

# **Anomaly Detection in large scale networks**

## **State of the art and trend**

**Kavé Salamatian, Lancaster University, ECODE Project**

# OUTLINE

- INTRODUCTION & METHODOLOGIES
- BASICS ANOMALY
- PARAMETRIC METHODS
- NON PARAMETRIC METHODS
- DISTRIBUTED ANOMALY DETECTION

# INTRODUCTION AND METHODOLOGIES

# NETWORK STATE

- NETWORK STATE RESULTS FROM
  - TRAFFIC DEMAND
    - TRAFFIC MATRIX
    - NOT OBSERVABLE DIRECTLY
  - CAPACITY OFFER
    - ROUTING MATRIX, LINK CAP., TRAFFIC ENGINEERING, ETC.
    - MONITORED BY SNMP, ETC.
- NETWORK MANAGER GOAL
  - TO DRIVE THIS EQUILIBRIUM TO THE BEST BENEFICIAL POINT 2
  - BY MANAGING CAPACITY OFFER
    - TRAFFIC ENGINEERING IS THE ART OF MANAGING OFFERED CAPACITY

# NETWORK MONITORING

- MONITORING ?
  - BEING ABLE TO SEPARATE
  - WHAT IS PREDICTABLE
    - EXPECTED, NORMAL, UNDER CONTROL, ...
  - WHAT IS NOT PREDICTABLE
    - UNEXPECTED, ABNORMAL, ...
- INTERPRETATION FRAMEWORK
  - ONLY WHAT IS UNPREDICTABLE HAVE A MEANING
  - WHAT CAN BE PREDICTED DOES NOT HAVE A INFORMATION

3

# CLASSES OF ANOMALY DETECTIONS

- DETERMINISTIC APPROACHES
- SIGNATURE BASED
- COMPLEXITY VS. EXHAUSTIVITY TRADE-OFF
- STATISTICAL APPROACHES
- PROBABILISTIC
- FALSE POSITIVE VS. DETECTION RATE TRADE-OFF

# DETERMINISTIC APPROACHES

- EACH OBSERVATION IS ASSUMED TO RESULT FROM A KNOWN CAUSALITY CHAIN
- ALL ANOMALIES ARE CHARACTERIZED BY CAUSALITY CHAIN THAT DESCRIBES THE SIGNATURE OF AN ANOMALY
- ANOMALY DETECTION CONSIST OF BACK TRACKING THE CAUSALITY SEQUENCE
- NEED **EXHAUSTIVE** ANOMALY SIGNATURE DATABASE
- HAVE TO CHECK ALL POSSIBLE SIGNATURE
- NOT SCALABLE

# STATISTICAL APPROACH

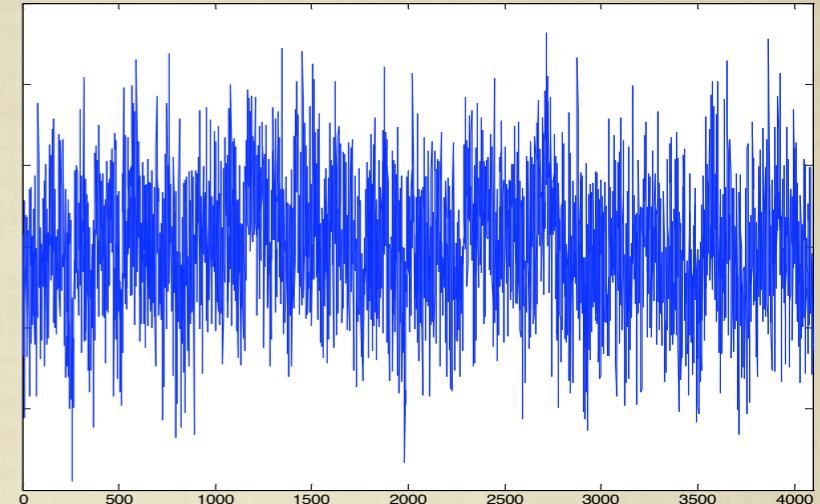
- A DATASET IS A UNIQUE SEQUENCE OF DETERMINISTIC OBSERVATION THAT IS NOT ANYMORE PERFECTLY REPRODUCIBLE
- STATISTICAL APPROACH ASSUMES THAT
  - OBSERVATION RESULTS FROM A RANDOM FUNCTION
    - THEY COME FROM A RANDOM CHOICE IN A (IN)FINITE SET OF «POSSIBLE» OBSERVATIONS.
  - IS THIS ASSUMPTION SOUND ?
- GIVES ACCESS TO AN ARSENAL OF PROBABILISTIC METHODS
  - STATIONARITY : INTRINSIC HYPOTHESIS
    - STATISTICAL PROPERTIES HOLD OVER TIME
    - BET ON FUTURE
    - POSSIBILITY OF RADICAL ERROR
  - ESSENTIAL FALSE ALARM VS. MISDETECTION TRADE-OFF

# EMPIRICAL MODELING CHALLENGE

- MEASUREMENT CONTAINS TWO COMPONENTS
  - A STRUCTURED COMPONENT THAT REFLECTS ESSENTIAL PROPERTIES OF THE PHENOMENON UNDER STUDY
  - A RANDOM COMPONENT THAT REPRESENT FLUCTUATIONS THAT ARE PUT ASIDE FROM THE MODEL
- MODELING CHALLENGE
  - SEPARATE STRUCTURE FROM RANDOMNESS AND CHARACTERIZE IT

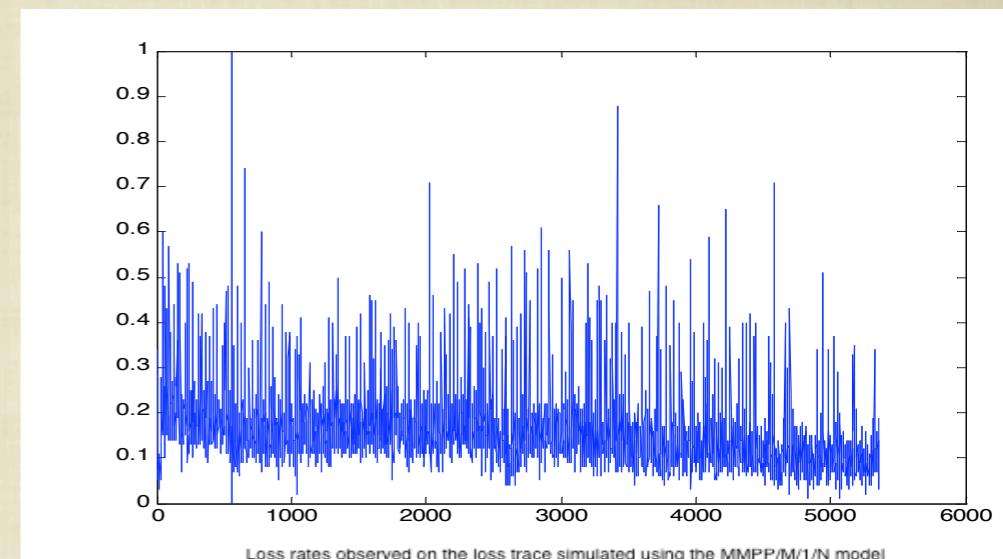
# INTERPRETATION

- TO RELATE EFFECT TO CAUSES



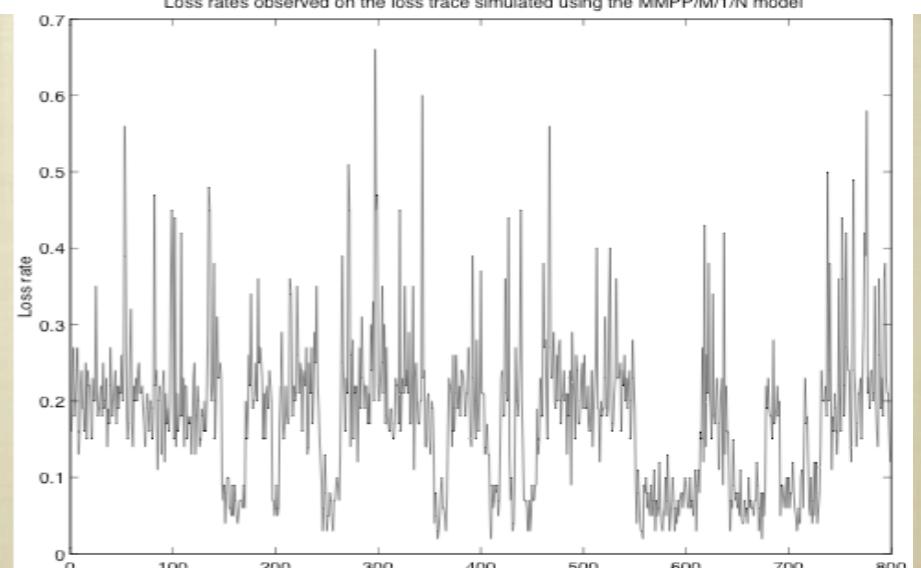
- ENABLE PREDICTION

- AT DIFFERENT TIME SCALE

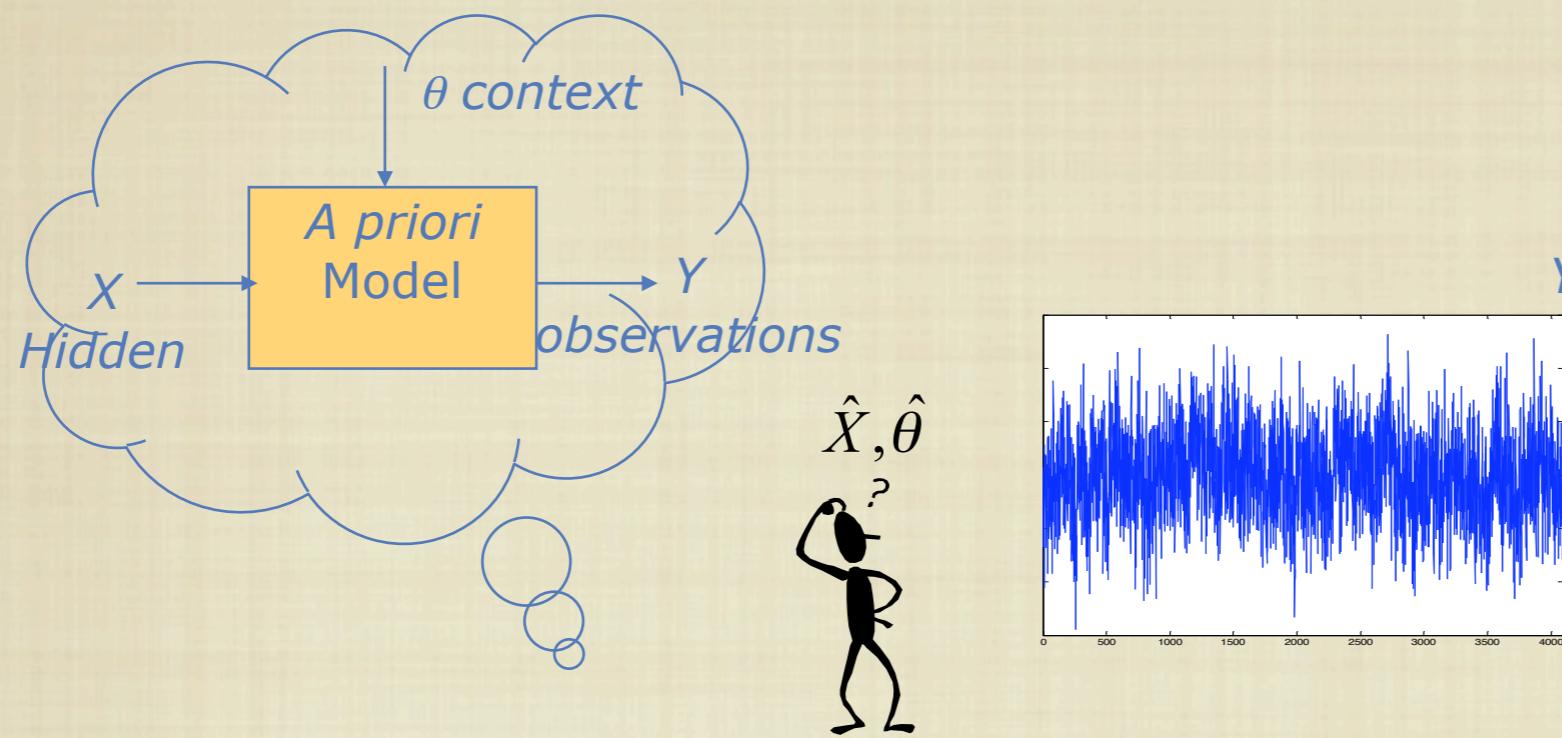


- REACTION

- WHEN CAUSES ARE KNOWN  
ONE CAN AFFECT THEM



# INTERPRETATION FRAMEWORK



- WHAT ARE THE HIDDEN CAUSES AND THE HIDDEN CONTEXT THAT LEADS TO THE OBSERVATION
- THE UNDERSTANDING OF THE PHENOMENON IS CONDENSED IN THE A PRIORI MODEL  $Y=M(x,\theta)$

