable detection, push data processing to the edge of network!
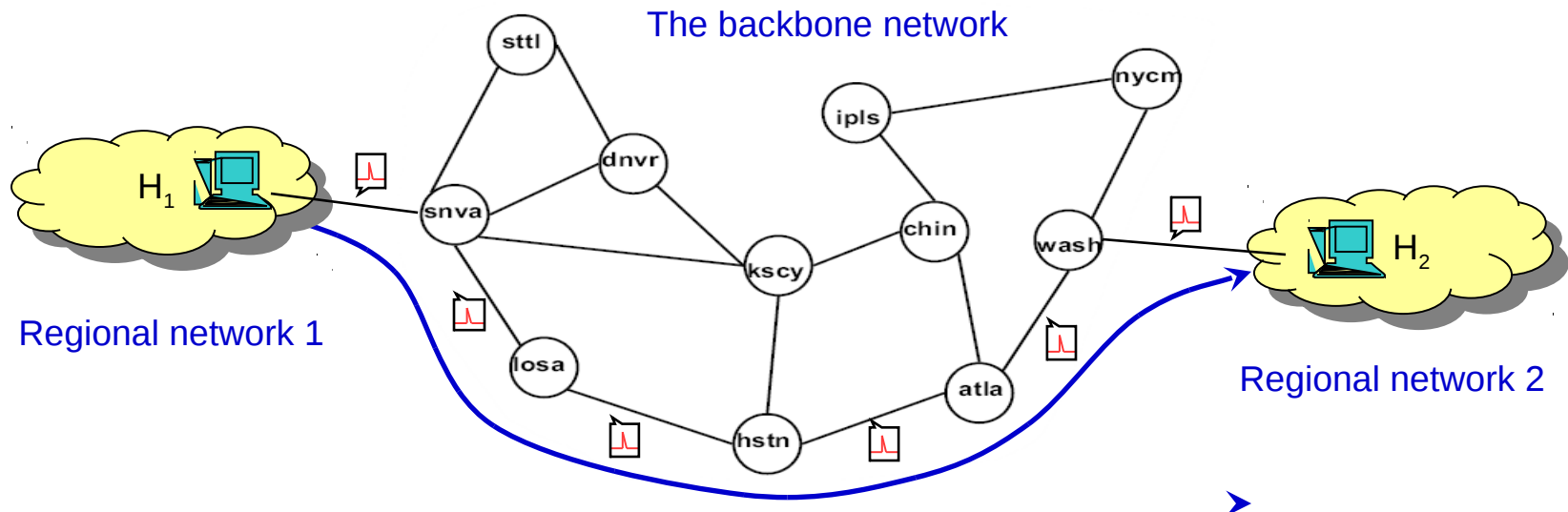
# We shall talk about:

- Lakhina et al.'s centralized algorithm
- Decentralized anomaly detection
- Slack determination
- Evaluation
- Open Discussion

# Towards Decentralized Detection

- Lakhina et al.: Distributed Monitoring & Centralized Computation
  - Stream-based data collection
  - *Periodically* evaluate detection function over collected data
  - Doesn't scale well in network size, timescale, detection delay
- Huang et al.: Decentralized Detection
  - *Continuously* evaluate detection function in a decentr. way
  - Low-overhead, rapid response, accurate and scalable
  - Detection accuracy controllable by a "tuning knob"
    - Provable guarantee on detection error (false alarm rate)
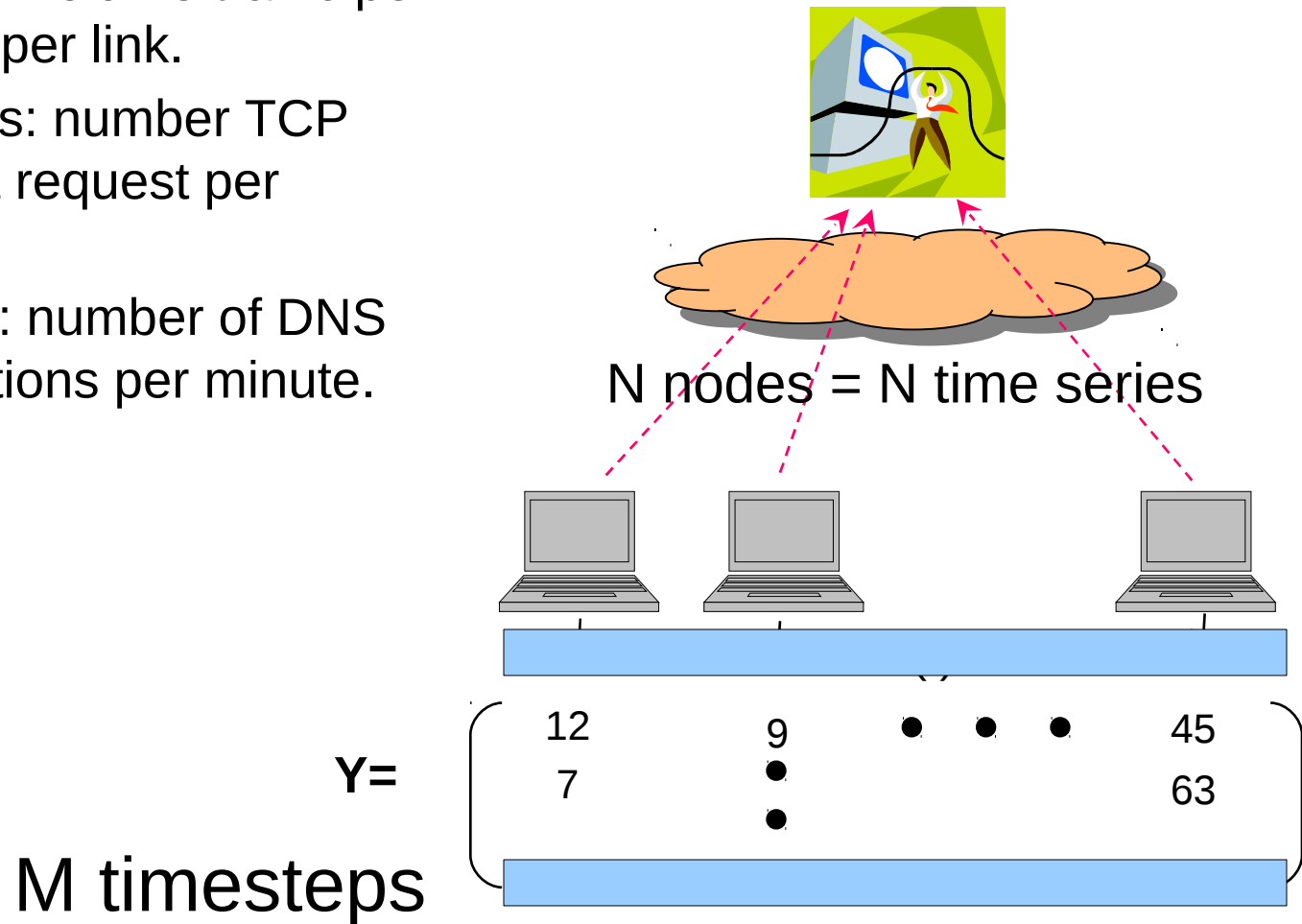    - Flexible tradeoff between overhead and accuracy

# Detection of Network-wide Anomalies

- A ***volume anomaly*** is a sudden change in an <span style="color:red">Origin-Destination flow</span> (*i.e.,* point to point traffic)

- Given <span style="color:red">link</span> traffic measurements, ***detect*** the volume anomalies

The backbone network
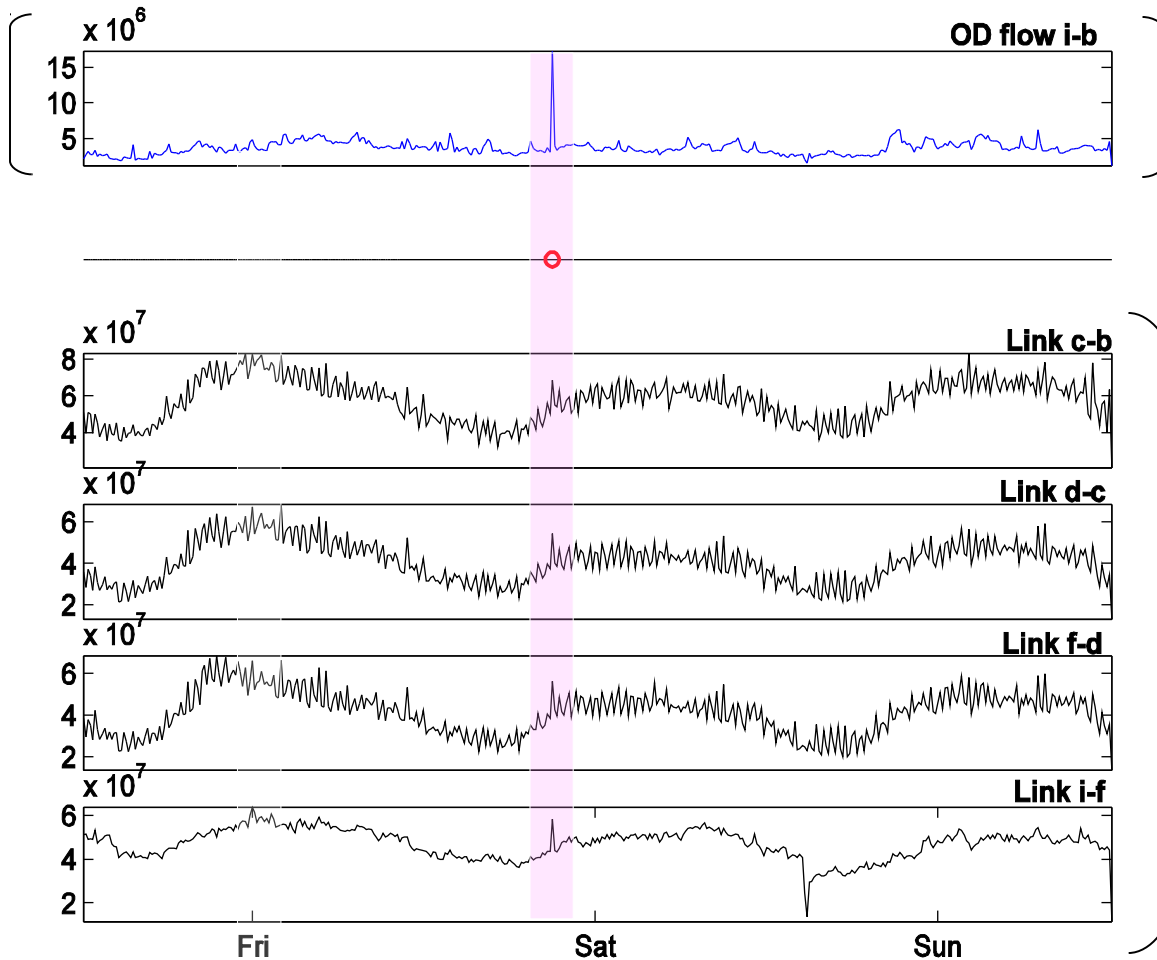
Regional network 1

Regional network 2

# The Data Collected by Monitors

- Routers: volume traffic per second per link.

- Firewalls: number TCP connect request per second.

- Servers: number of DNS transactions per minute.

N nodes = N time series

$$Y= \begin{bmatrix} 12 & 9 & \bullet\ \bullet\ \bullet & 45 \\ 7 & \bullet & & 63 \\ & \bullet & & \end{bmatrix}$$
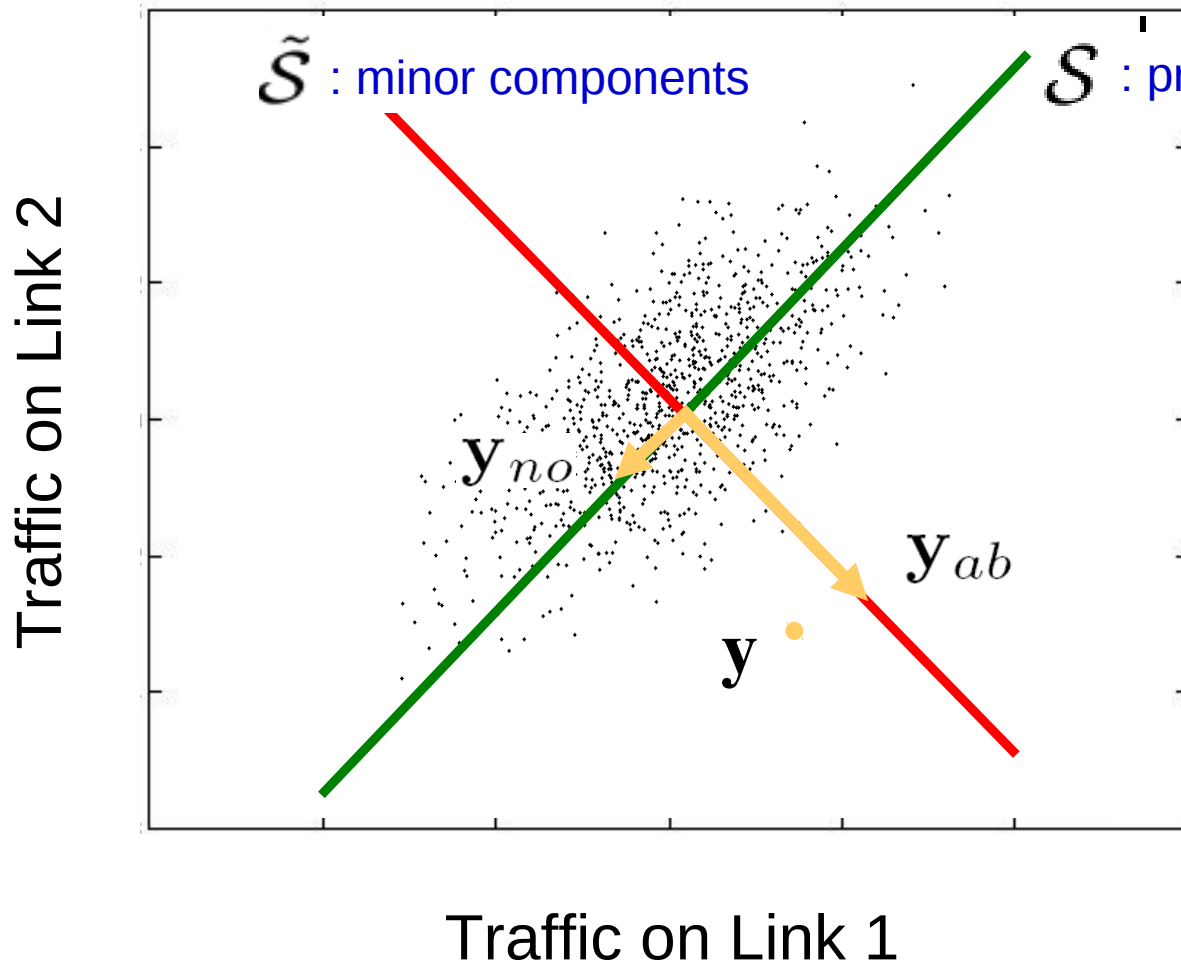
M timesteps

# Flow vs. Link (Lakhina et al.)



Anomalies in (unobserved) flow data

Observed network link data = aggregate of application-level flows
Each link is a dimension

Finding anomalies in high-dimensional, noisy data is difficult!

# Principal Component Analysis (PCA)



$\tilde{\mathcal{S}}$ : minor components    $\mathcal{S}$ : principal components

Principal components are top eigenvectors of covariance matrix.

$$Y\,Y^{T}$$

They are also directions of maximal variance.
They form the subspace projection matrices $C_{no}$ and $C_{ab}$

$$\mathbf{y}_{no} = \mathbf{C}_{no}\mathbf{y}$$
$$\mathbf{y}_{ab} = \mathbf{C}_{ab}\mathbf{y}$$

Traffic on Link 2

Traffic on Link 1

Anomalous traffic usually results in a large value of $\mathbf{y}_{ab}$
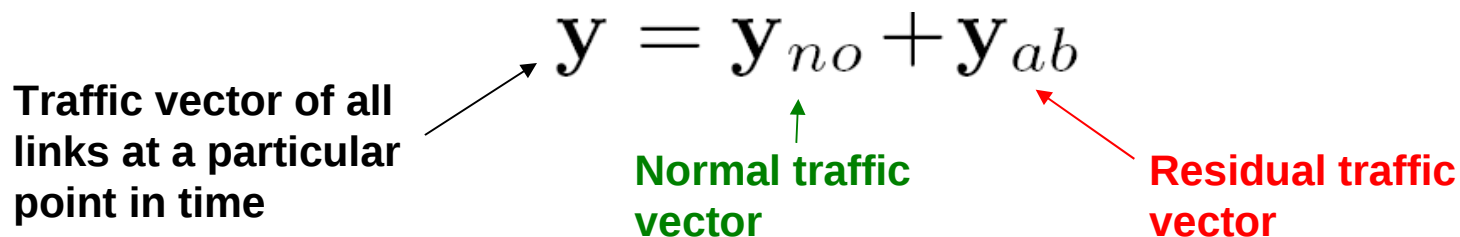
# The Subspace Method (Lakhina'04)

- An approach to separate normal from anomalous traffic based on Principal Component Analysis (PCA)
- Normal Subspace $\mathcal{S}$:  space spanned by the top $k$ principal components
- Anomalous Subspace $\tilde{\mathcal{S}}$:  space spanned by the remaining components
- Then, decompose traffic on all links by projecting onto $\mathcal{S}$ and $\tilde{\mathcal{S}}$ to obtain:

$$\mathbf{y} = \mathbf{y}_{no} + \mathbf{y}_{ab}$$

**Traffic vector of all links at a particular point in time**

**Normal traffic vector**

**Residual traffic vector**

# Link Traffic Variance of Principle Components

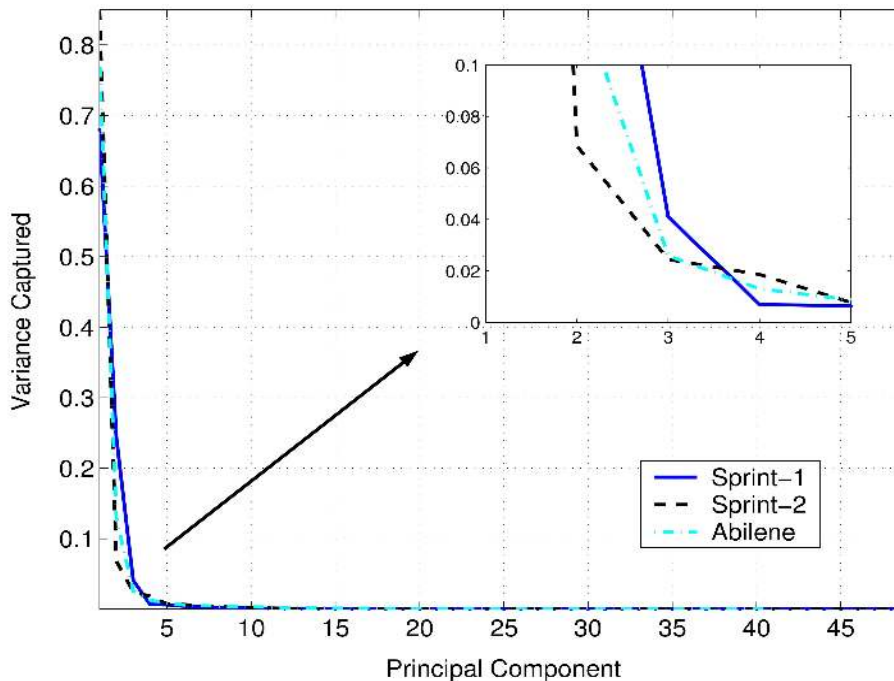- **Link matrices have low dimensionality**



Figure 2: Fraction of total link traffic variance captured by each principal component.