# INTERPRETATION FRAMEWORK

- TWO STATISTICAL INVERSE PROBLEMS
- MODELING PROBLEM
  - WHAT IS THE CONTEXT PARAMETER  $\theta$
- INTERPRETATION PROBLEM
  - WHAT IS THE HIDDEN INPUT  $X$

# STATISTICAL ANOMALY DETECTION BASICS

# ANOMALIES ???

- WHAT IS AN ANOMALY?
- WHAT IS NORMAL?
- WHAT IS LIKELY?
- ASSUMPTIONS
  - ANOMALY ARE ANOMALOUS
    - DIFFERENT FROM THE NORM
  - ANOMALIES ARE RARE
  - ANOMALOUS ACTIVITY IS MALICIOUS
    - REVERSE ASSUMPTION : ANOMALY REPRESENTS ATTACKS OR MALICIOUS BEHAVIOR
  - ARE THESE ASSUMPTIONS CORRECT ???

# WHAT IS LIKELY ?

## ■ TYPICAL SET

■ THE SET OF OBSERVED SEQUENCE  $Y_k^{n+k}$  SUCH THAT

$$2^{-n(\bar{H}(Y_k^{k-n})-\varepsilon)} \leq \Pr\left\{Y_k^{k+n}\right\} \leq 2^{-n(\bar{H}(Y_k^{k-n})+\varepsilon)}$$

$$\Pr\left\{A_\varepsilon^{(n)}\right\} \geq 1 - \varepsilon$$

## ■ FROM AEP

$$(1 - \varepsilon)2^{-n(\bar{H}(Y_k^{k-n})+\varepsilon)} \leq |A_\varepsilon^{(n)}| \leq 2^{-n(\bar{H}(Y_k^{k-n})+\varepsilon)}$$

# UNIVERSAL ANOMALY DETECTOR !

- A SIMPLE ANOMALY DETECTOR
- CALCULATE THE LIKELIHOOD  $\Pr\{Y_k^{k+n}\}$  OF AN OBSERVED SEQUENCE
- CHECK IF  $2^{-n(\bar{H}(Y_k^{k-n})-\varepsilon)} \leq \Pr\{Y_k^{k+n}\} \leq 2^{-n(\bar{H}(Y_k^{k-n})+\varepsilon)}$
- IF NOT, IT IS AN ANOMALY
- NON PARAMETRIC AND UNIVERSAL ANOMALY DETECTOR
- HOW TO CALCULATE THE LIKELIHOOD ?
  - NEED SIMPLIFYING HYPOTHESIS

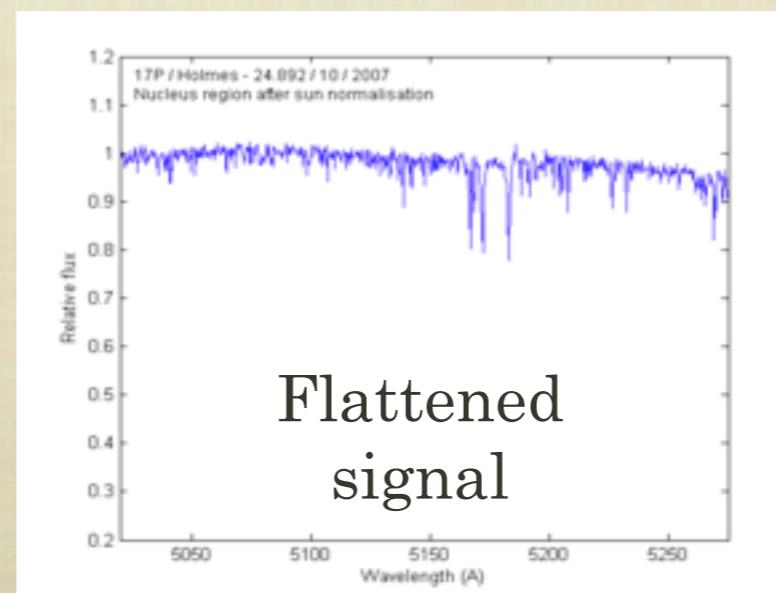
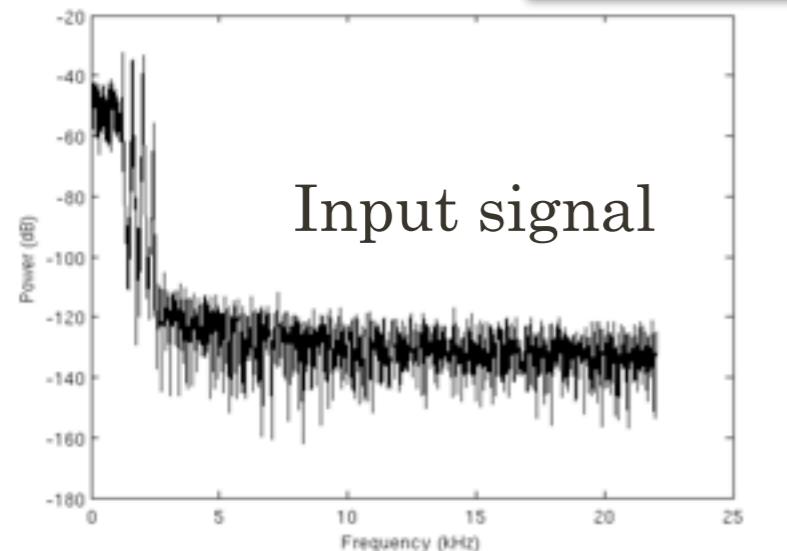
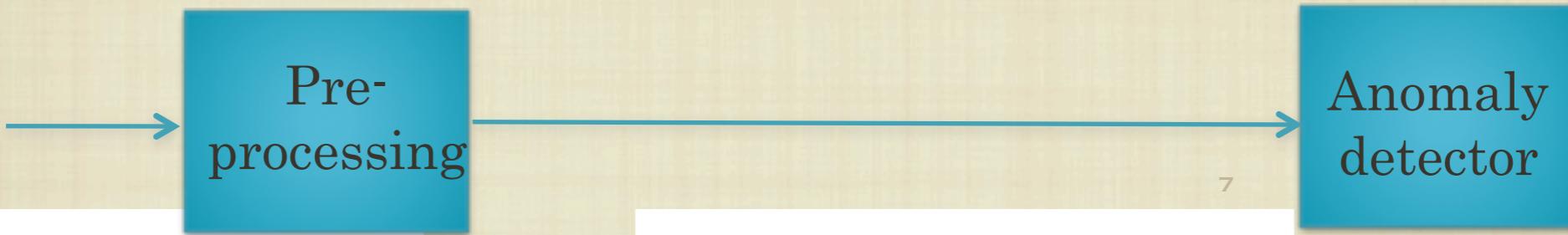
# FILTERING

- WHAT IF THE OBSERVATIONS ARE INDEPENDENT !

$$\Pr\{Y_k^{k+n}\} = (\Pr\{Y_k\})^n$$

- WE HAVE TO TRANSFORM THE OBSERVATION SO THAT IT BECOME INDEPENDENT !

- BASIC SIGNAL PROCESSING QUESTION



# ANOMALY DETECTOR STRUCTURE

- TWO GENERIC STEPS

- ENTROPY REDUCTION/FILTERING

- REMOVE DEPENDENCIES IN DATA

- COMPRESS THE DATA: INFORMATION BOTTLENECK

- REMOVE PREDICTABLE COMPONENT

- DECISION/DETECTION

- APPLY THE TYPICALITY TEST ON THE RESULTING HOPEFULLY INDEPENDENT SIGNAL

- CLASSICAL TEST IN THE NEYMAN-PEARSON FRAMEWORK

# TAXONOMY OF APPROACHES

- PARAMETRIC
  - ASSUME A PARAMETRIC DEPENDENCIES MODEL
- NON PARAMETRIC
  - BROKE THE DEPENDENCIES IN A UNIVERSAL WAY

# PARAMETRIC APPROACHES

- PARAMETRIC APPROACHES

- ASSUMES A PARAMETRIC STRUCTURE FOR DEPENDENCIES

- NORMAL BEHAVIOR MODEL

- CALIBRATE THE MODEL

- MLE ESTIMATION

- FILTER THE MODEL PREDICTION FROM OBSERVATION

- RESULTS IN AN IID GAUSSIAN INNOVATION PROCESS

- APPLY A NEYMAN-PEARSON STATISTICAL TEST WITH  $H_0$ : THE OBSERVATION IS COMPATIBLE WITH A ZERO MEAN WITH GIVEN VARIANCE

# NON PARAMETRIC APPROACH

- APPLY A DE-CORRELATING TRANSFORM
  - RANDOM PROJECTION, HASHING, SKETCH, ETC...
- LEARN THE RESULTING DISTRIBUTION OF NORMAL BEHAVIOR
  - CLUSTERING, SVM, ETC...
- DETECT ANOMALIES BY CHECKING CLUSTER MEMBERSHIP