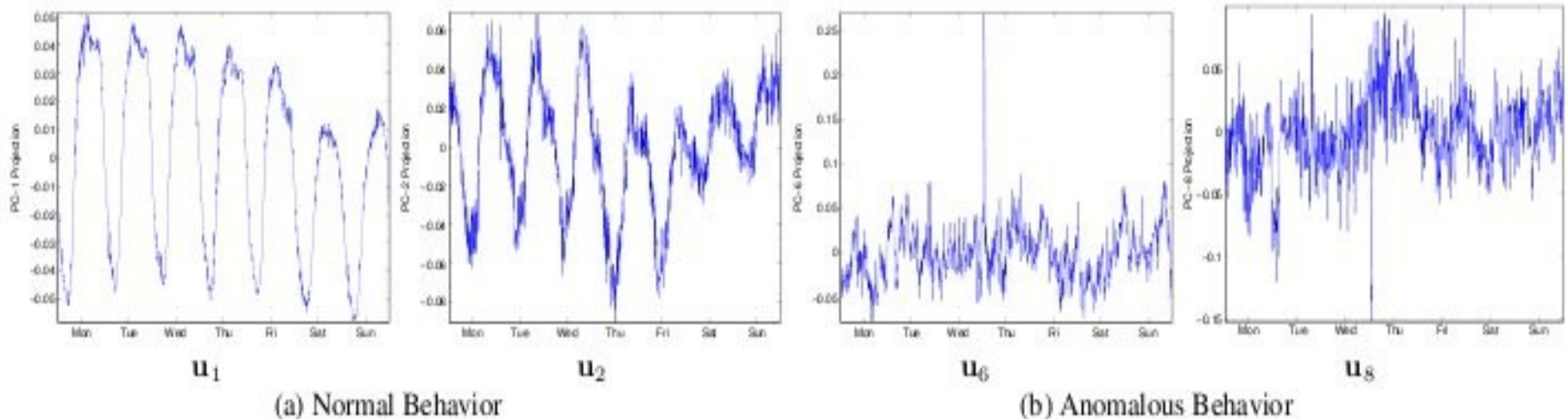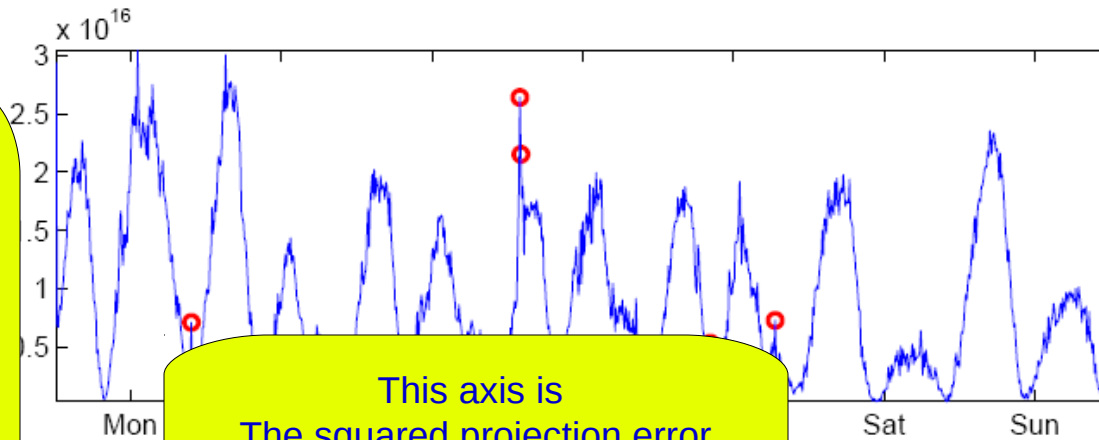# Projections onto Principle Components – normal and abnormal traffic variation



Figure 3: Projections onto principal components showing normal and anomalous traffic variation.
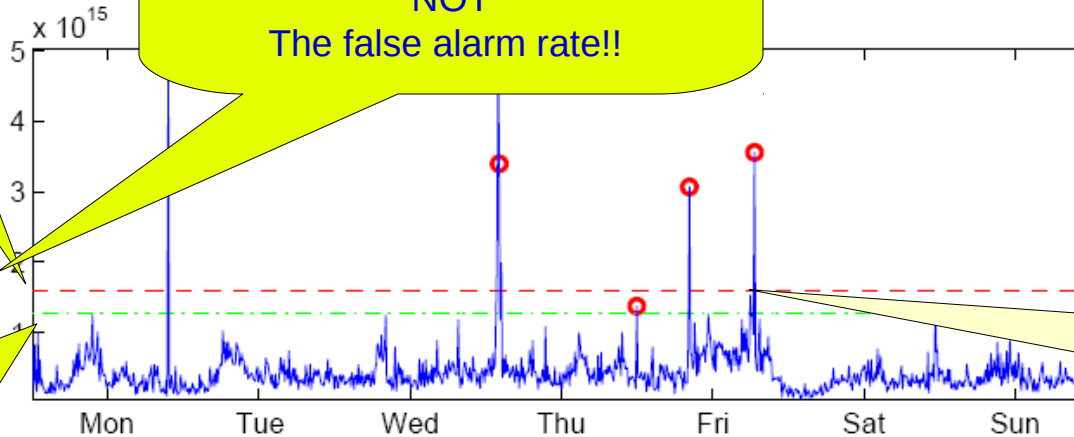
# Detection Illustration



Value of $\|\mathbf{y}\|^2$ over time **(all traffic)**

Value of $\|\mathbf{C}_{ab}\mathbf{y}\|^2$ over time

1-α=99.9
99.9% of Alarms Are Real, But More Anomalies Go undetected

This axis is The squared projection error NOT The false alarm rate!!

$Q_\alpha$

1-α=99.5
Only 99.5% of Alarms Are real But many Anomalies Are detected

This small spike is not an anomaly we wished to detect

Red dots: anomalies     Blue curve: traffic data

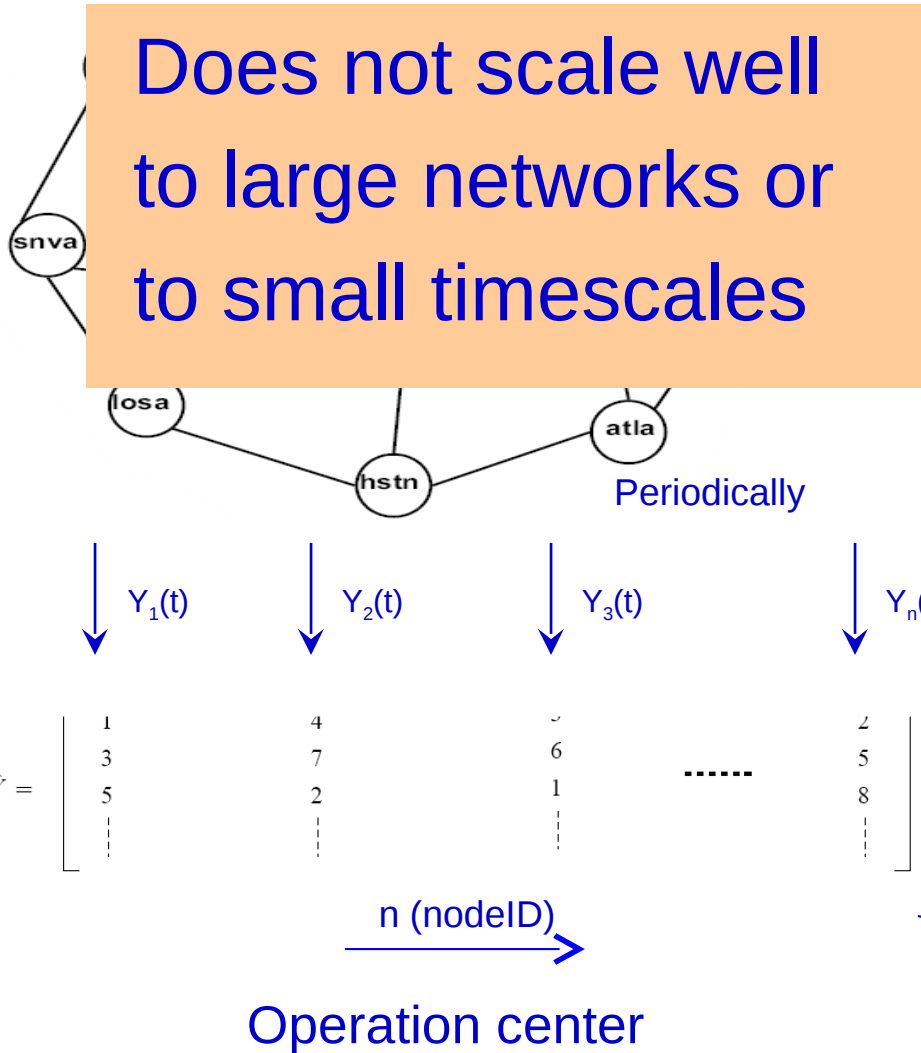Lakhina et al.,     Huang et al., presented by Agmon Ben-Yehuda     13

# Detection Threshold

$$\|C_{ab}\, y\|^2 > Q_\alpha$$

- $Q_\alpha$ is a threshold on the Squared Projection Error (SPE). It guarantees a false alarm rate of less than α.
- Jackson & Mudholkar: computed threshold based on the abnormal eigenvalues of the covariance matrix.
  - No matter where the distinction is made (how many components are considered normal).
  - No matter what the mean amount of traffic is.
  - For multivariate Gaussian distribution only.
- Jensen & Solomon: In practice, holds for different distributions.
- Lakhina et al. Believe traffic is multivariate Gaussian.
  - but have not verified this.

# The Centralized Algorithm

Does not scale well to large networks or to small timescales

- Data matrix *Dat*

    1) Each link produces a column of m data over time.

    2) n links produce a row data y at each time instance.

Detection by Squared Prediction Error (SPE):

$$\|C_{ab}\, y\|^2 > Q_\alpha$$

Projection $C_{ab}$     Threshold $Q_\alpha$
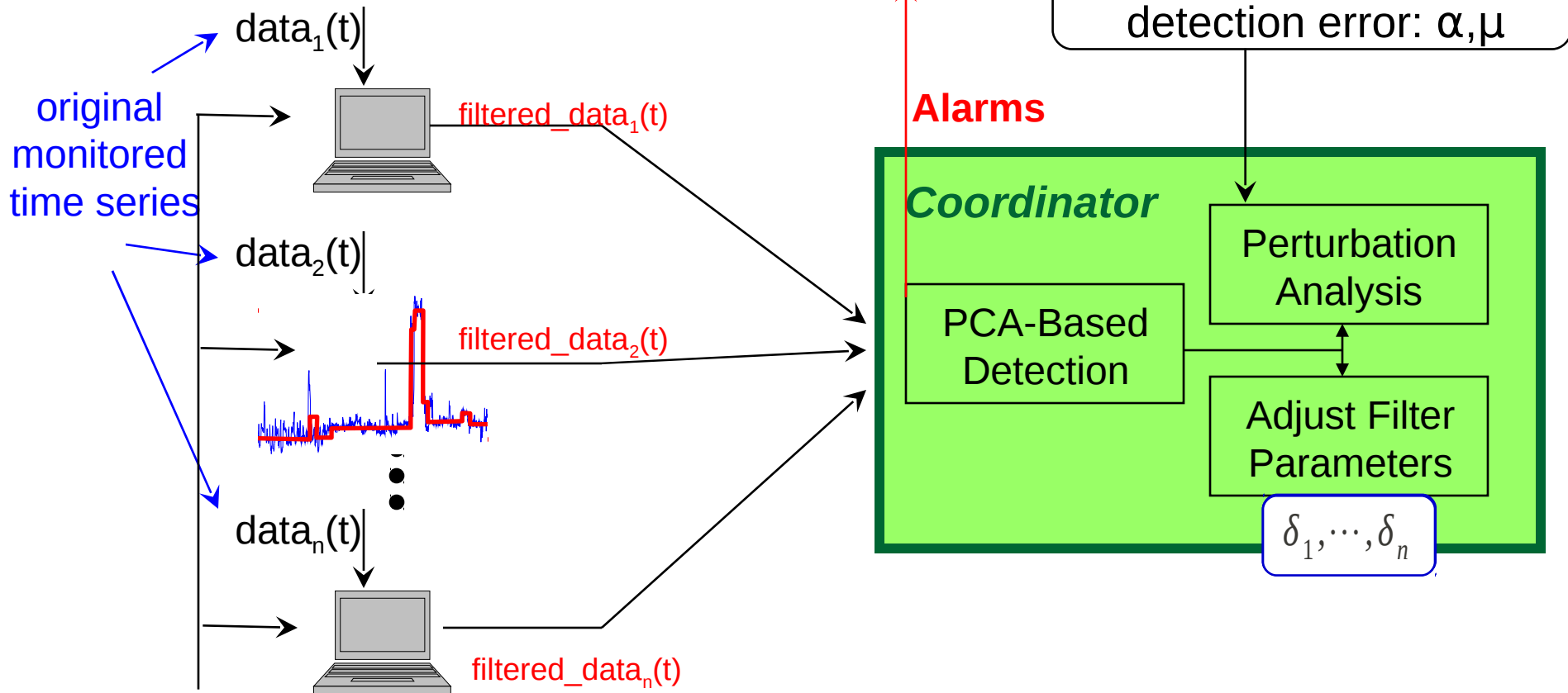
Periodically

Y₁(t)   Y₂(t)   Y₃(t)   Yₙ(t)

$Y = \hat{Y} =$

| 1 | 4 | 6 | 2 |
| 3 | 7 | 1 | 5 |
| 5 | 2 | | 8 |

......

m (timestep)

Eigen vectors     Eigen values

n (nodeID)

PCA on $Y$

Operation center

# Huang et al.: In-Network Detection Framework

**Distr. Monitors**

$data_1(t)$

original monitored time series

filtered_data$_1$(t)

$data_2(t)$

filtered_data$_2$(t)

$data_n(t)$

filtered_data$_n$(t)

**Alarms**

user input: required detection error: α,μ

*Coordinator*

PCA-Based Detection

Perturbation Analysis

Adjust Filter Parameters

$\delta_1, \cdots, \delta_n$

# The Communication and Error Tradeoff

Approximate Info. $\longleftarrow$ PCA on $\hat{\hat{Y}}$ $\longleftarrow$ $\hat{Y}$

$$\|\hat{C}_{ab}\,\hat{y}\|^2 > \hat{Q}_\alpha$$
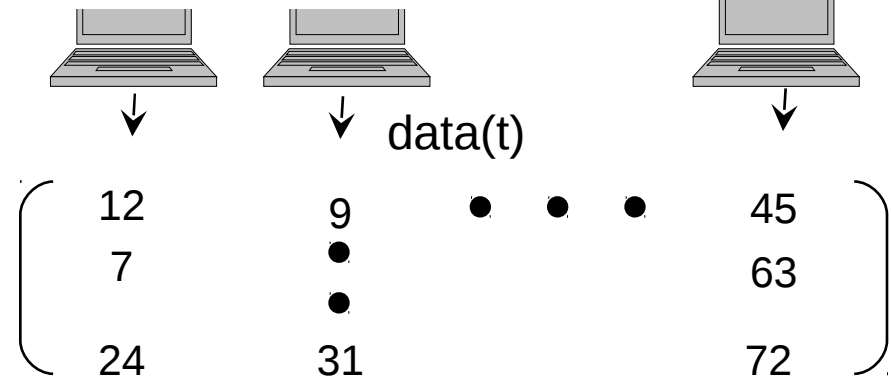
The bigger the filtering parameter $\delta_i$,

the less the communication overhead,

but the more the detection error!

$$\|C_{ab}\,y\|^2 > Q_\alpha$$

PCA on $Y$ $\longleftarrow$ **Y=**

:ered_data(t)

data(t)

$$\begin{pmatrix} 12 & 9 & \bullet\ \bullet\ \bullet & 45 \\ 7 & \bullet & & 63 \\ & \bullet & & \\ 24 & 31 & & 72 \end{pmatrix}$$

The coordinator computes a set of good $\delta_1, \ldots, \delta_n$ to manage this difference.

# The Protocol At Monitors

- Monitor i updates information if
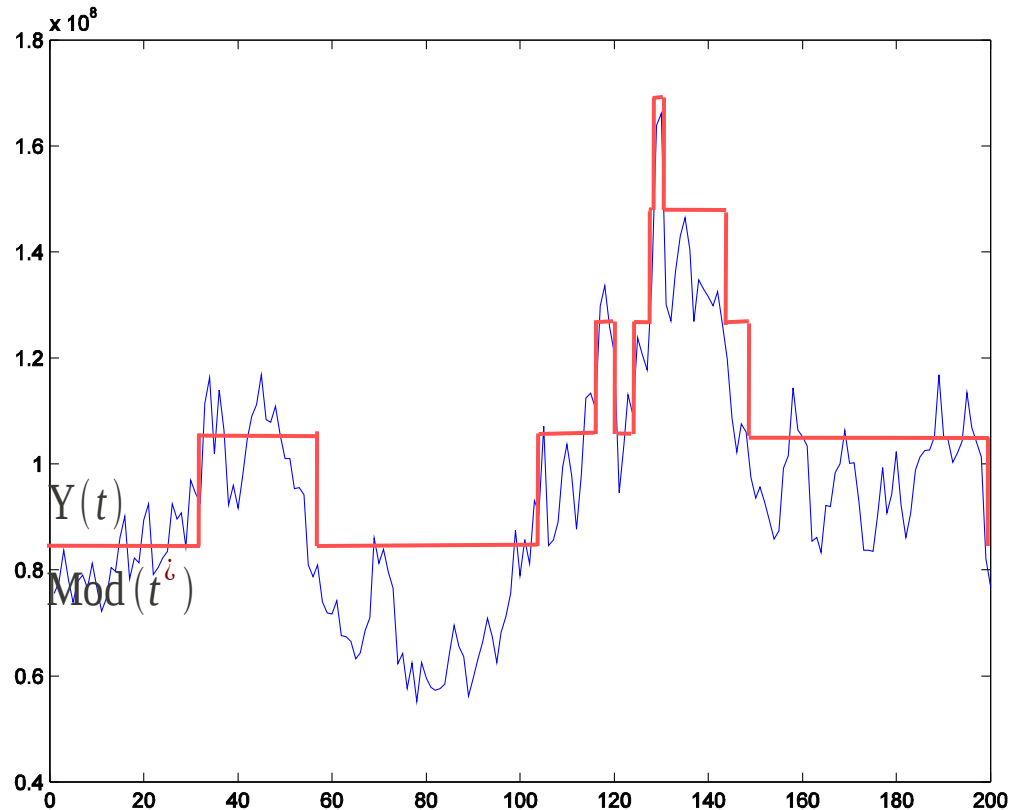
$$\left|Y_i(t) - \mathrm{Mod}_i(t^{\iota})\right| > \delta_i$$

$\delta_1, \cdots, \delta_n$ are the *filtering parameters*

- $\mathrm{Mod}_i(t^{\iota})$ can be based on any prediction **mod**el built on historical data.