# PARAMETRIC TECHNIQUES

# ANOMALY DETECTION STEPS

■ MODELING

■ FILTERING

■ DECISION

# MODEL STRUCTURE

- FIRST STEP IS TO DEFINE A CORRELATION/MODEL STRUCTURE

- TEMPORAL, SPATIAL CORRELATIONS

- THE TRAFFIC IS A DYNAMIC ENTITY

- GENERIC DYNAMICAL MODEL

- LINEAR APPROXIMATION

- LTI MODEL

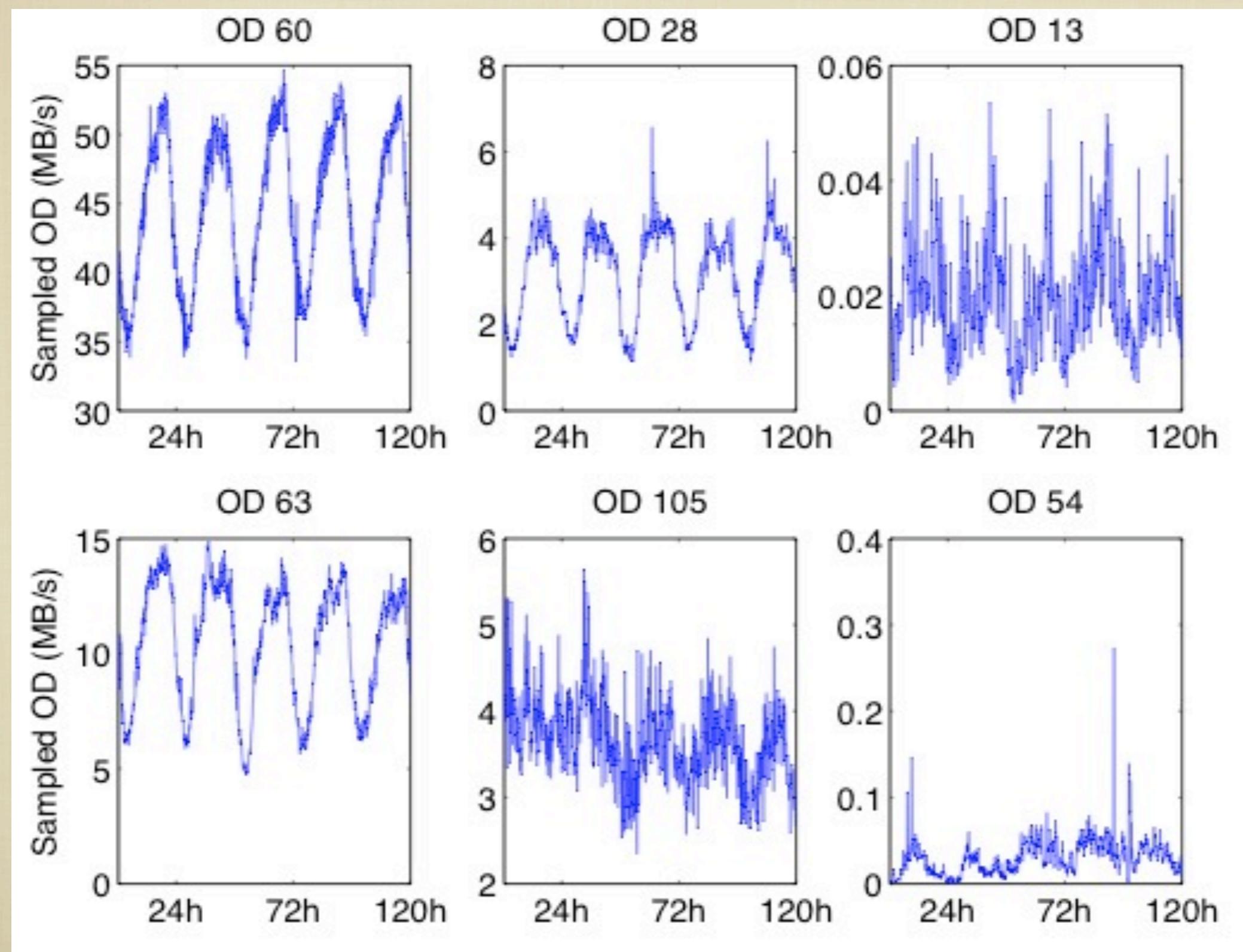
$$\begin{cases} X_{k+1} = f(X_k) + g(U_k) \\ Y_k = h(X_k) + i(U_k) \end{cases}$$

$$\begin{cases} X_{k+1} = AX_k + BU_k \\ Y_k = CX_k + DU_k \end{cases}$$

$$X_k = (x_k^1, x_{k-1}^1, \dots, x_{k-T}^1, x_k^2, x_{k-1}^2, \dots, x_{k-T}^2, x_k^3, \dots, x_{k-1}^L, \dots, x_{k-T}^L)$$

- A WELL BEHAVED NON-LINEAR PROCESS WITH ANY MEMORY CAN BE APPROXIMATED BY A LARGE ENOUGH LINEAR MODEL

# TRAFFIC DYNAMIC



# STOCHASTIC MODELS

- STATE MODEL:

$$\begin{cases} X(t+1) = C * X(t) + W(t) \\ Y(t) = A * X(t) + V(t) \end{cases} \begin{cases} W \sim N(0, Q) \\ V \sim N(0, R) \end{cases}$$

- HOW TO CALIBRATE C, Q AND R?

- MAXIMUM LIKELIHOOD METHOD (EM)

- FIND VALUES OF C, Q, R THAT MAXIMIZE THE LIKELIHOOD OF THE OBSERVATIONS

- PCA METHOD

- FIND VALUES OF C, Q, R THAT GIVE THE BEST K-DIMENSIONAL APPROXIMATION

# PCA METHOD

## ■ PCA THEOREM

$$\mathbf{X} = \sum_{i=1}^K Y_i \phi_i$$

## ■ KARHUNEN-LOEVE THEOREM

$$X_l(t) = \sum_{i=1}^K \sum_{j=1}^{\infty} Y_{i,j}^l \Phi_{i,j}(t)$$

## ■ KLT/PCA BASED MODEL

$$\hat{X}_l(kT) = \sum_{i=1}^L \sum_{j=1}^M Y_{i,j}^l \Phi_{i,j}^k.$$

# PCA METHOD

## ■ HOW TO DERIVE THE BASIS

### ■ PCA

$$\Sigma \phi_i = \lambda_i \phi_i$$

### ■ KLT

$$\sum_{i=1}^K \int_a^b \sigma_{i,l}(s) \Phi_{i,j}(s-t) ds = \lambda_{l,j} \Phi_{l,j}(t), \quad j > 0.$$

### ■ GALERKIN METHOD : APPLY PCA TO

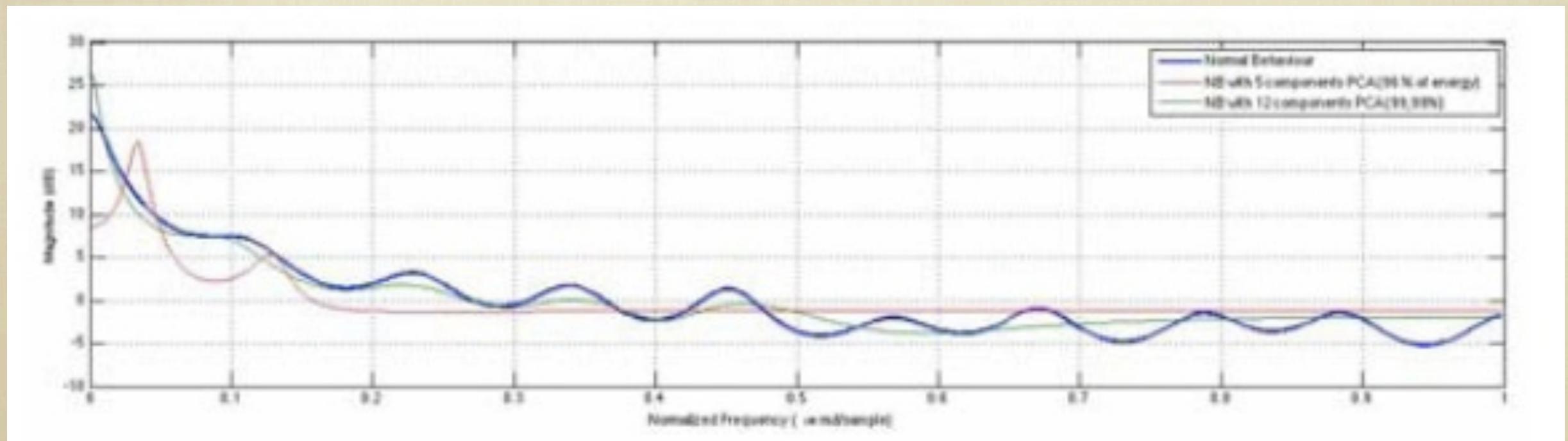
$$\mathbf{X} = \begin{pmatrix} x_1(1) & \dots & x_1(n-N) \\ x_1(2) & \dots & x_1(n-N+1) \\ \vdots & \ddots & \vdots \\ x_1(N) & \dots & x_1(n) \\ x_2(1) & \dots & x_2(n-N) \\ \vdots & \ddots & \vdots \\ x_2(N) & \dots & x_2(n) \\ \vdots & \ddots & \vdots \\ x_K(1) & \dots & x_K(n-N) \\ \vdots & \ddots & \vdots \\ x_K(N) & \dots & x_k(n) \end{pmatrix}$$

# DEMONSTRATION

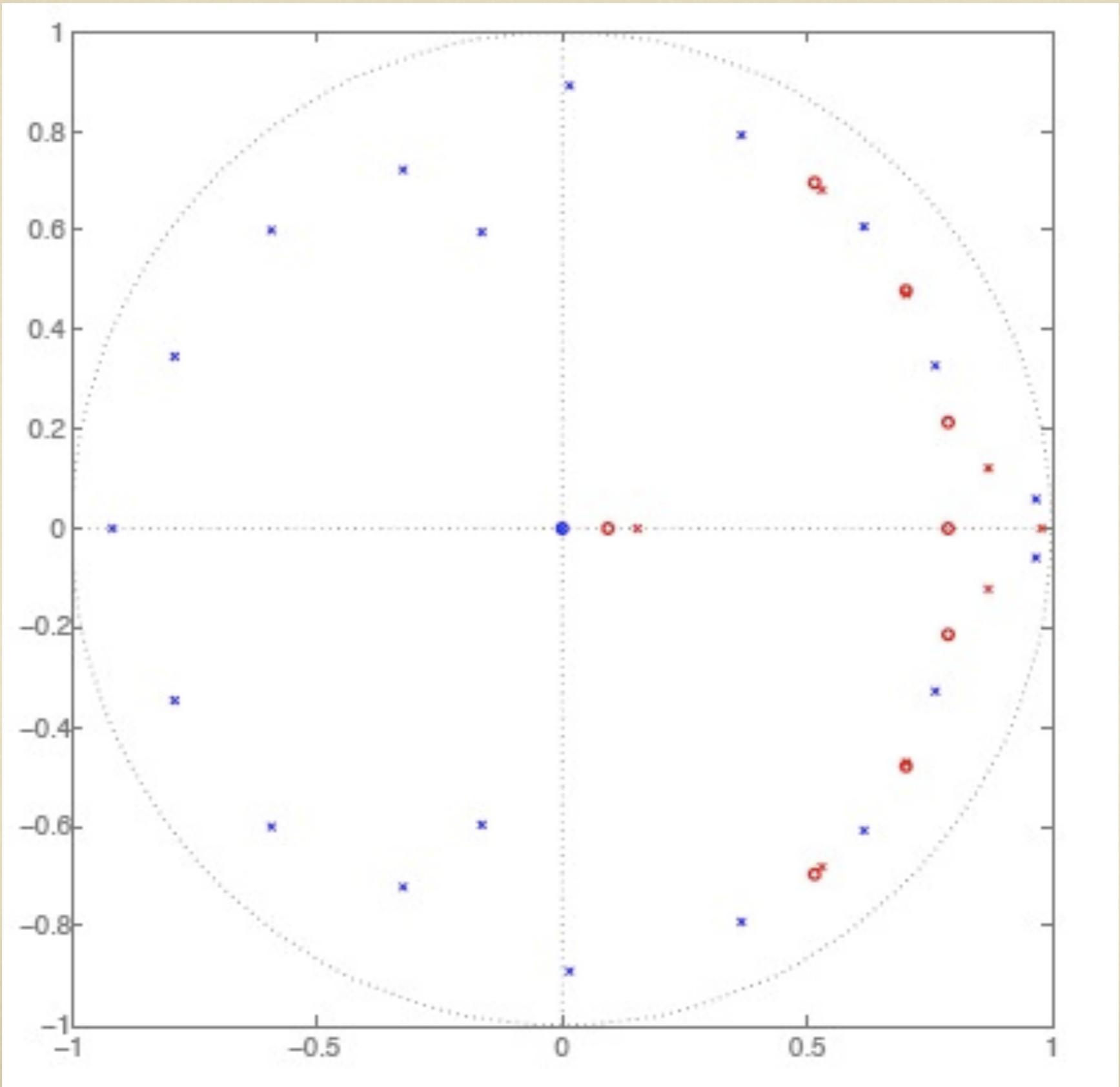
■ MATLAB

# NORMAL BEHAVIOR ?

- WHAT IS THE NORMAL BEHAVIOR CHARACTERIZED BY THE MODEL
- TIME AND FREQUENCY REPRESENTATION ARE EQUIVALENT



# PCA MODEL



# FILTERING

- FILTER SEPARATE THE SIGNAL SPACE INTO TWO COMPONENTS
  - SIGNALS THAT PASS THROUGH THE FILTER
  - SIGNALS THAT ARE REJECTED BY THE FILTER
- ANOMALY DETECTION FILTER
  - PASS EVERYTHING IS COMPATIBLE WITH THE NORMAL BEHAVIOR
  - BLOCK EVERYTHING IS NOT COMPATIBLE WITH THE NORMAL BEHAVIOR