  - The prediction model is known to both monitor and coordinator.

  - For example, the average of last 5 communicated signal values.
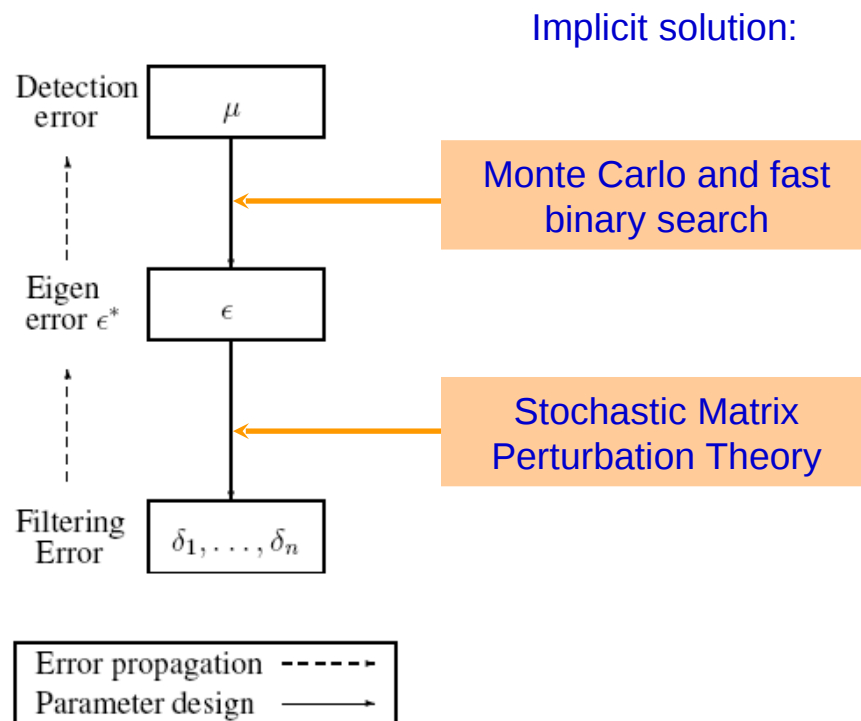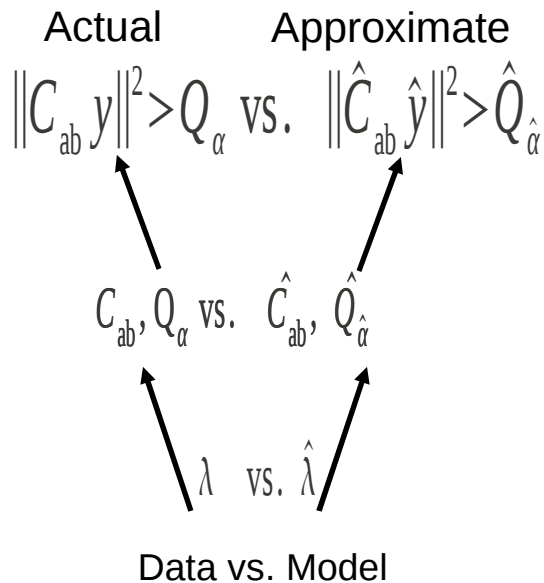
# The Protocol At Monitors



◻ Simple but enough to achieve 10x data reduction

# The Protocol at the Coordinator

- Create new time data from communication and predictions

- Update (cyclic) matrix: add new data, lose oldest

- Re-compute PCA (residual projection matrix, threshold)

- Detect anomalies, fire warnings

- Update slacks when needed (no details...)

# Parameter Design and Error Control

- Users specify an upper bound on false alarm rate, then we determine the filtering parameters $\delta$'s

Actual            Approximate

$$\|C_{ab}\, y\|^2 > Q_\alpha \quad vs. \quad \|\hat{C}_{ab}\, \hat{y}\|^2 > \hat{Q}_{\hat{\alpha}}$$

$$C_{ab}, Q_\alpha \quad vs. \quad \hat{C}_{ab}, \hat{Q}_{\hat{\alpha}}$$

$$\lambda \quad vs. \quad \hat{\lambda}$$

Data vs. Model

Implicit solution:



Detection error $\mu$

Eigen error $\epsilon^*$ $\epsilon$

Filtering Error $\delta_1, \ldots, \delta_n$

Error propagation  ------
Parameter design  ——→

Monte Carlo and fast binary search

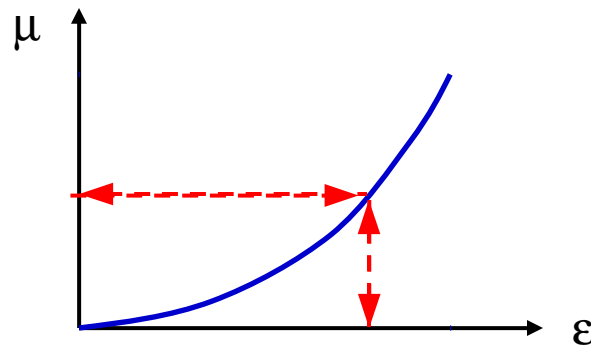Stochastic Matrix Perturbation Theory

Eigen error: $L_2$ norm of the difference between

the approximate eigenvalues and the actual ones

# Parameter Design and Error Control (II)

- ## Detection Error μ → Eigen-Error ε

  - ❑ Monte Carlo simulation to find the mapping from ε to μ

    

  - ❑ For the given μ, a fast binary search to find an ε

# From Eigen-Error to detection Deviation

$$\Pr\left[\|C_a y\|^2 > Q_\alpha\right] = \Pr\left[X > c_\alpha\right] = \alpha,$$
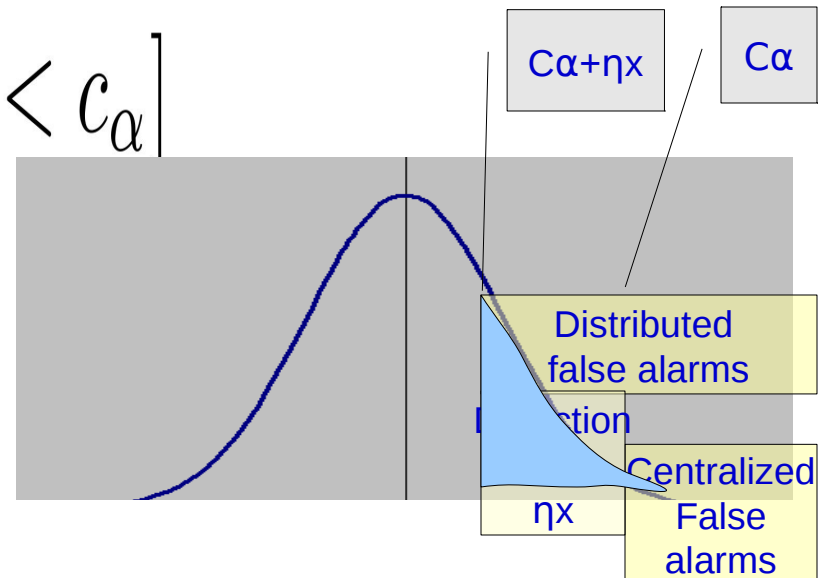
Normalized form of
$$\|C_{ab}\, y\|^2$$
(Jensen & Solomon)

(1-α)-th  percentile

$$\mu = \Pr\left[c_\alpha - \eta_X < N(0,1) < c_\alpha\right]$$

Cα+ηx

Cα

Upper bound on
$$\left|\hat{X} - X\right|$$
Estimated using max of
Monte Carlo results

Distributed
false alarms

...ction

Centralized
False
alarms

ηx

# Parameter Design and Error Control (III)

## Eigen-Error $\varepsilon$ $\rightarrow$ Filtering parameters $\delta$s

- **Error Matrix:** $W = Y - \hat{Y}$

- **Elements of column vector $W_i$ bound by $\delta_i$**

- **Assumptions:**
  - $W_i$ are independent, radially symmetric random vectors
  - For each i, all elements of a column vector are i.i.d random variables with mean 0 and variance $\sigma^2$

- **The variance $\sigma^2$ is a function of the slacks $\delta_i$**

# Parameter Design and Error Control (III)

Theorem: Setting $\delta_i$ to satisfy:

**Tolerable Eigen-Error**

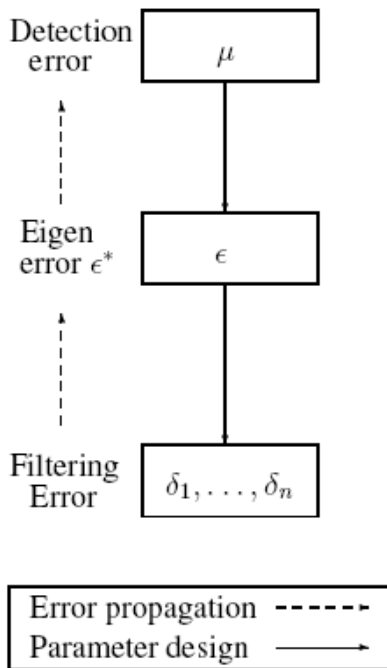**Average of Perturbed eigenvalues**

$$2\sqrt{\frac{\bar{\lambda}}{m} \cdot \sum_{i=1}^{n} \sigma_i^2} + \sqrt{\left(\frac{1}{m} + \frac{1}{n}\right) \sum_{i=1}^{n} \sigma_i^4} = \epsilon$$

Guarantees $\epsilon^{\dot{c}} < \epsilon$ with high probability.

Detection error: $\mu$

Eigen error $\epsilon^*$: $\epsilon$

Filtering Error: $\delta_1, \ldots, \delta_n$

Error propagation — — — —
Parameter design ———→

Absent:
A connection between local variances and local slacks



WHO'S ABSENT?

Is it *You*?

# Slack Allocation Methods

1. Homogeneous slack allocation: uniform distribution of errors in range $[-\delta_i, \delta_i]$

- $\sigma_i = \dfrac{\delta_i^2}{3}$ , results in closed expression for $\delta$

2. Homogeneous slack allocation: local variance estimation

- $\sigma_i = \sigma_i(\delta)$ , monitors approximate locally by fitting an (e.g., quadratic) function according to a recent window of data. Approximation sent to coordinator.

3. Heterogeneous slack allocation.

- Assume uniform distribution of errors in range

- Minimize communication; Solve using Lagrange multipliers.

# Evaluation:Accuracy and Cost

- **Given user-specified false alarm rate, evaluate the actual detection accuracy and communication overhead**

- **Experiment setup**
  - Abilene backbone network data of one week:
    - 121 flows, 41 links, 1008 10 minute periods
  - Traffic matrices of size 1008 X 41
  - Set uniform slack $\delta_i = \delta$ for all monitors
  - Injected: 60 small "bursts" +60 large "anomalies"
  - Threshold corresponding 0.5% false alarm rate
  - How many experiments (repetitions)?

# Evaluation Metrics

- False alarm rate = false alarms/ bursts
- Missed detection rate = missed detections/anomalies
- Cost = num/(n*m) = messages per monitor per sampled time points
  - num = all exchanged messages
  - n = number of monitors
  - M = number of time series points
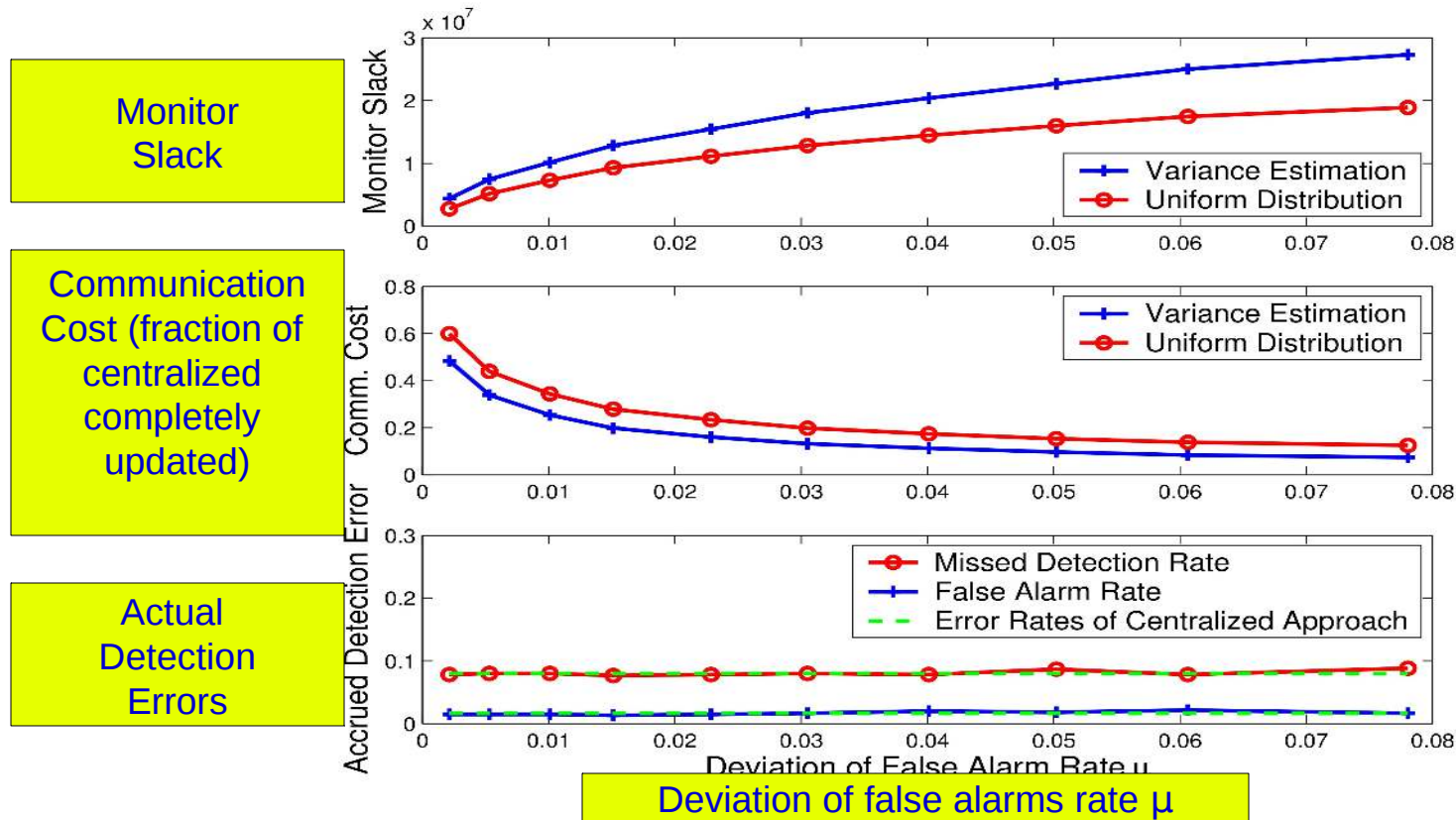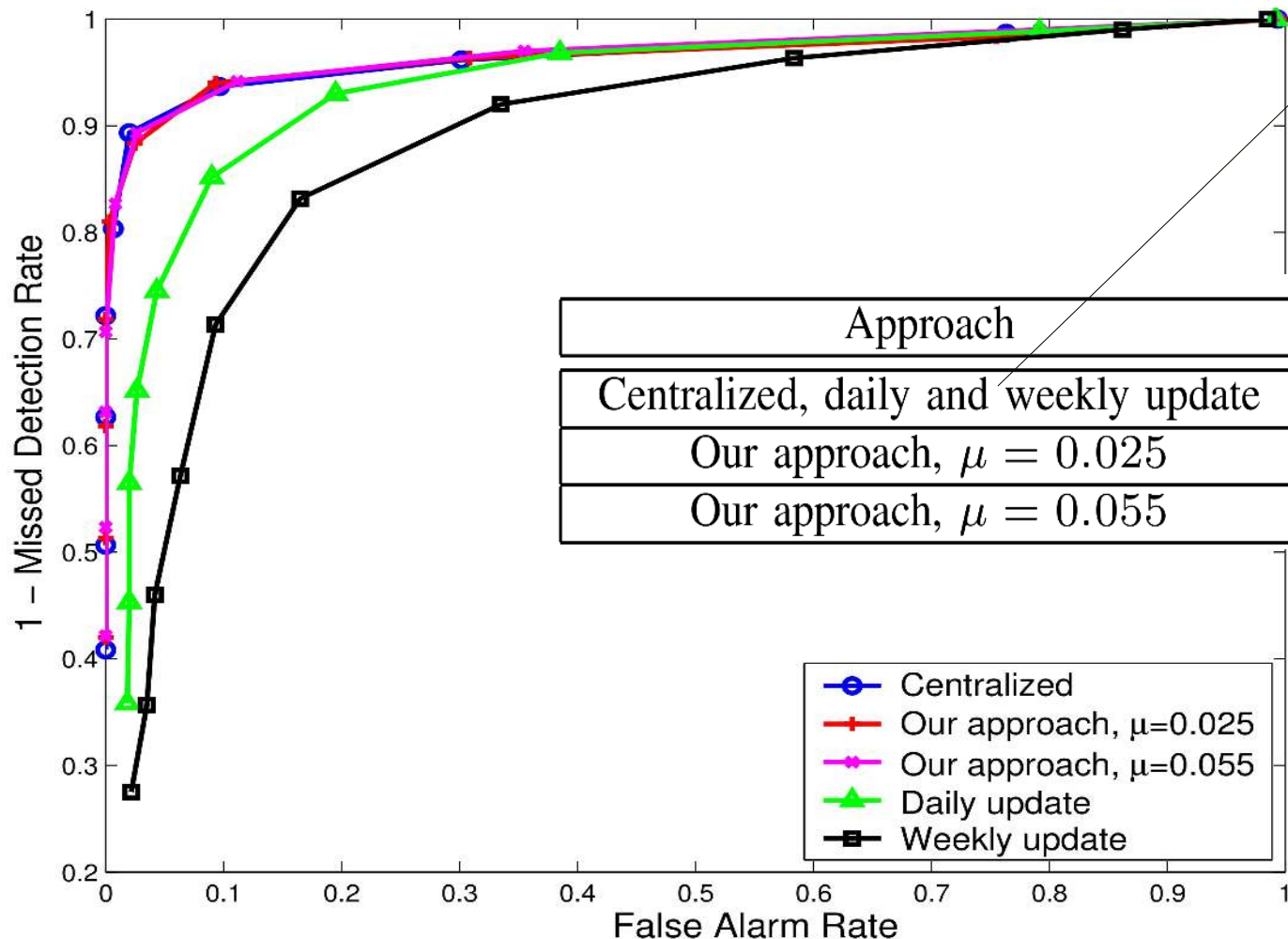
# Evaluation Results



Fig. 6. Monitor slacks, communication cost and accrued detection. The dashed line is the detection error of centralized approach with complete data.