# KALMAN FILTER

## ■ FILTERING WHAT IS COMPATIBLE WITH THE MODEL

### ■ TWO STEPS

#### ■ PREDICTION

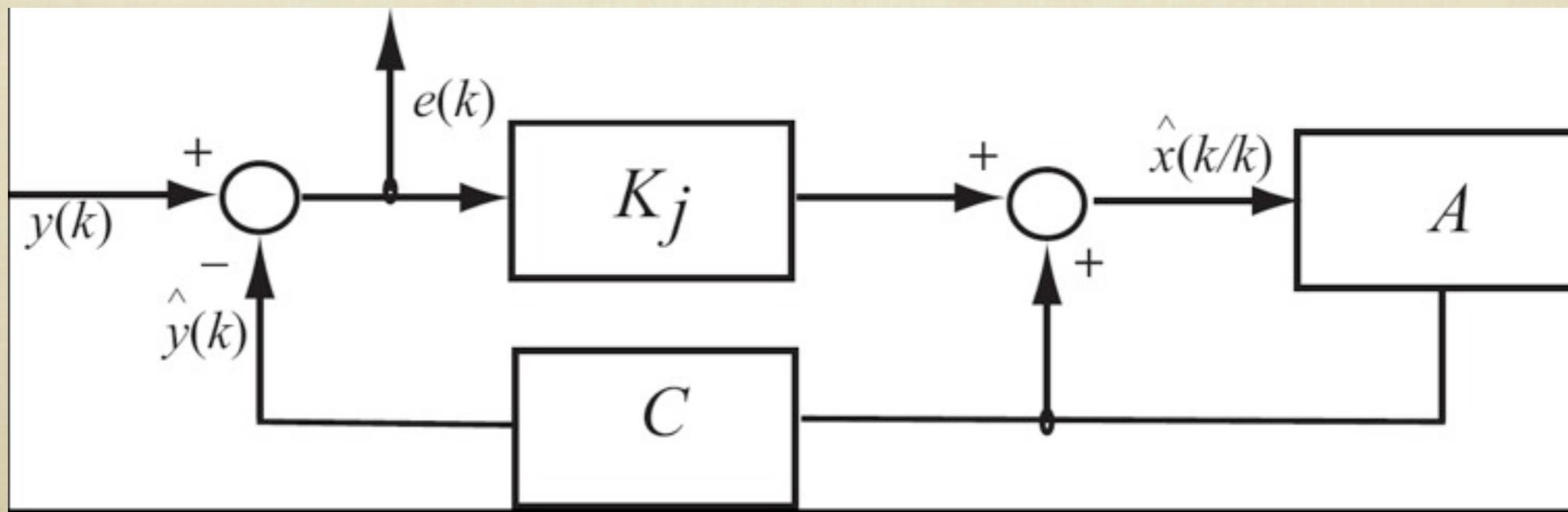
$$\hat{X}^-(t+1) = C\hat{X}(t)$$

#### ■ CORRECTION

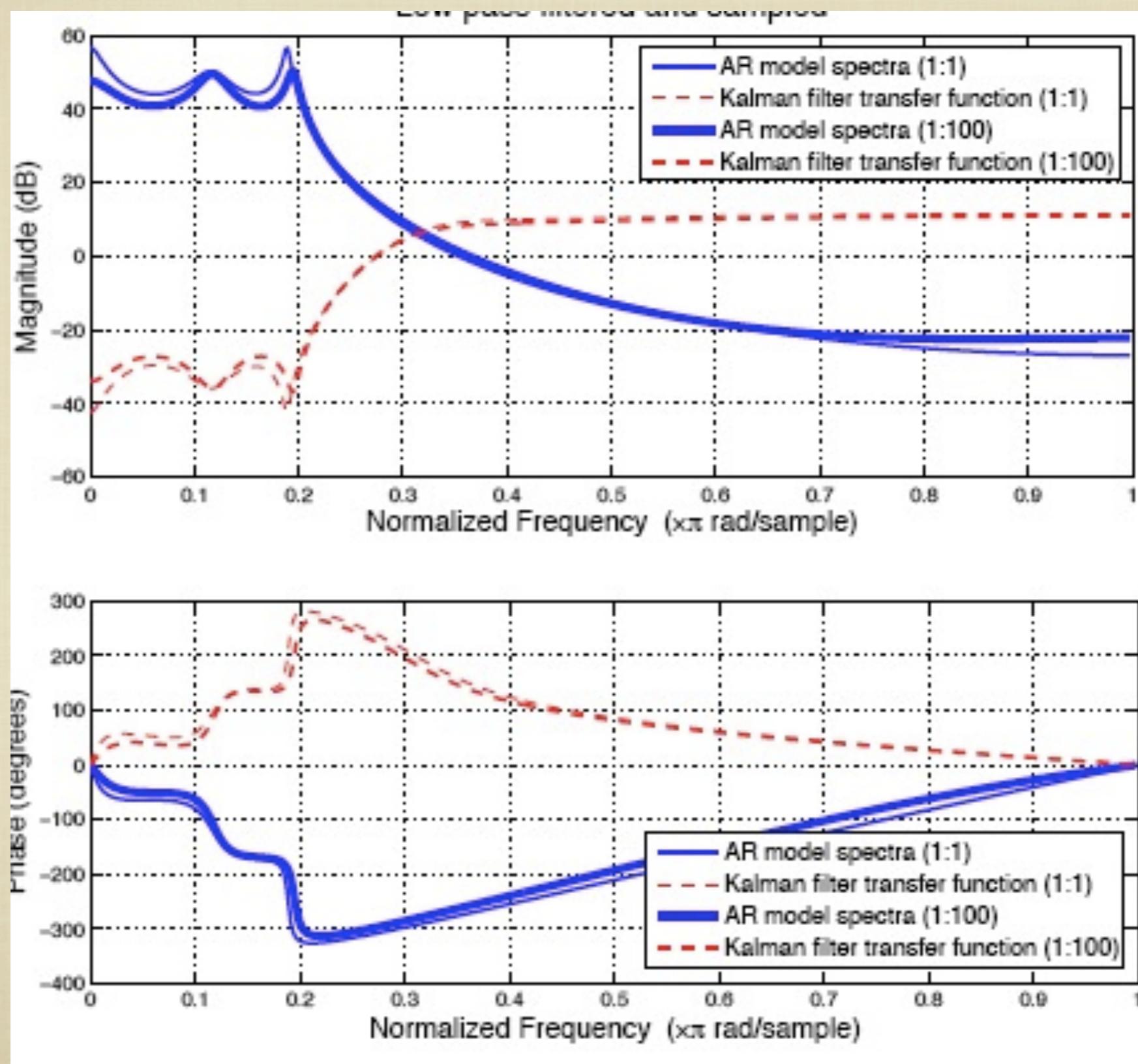
$$\hat{X}(t+1) = \hat{X}^-(t+1) + K(t+1)(Y(t+1) - A\hat{X}^-(t+1))$$

#### ■ INNOVATION

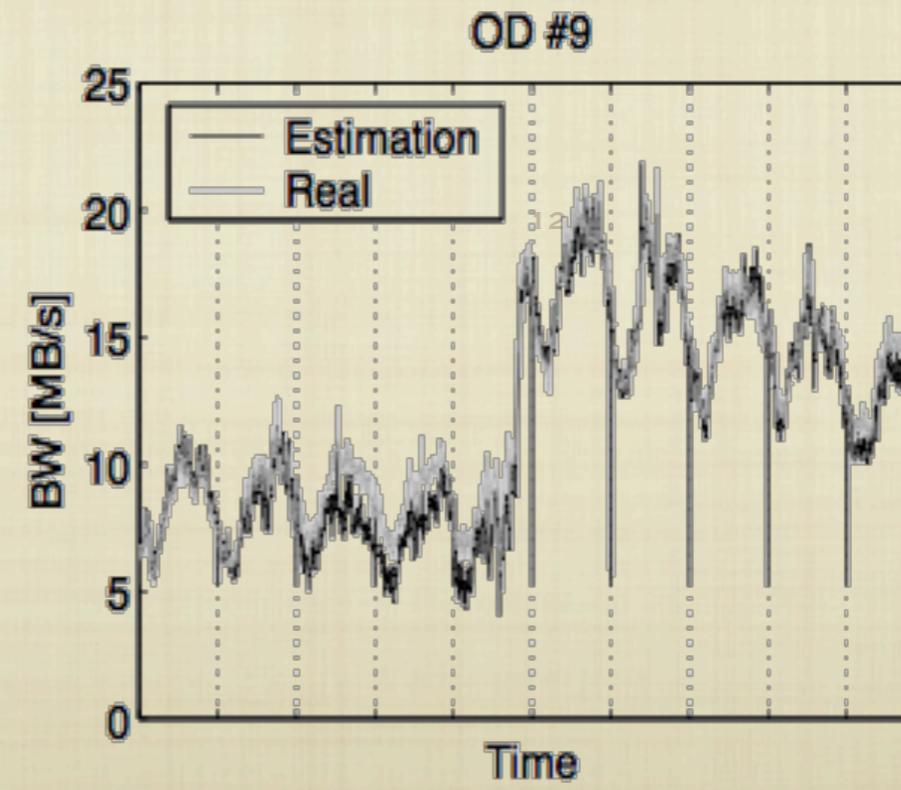
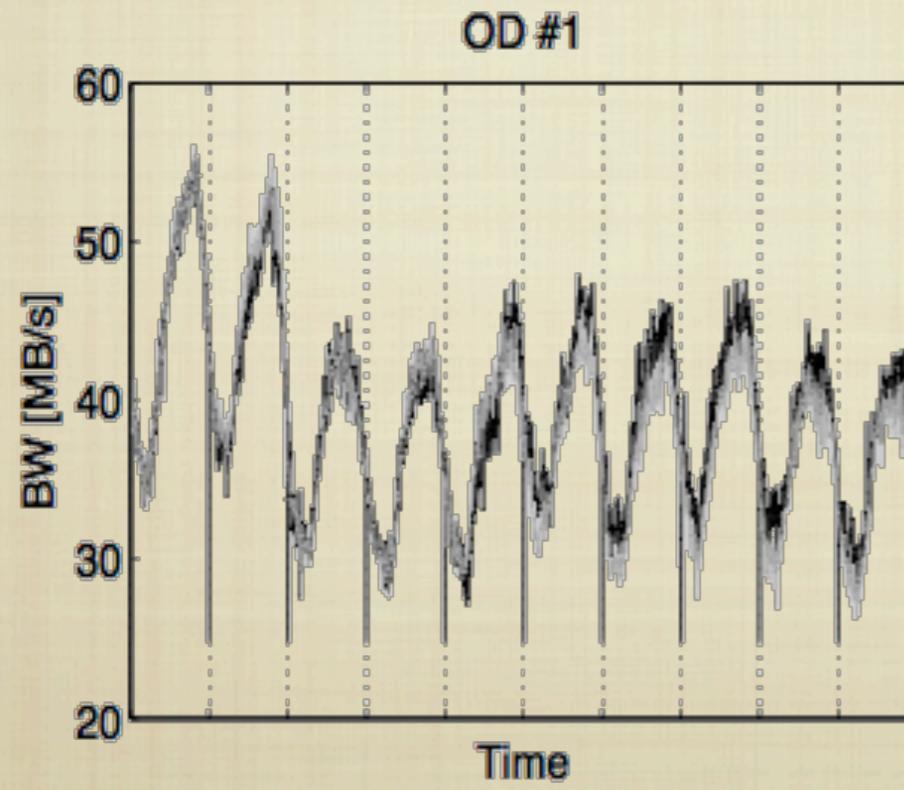
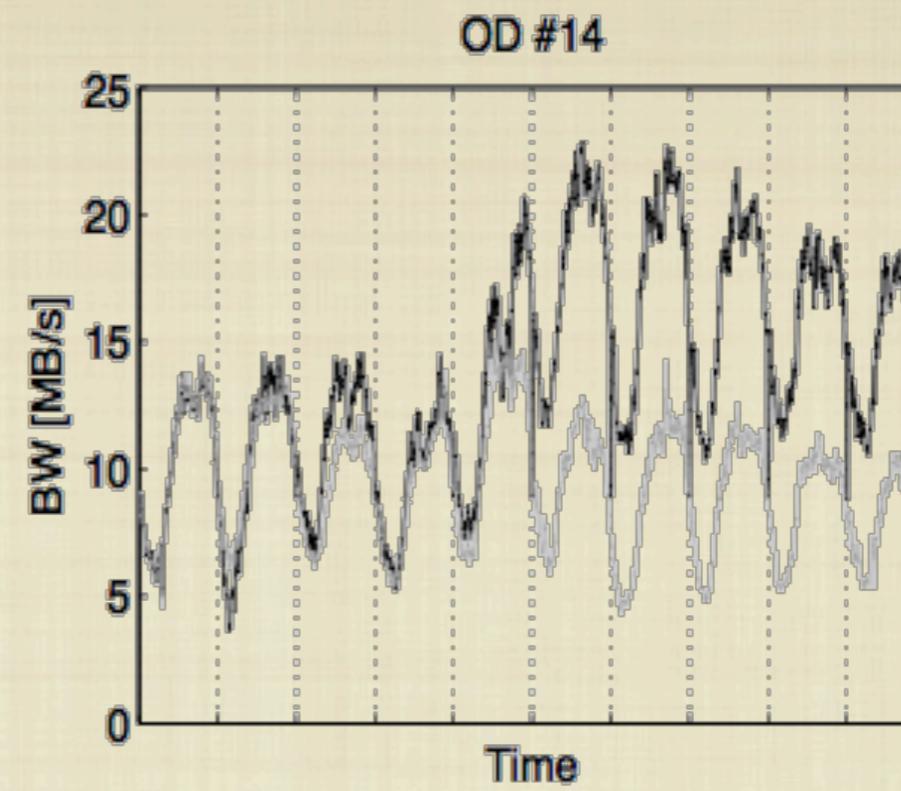
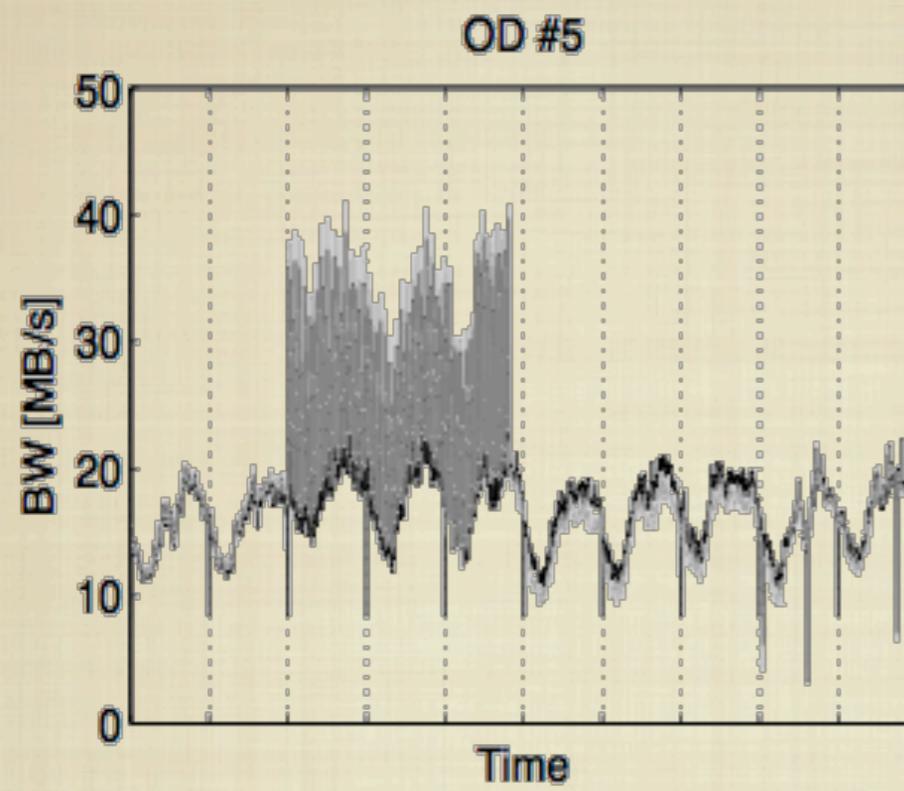
$$\eta(t+1) = Y(t+1) - A\hat{X}^-(t+1)$$