# Observations

- Homogenous variance estimation outperforms Homogenous Uniform, but not by much (5%-10%).

- Homogenous Uniform method is simple.

- Homogenous Uniform might be "good enough".

- 80%-90% of the transferred data can be saved without hurting performance.

# ROC – Receiver Operating Characteristic Curve



Update rate
Of PCA
(data is always full)

| Approach | Communication Cost |
|---|---|
| Centralized, daily and weekly update | 1.000 |
| Our approach, $\mu = 0.025$ | 0.159 |
| Our approach, $\mu = 0.055$ | 0.097 |

Legend:
- Centralized
- Our approach, $\mu=0.025$
- Our approach, $\mu=0.055$
- Daily update
- Weekly update

X axis: False Alarm Rate
Y axis: 1 – Missed Detection Rate
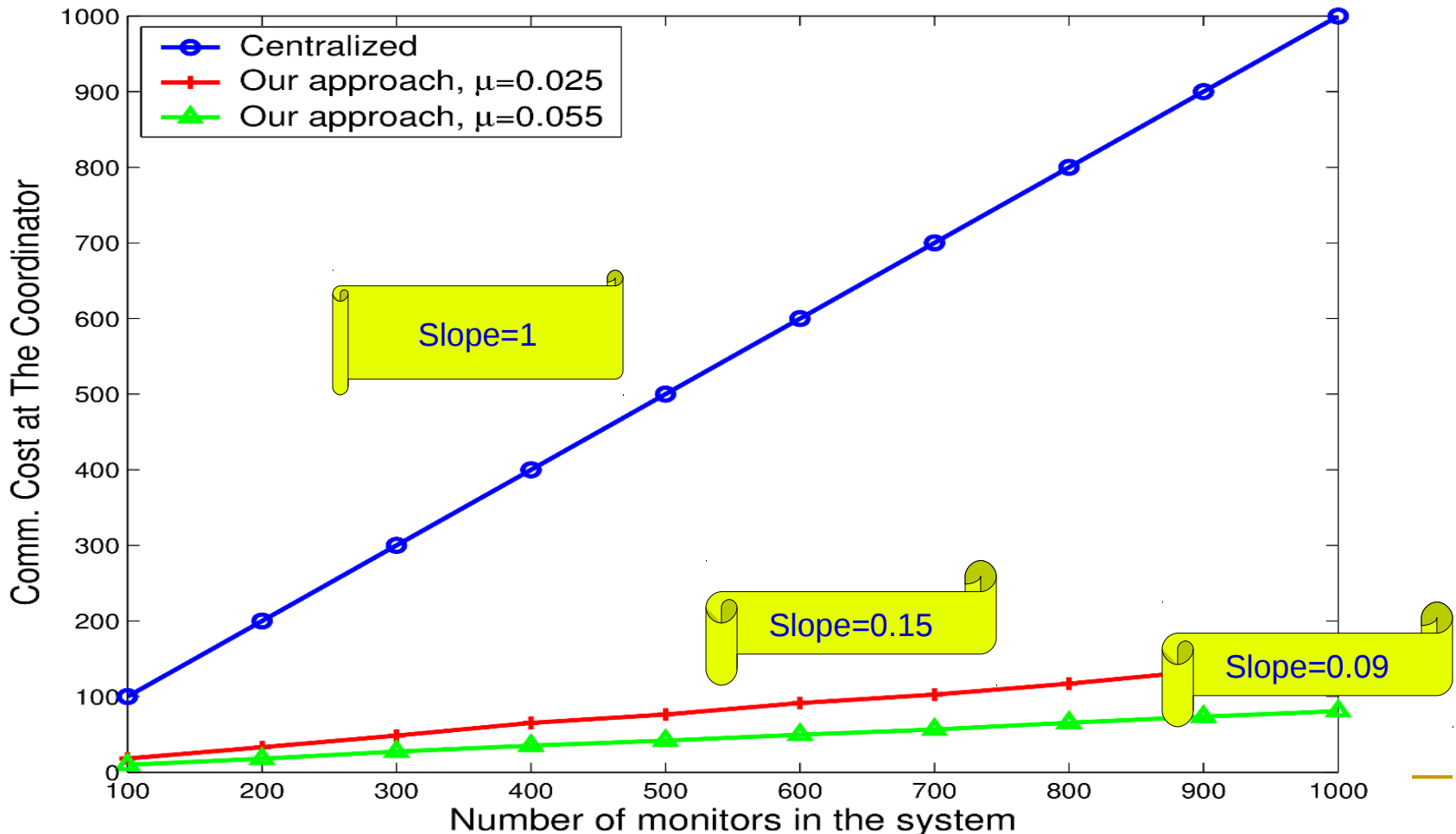
# Evaluation of Scalability

- BRITE topology generator
- 100-1000 links
- Up to 500*500 Origin-Destination flows
- 4 weeks of realistic data, based of statistical characteristics of Abilene
- In each experiment on n nodes: 5 repetitions, on n randomly picked nodes.

# Graceful Scalability by number of monitors: coordinator communication

# Summary

- A communication-efficient framework that
  - detects anomalies at desired accuracy level
  - with minimal communication cost
- A distributed protocol for data processing
  - Local monitors decide when to update data to coordinator
  - Coordinator makes global decision and feedback to monitors
- An algorithmic framework to guide the tradeoff between communication overhead and detection accuracy

# Discussion

## *References*

[Huang'07] *Communication-Efficient Online Detection of Network-Wide Anomalies.* L. Huang, X. Nguyen, M. Garofalakis, J. Hellerstein, M. Jordan, A. Joseph and N. Taft. To appear in INFOCOM'07.

[Lakhina'04] *Diagnosing Network-Wide Traffic Anomalies.* A. Lakhina, M. Crovella and C. Diot. In SIGCOMM '04.

[Jensen & Solomon] *A Gaussian approximation for the distribution of definite quadratic forms.* D.R. Jensen and H. Solomon, In J. Amer. Stat. Assoc., 67:898-902 (1972).

[Jackson & Mudholkar] *Control Procedures for Residuals Associated with Principal Component Analysis.* J. E. Jackson and G. S. Mudholkar, Technometrics, pages 341–349, 1979.
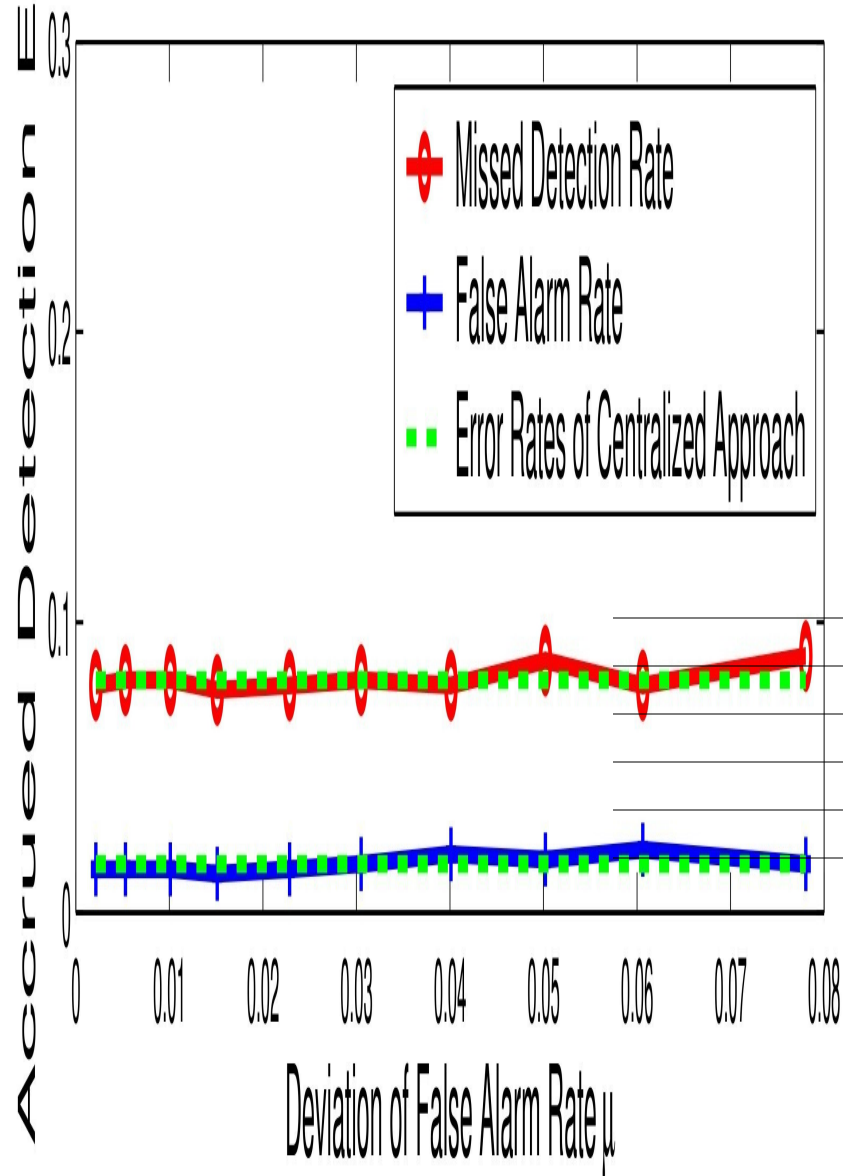
# Weaknesses (My Opinion)

- Symmetry + Independence
- Experiments

# symmetry + independence

- Is the symmetry + independence assumption valid?

- Correlation may result from simultaneous errors upon surprising data changes, or from (cyclic?) bursts induced by the updating algorithm.

# Experiments



Single Experiment
Quantization error:
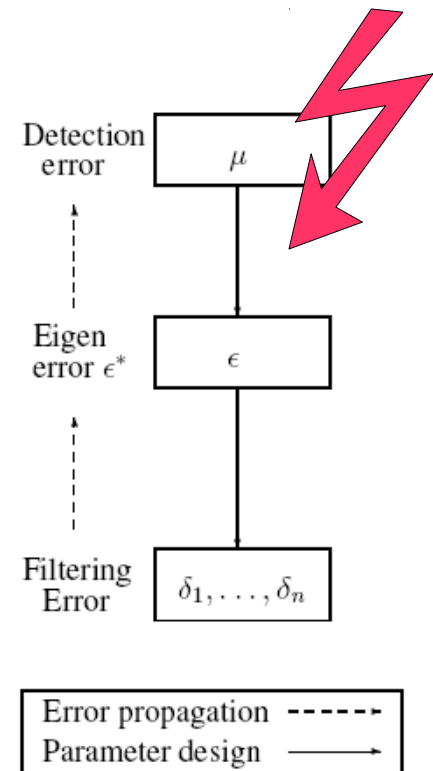1/60=0.016
(means one alarm)

# Experiments: Lack of Trend

- Experiments do not show a statistically significant trend (dependency) of "tolerated deviation from false alarm rate"and actual false alarm rate.
- Estimations are too loose, or
- Experiments are too synthetic
- Between the lines: user is expected to trust experiment results.



Detection error $\mu$

Eigen error $\epsilon^*$ $\epsilon$

Filtering Error $\delta_1, \ldots, \delta_n$

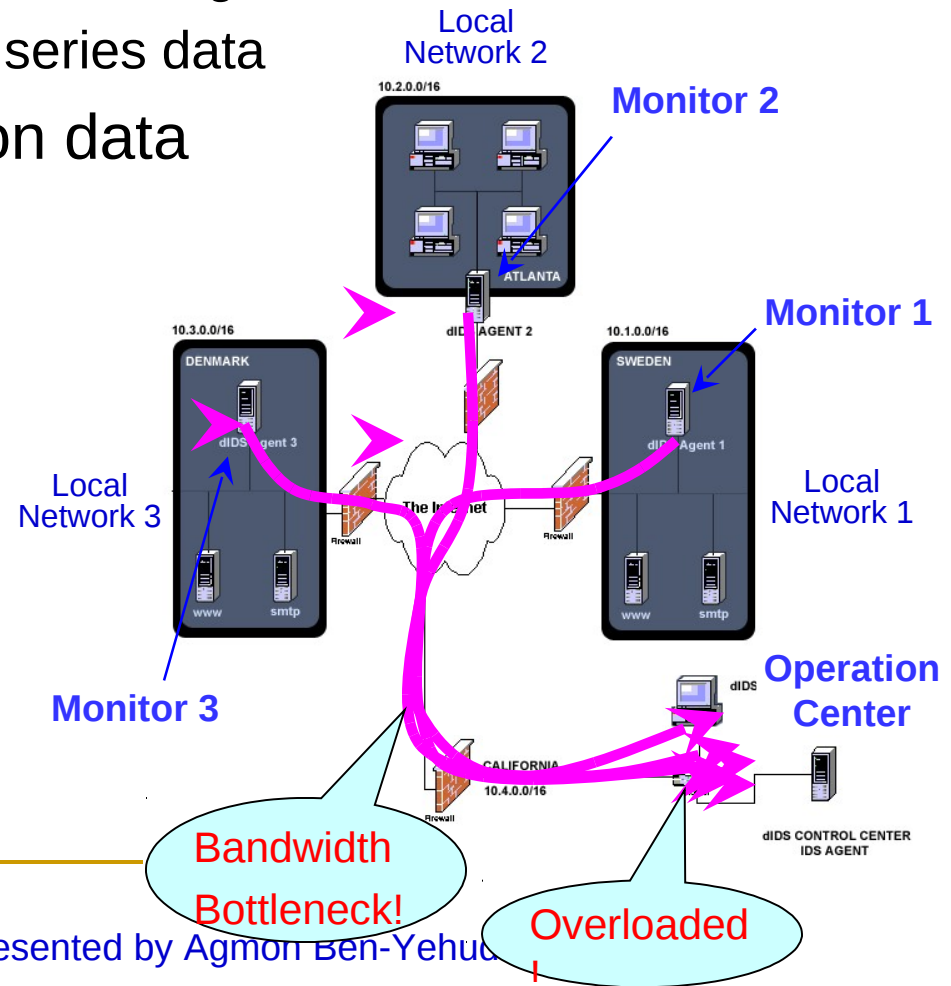Error propagation - - - - - ▸
Parameter design ⟶

# My Summary

- The decentralized algorithm works well in practice according to insufficient experiments.

- The tuning knob was not proved to work in experiments (to be connected to practical accuracy guarantees).

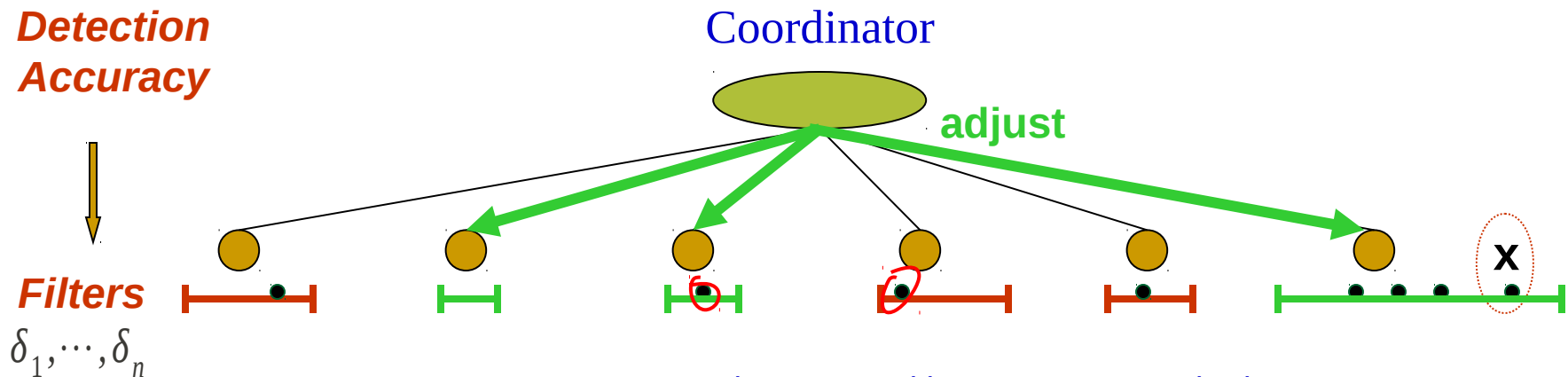- Noisier experiments are needed.

# Backup Slides

# Traditional Distributed Monitoring

- Large-scale network monitoring and detection systems
  - Distributed and collaborative monitoring boxes
  - Continuously generating time series data
- Existing research focuses on data streaming
  - *Centrally* collect, store and aggregate network state
  - Well suited to answering approximate queries and continuously recording system state
  - Incur high overhead!

Local Network 2

10.2.0.0/16

**Monitor 2**

ATLANTA

dIDS AGENT 2

10.3.0.0/16

DENMARK

dIDS Agent 3

Local Network 3

www    smtp

**Monitor 3**

10.1.0.0/16

SWEDEN

dIDS Agent 1

**Monitor 1**

Local Network 1

www    smtp

The Internet

Firewall

Firewall

CALIFORNIA

10.4.0.0/16

Firewall

dIDS

**Operation Center**

dIDS CONTROL CENTER
IDS AGENT

Bandwidth Bottleneck!

Overloaded!

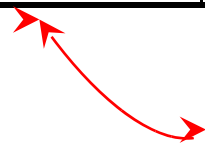Huang et al., presented by Agmon Ben-Yehuda

# Our Distributed Processing Approach

- A coordinator
  - Is aggregation, correlation and detection center
- A set of distributed monitors
  - Each produces a time series signals
  - Processes data locally, only sends needed info. to coordinator
  - No communication among monitors
  - *Coordinator tells monitors the level of accuracy for signal updates*

# Performance

| μ | Missed Detections | | False Alarms | | Data Reduction | |
|---|---|---|---|---|---|---|
| | Week 1 | Week 2 | Week 1 | Week 2 | Week 1 | Week 2 |
| 0.01 | 0 | 0 | 0 | 0 | 75% | 70% |
| 0.03 | 0 | 1 | 1 | 0 | 82% | 76% |
| 0.06 | 0 | 1 | 0 | 0 | 90% | 79% |

error tolerance = upper bound on error

Data Used: Abilene traffic matrix, 2 weeks, 41 links.