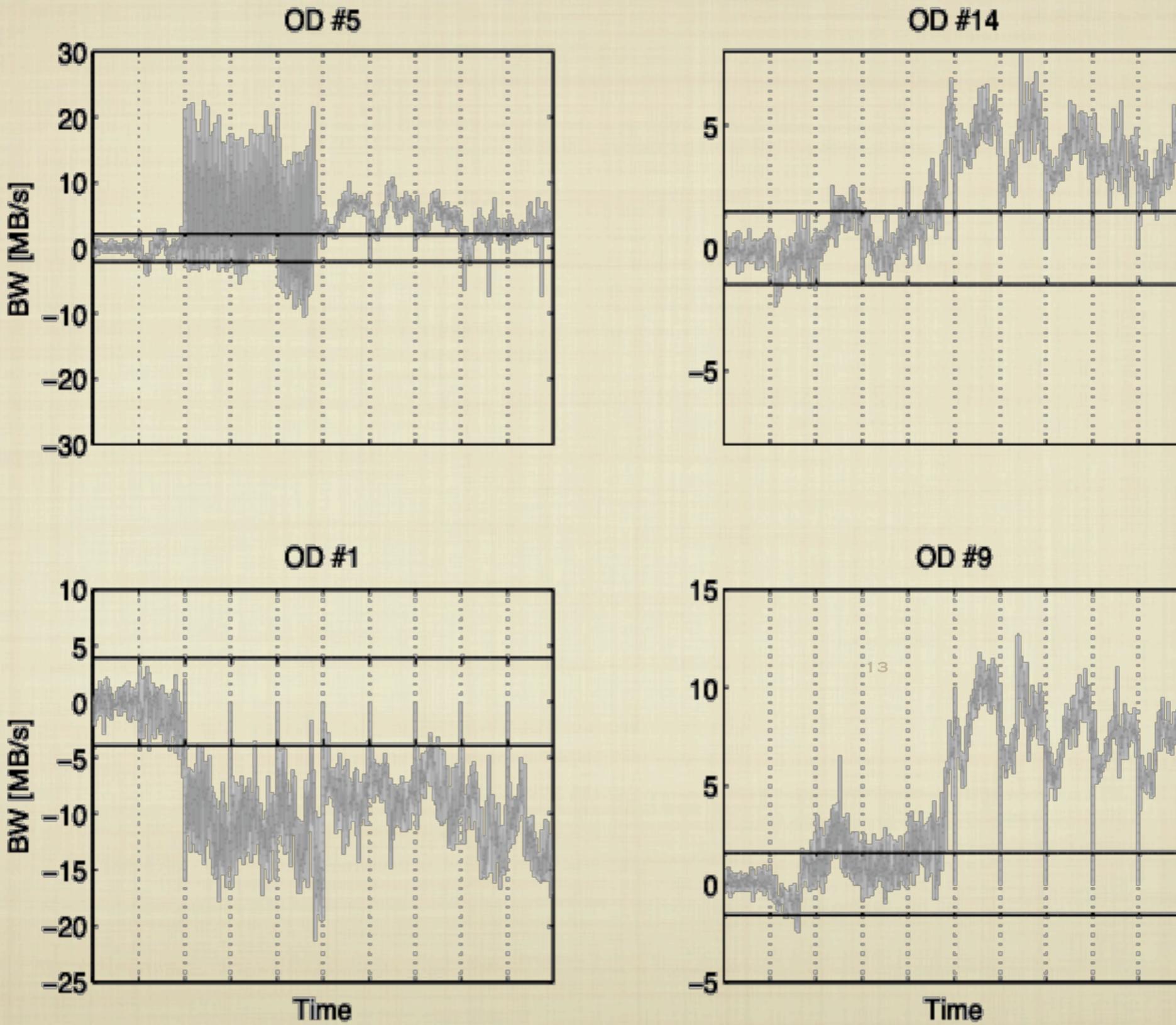
# KALMAN FILTER



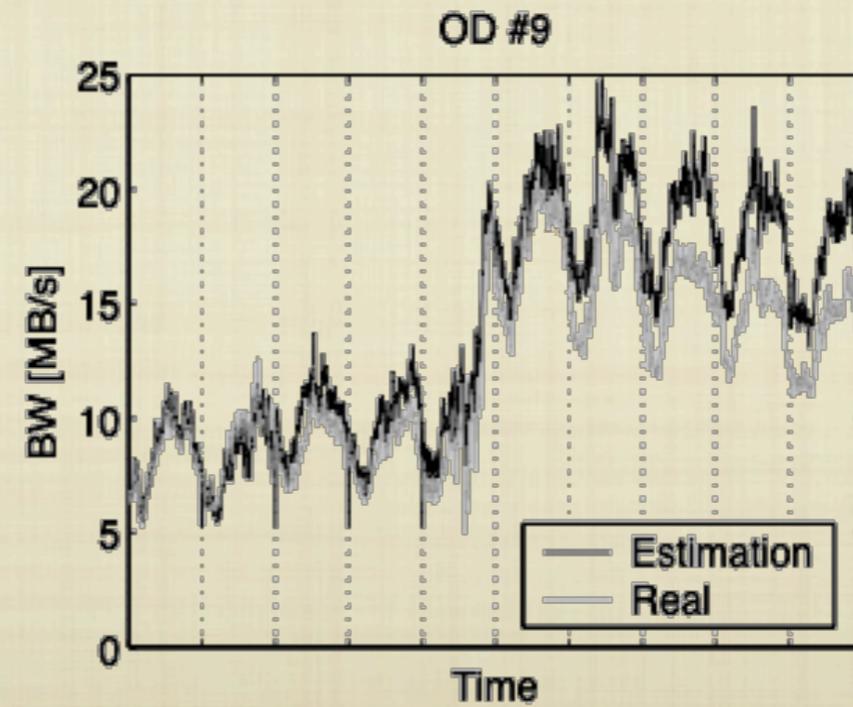
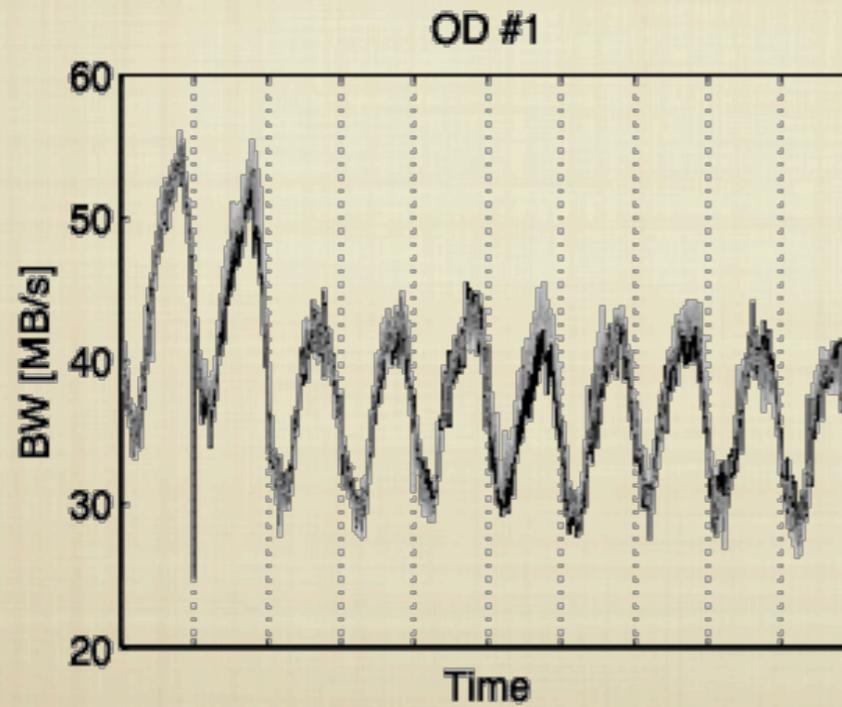
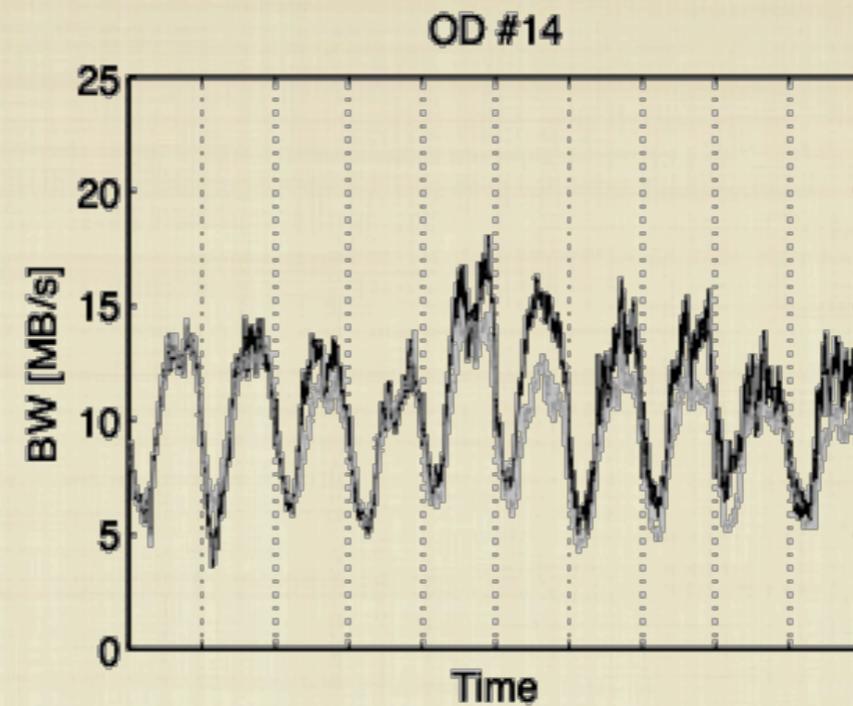
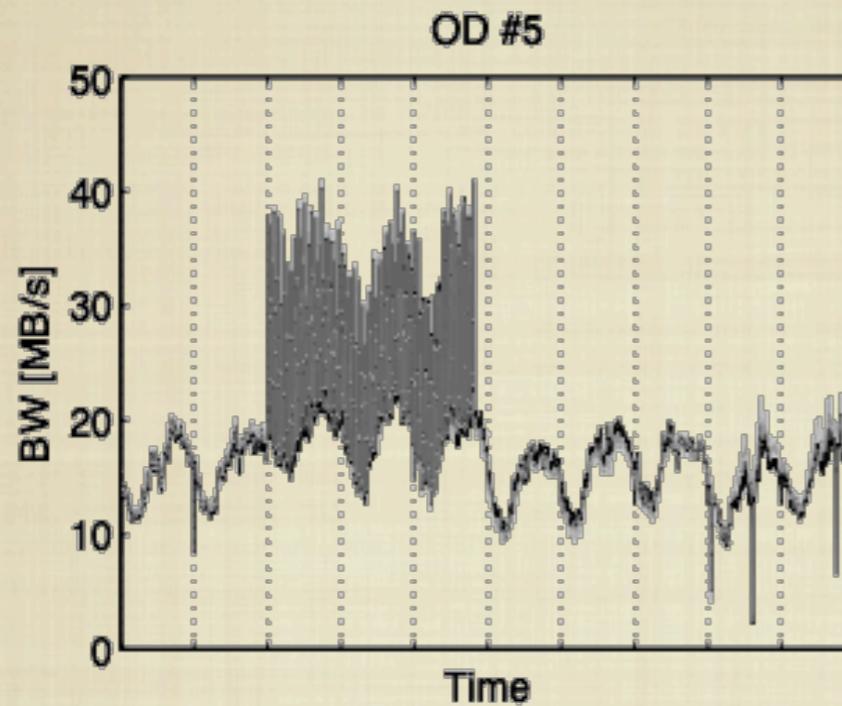
# KALMAN FILTER RESULTS



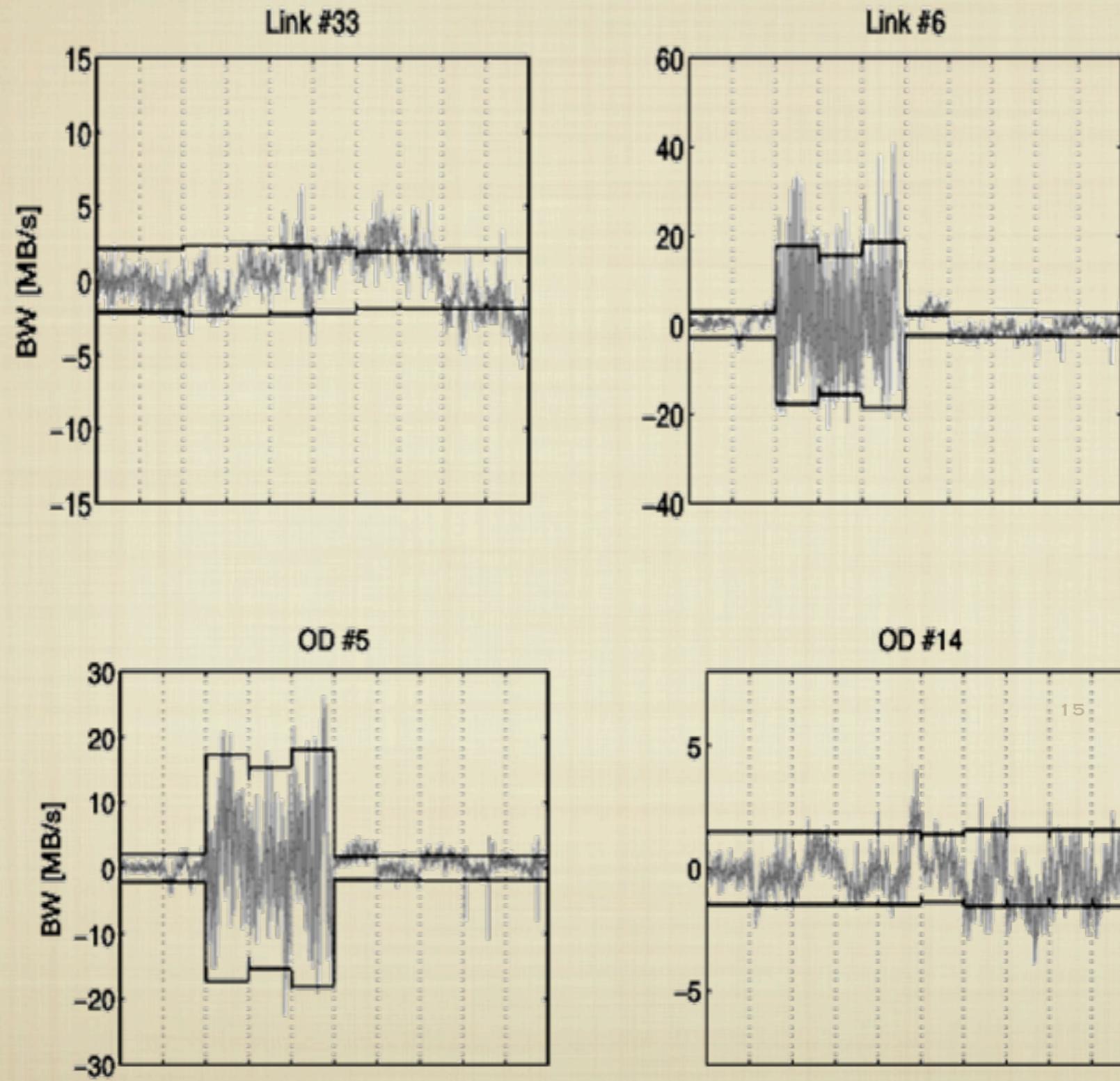
# INNOVATION



# RECALIBRATION !

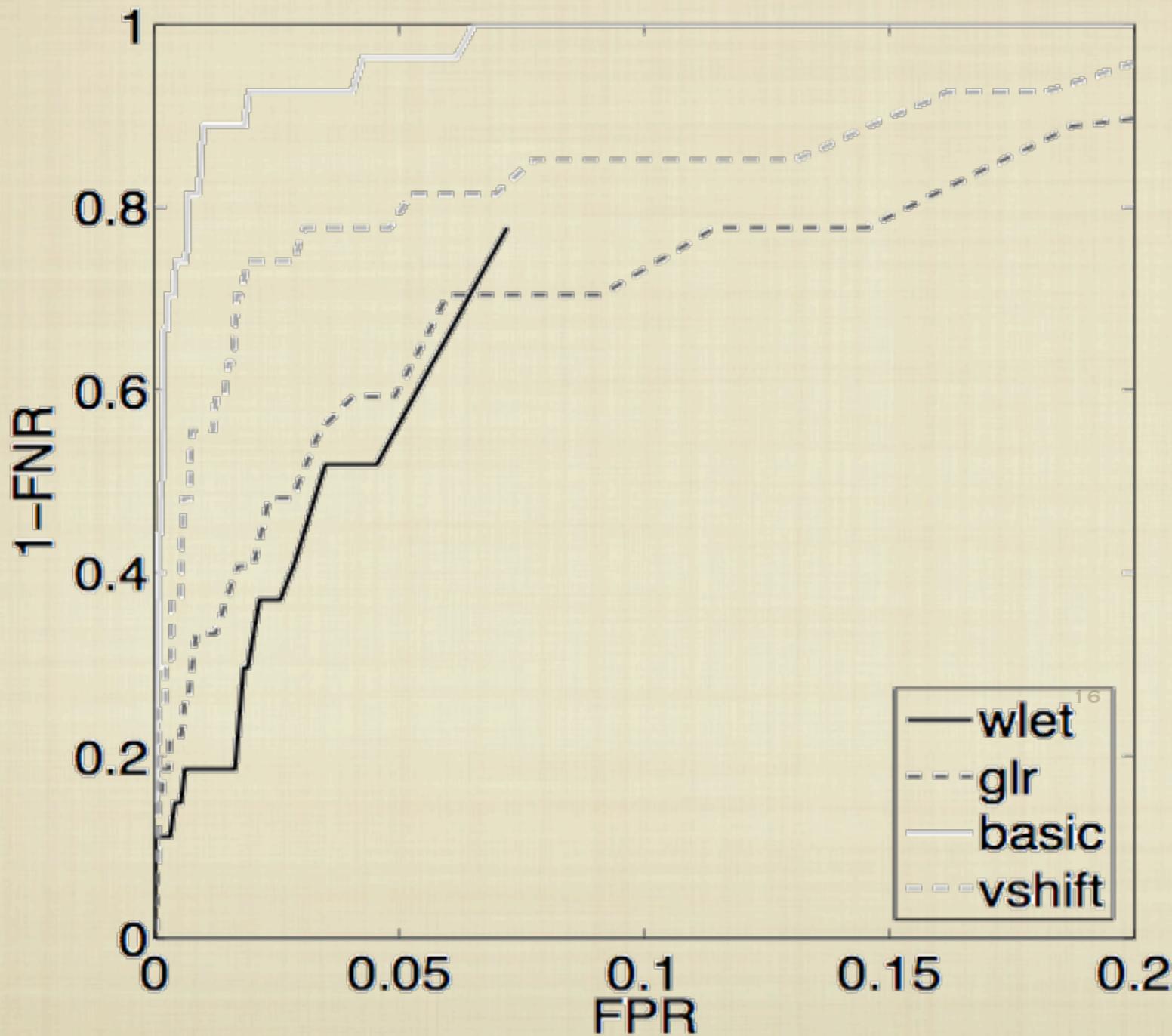


# INNOVATION AFTER RECALIBRATION

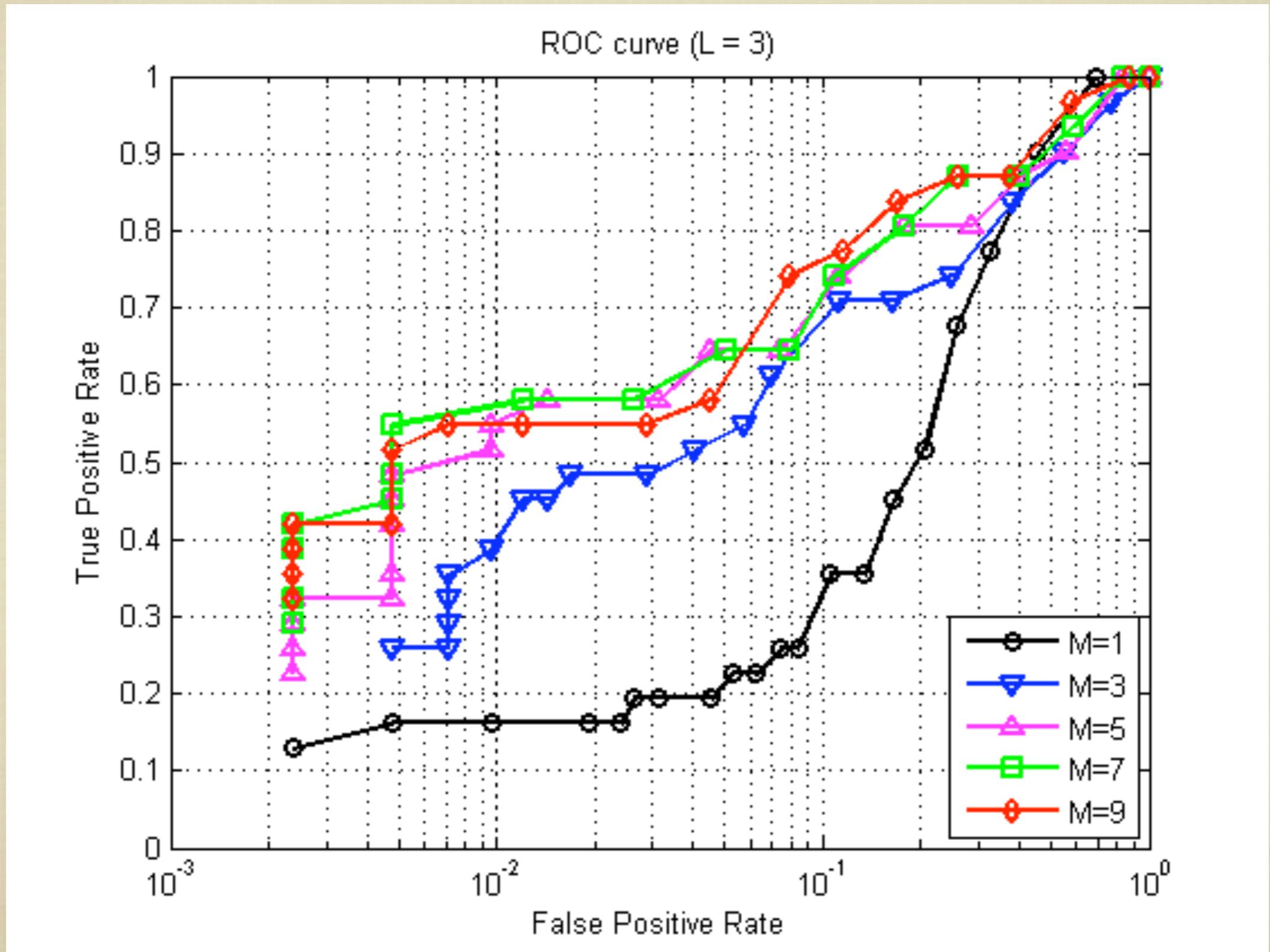


# ROC CURVE

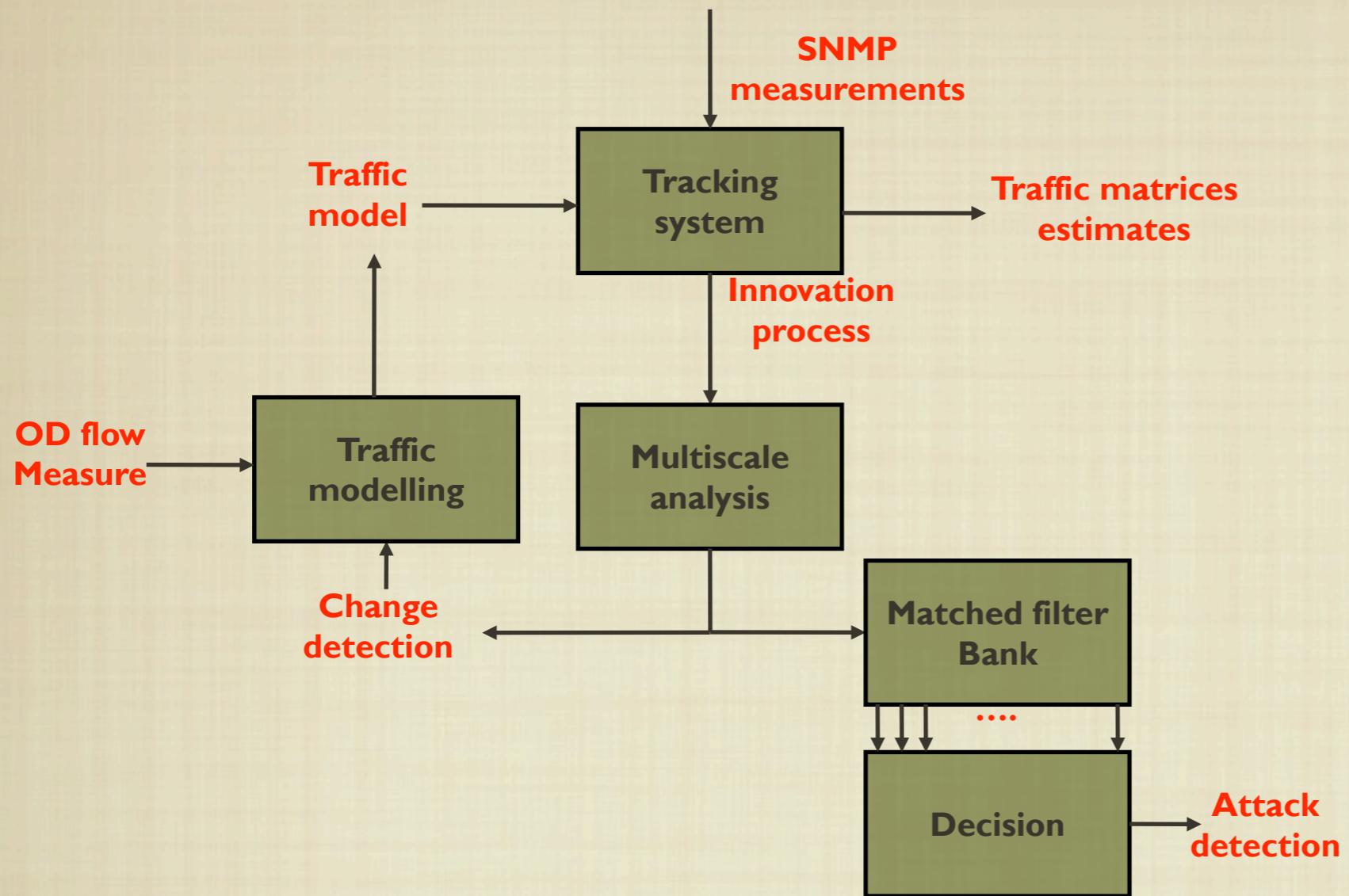
Receive Operating Characteristics (ROC) Curve



# ROC CURVE PCA



# ANOMALY



17

# NON PARAMETRIC TECHNIQUES

# ANOMALY DETECTION

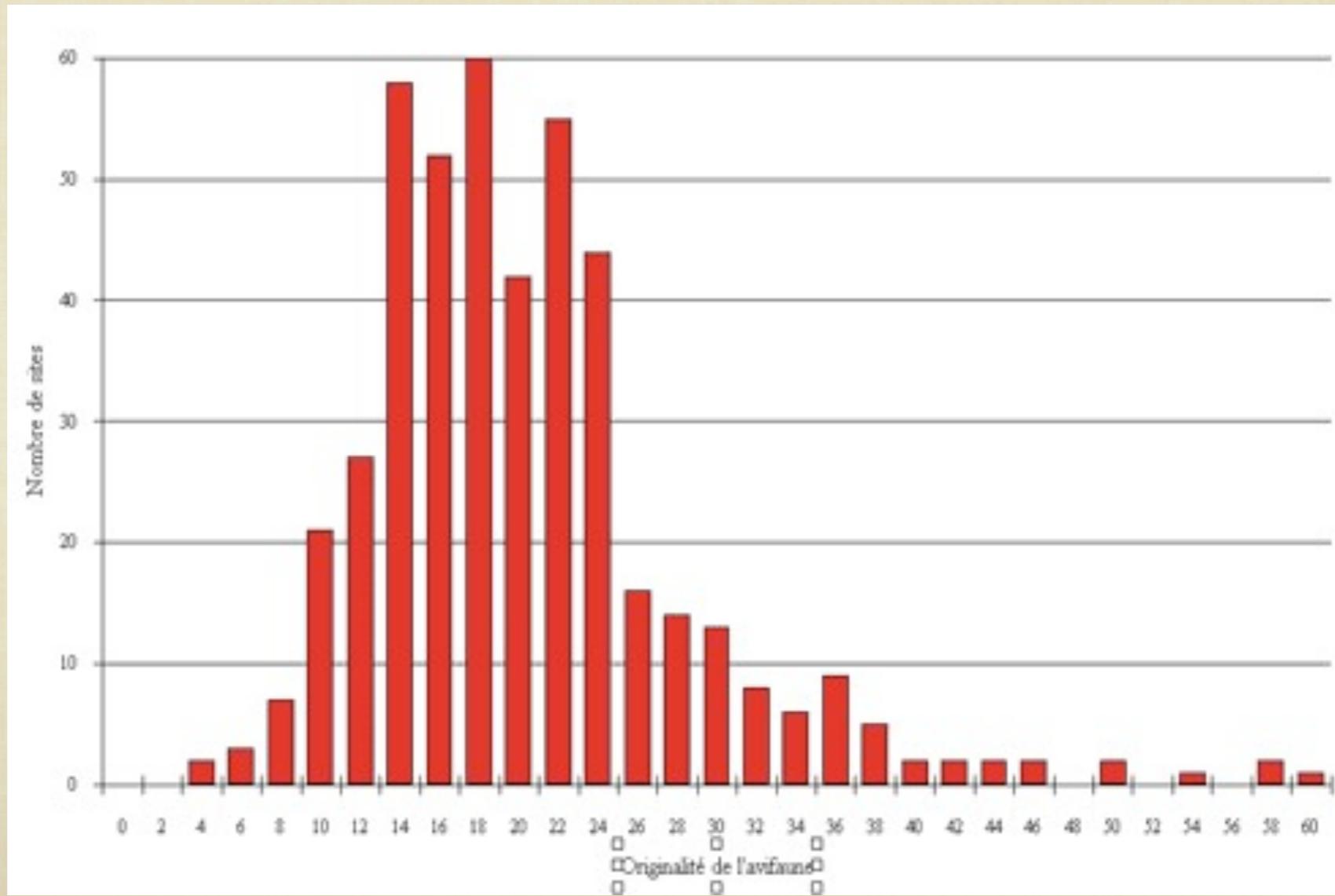
- THE NORMAL BEHAVIOUR OF A NETWORK SHOULD BE ASSESSED AS A GLOBAL JOINT DISTRIBUTION
  - AN ANOMALY WILL HAVE EFFECTS AT DIFFERENT PLACE OF THE NETWORK
- BASIC UNDERLYING QUESTION HOW TO INFER THE JOINT DISTRIBUTION IN A DISTRIBUTED WAY ?
  - IN AN EFFICIENT WAY

# DISTRIBUTED DENSITY

- **NAIVE, COSTLY SOLUTION**
  - COLLECT ALL DATA SOURCES TO A CENTRAL ANALYSIS POINT. THEN PERFORM DENSITY ESTIMATION
  - INFEASIBLE ESPECIALLY IF THE PROCEDURE HAS TO BE REPEATED PERIODICALLY
- **DISTRIBUTED DENSITY ESTIMATION**
  - EACH DATA SOURCE MAKES A LOCAL ANALYSIS
  - SEND THE RESULTS TO A CENTRAL POINT
    - AGGREGATION OF LOCAL ANALYSES TO OBTAIN A GLOBAL ANALYSIS OF DATA

# CLASSICAL DISTRIBUTION ESTIMATION

- SIMPLEST DENSITY ESTIMATOR : MULTIDIMENSIONAL HISTOGRAM

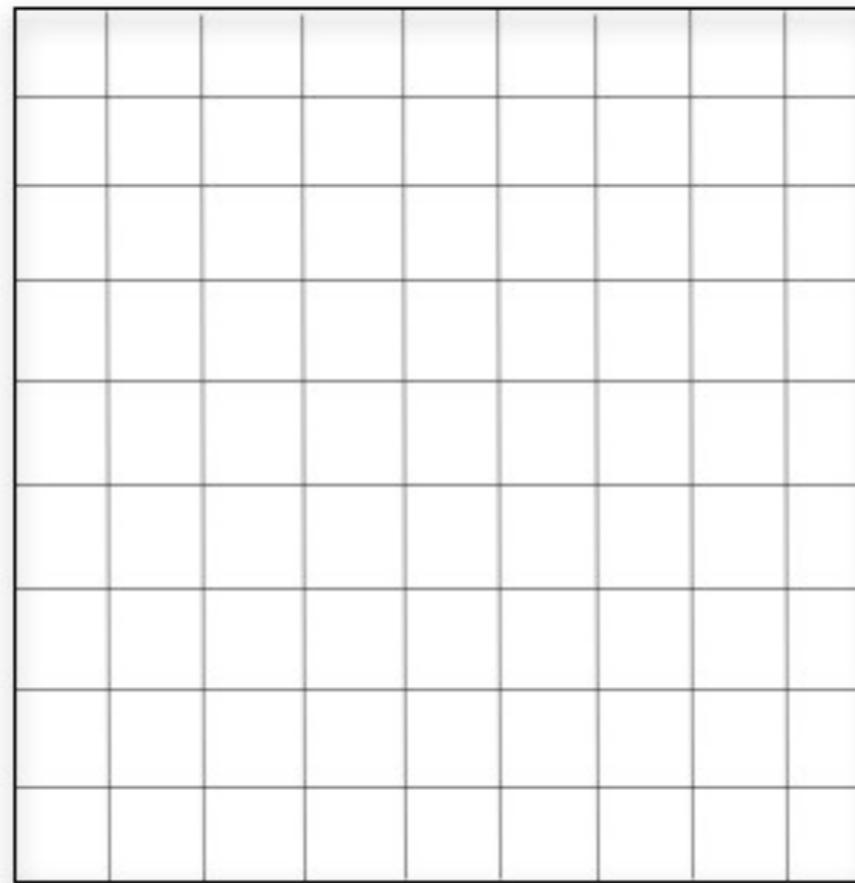


# RECTANGULAR GRID ISSUES

- **QUALITY OF THE BINNING DEPENDS OF**
  - **BIN SIZE**
  - **ORIGIN**
  - **GRID ORIENTATION**
- **FOLLOWING PRINCIPAL DIRECTIONS**
- **A SINGLE GRID COULD EITHER BE QUITE GOOD OR QUITE BAD, DEPENDING HOW MUCH IT IS ORIENTED LIKE THE DATA.**

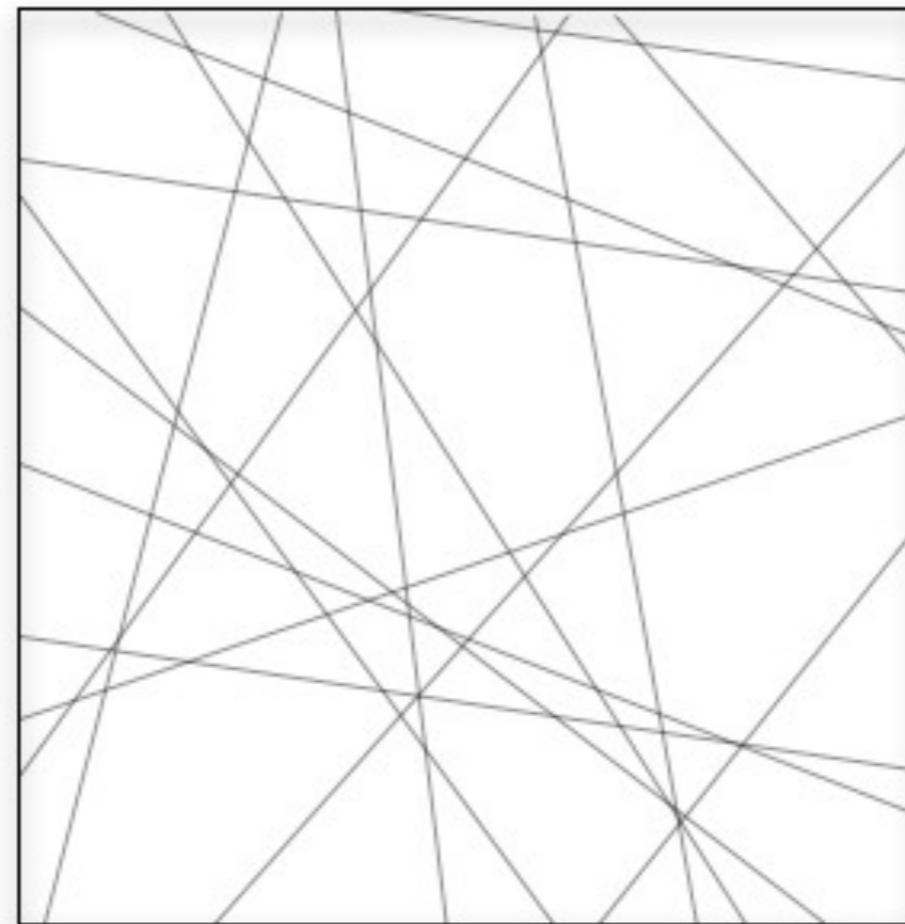
# MOSAIC-BASED HISTOGRAM

- A RECTANGULAR BINNING CAN BE VIEWED AS A SET OF REGULARLY-SPACED HYPERPLANES ON EACH DIMENSION, EACH BIN BEING AN HYPER-RECTANGLE DELIMITED BY HYPERPLANES



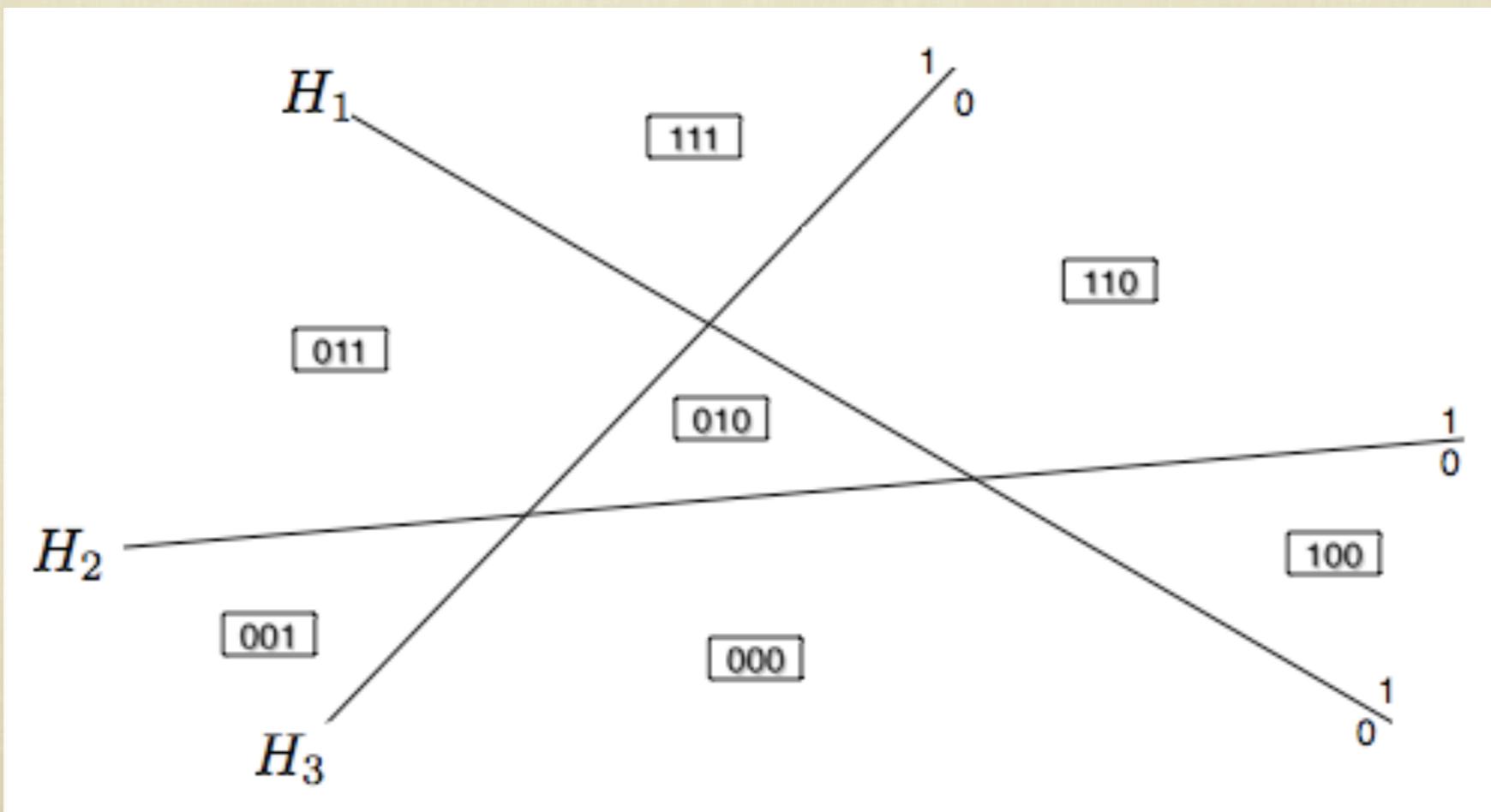
# RANDOM MOSAIC

- THROW RANDOMLY THE SAME NUMBER  $M$  OF HYPERPLANES TO OBTAIN A MOSAIC COVERING THE DATA SPACE



# MOSAIC-BASED BINNING

- WE CAN ASSOCIATE A BINARY NUMBER TO EACH OF THE REGIONS (CONVEX POLYTOPES) IN A NATURAL WAY



# MOSAIC-BASED BINNING

- THE BINNING  $Y_i$  OF THE DATA SAMPLE  $X_i$  ON EACH LOCAL SITE IS MADE VERY SIMPLY

$$Y_i = \operatorname{sgn}(A \cdot X_i - b)$$

WHERE A AND B ARE RANDOM MATRICES.

- ONLY NEED TO TRANSMIT NUMBER OF HYPERPLANES AND RANDOM SEED TO EACH LOCAL SITE !

# MOSAIC-BASED BINNING

- THE BINARY CODEWORDS ARE TRANSMITTED TO THE CENTRAL NODE
- CENTRAL NODE MUST COMPUTE THE VOLUME OF EACH NON-EMPTY POLYTOPE OF THE TESSELLATION
- FINALLY WE CAN BUILD THE HISTOGRAM

# DENSITY ESTIMATION

- LET THE SYSTEM OF CELLS BE  $\{C_j\}$  THEN THE ESTIMATE AT A POINT  $t \in \mathbb{R}^d$  IS

$$\hat{f}(t) = \frac{\mathbf{1}(X \in C_j)}{N_X \cdot Vol(C_j)}, \quad t \in C_j$$

# BACK TO ANOMALY DETECTION

- THIS PERFORMS BETTER THAN A FIXED GRID WHEN THE LOCAL DISTRIBUTIONS DIFFER
- ANOMALY DETECTION
  - FIX A MOSAIC
  - PERIODICALLY UPDATE THE DENSITY ESTIMATE  $f_t(\cdot)$
  - IF THE CURRENT DENSITY ESTIMATE  $f_0(\cdot)$  IS “TOO FAR” FROM THE DENSITY OF THE “NORMAL” STATE  
    => ANOMALY

# DESIRED QUALITIES OF THE ESTIMATOR

- **QUALITY OF THE GLOBAL AGGREGATE STRAIGHTLY DEPEND ON THE QUALITY OF LOCAL ANALYSES**
- **LOCAL ANALYSES : QUALITY / SIZE TRADEOFF**
- **A GOOD DISTRIBUTED ESTIMATOR MUST ALSO ALLOW TUNABILITY OF THIS TRADEOFF**

# ANOMALY DETECTION

- “TOO FAR” MAY BE THE QUADRATIC ERROR BETWEEN THE DISTRIBUTIONS

$$\int (f_t(u) - f_0(u))^2 du$$

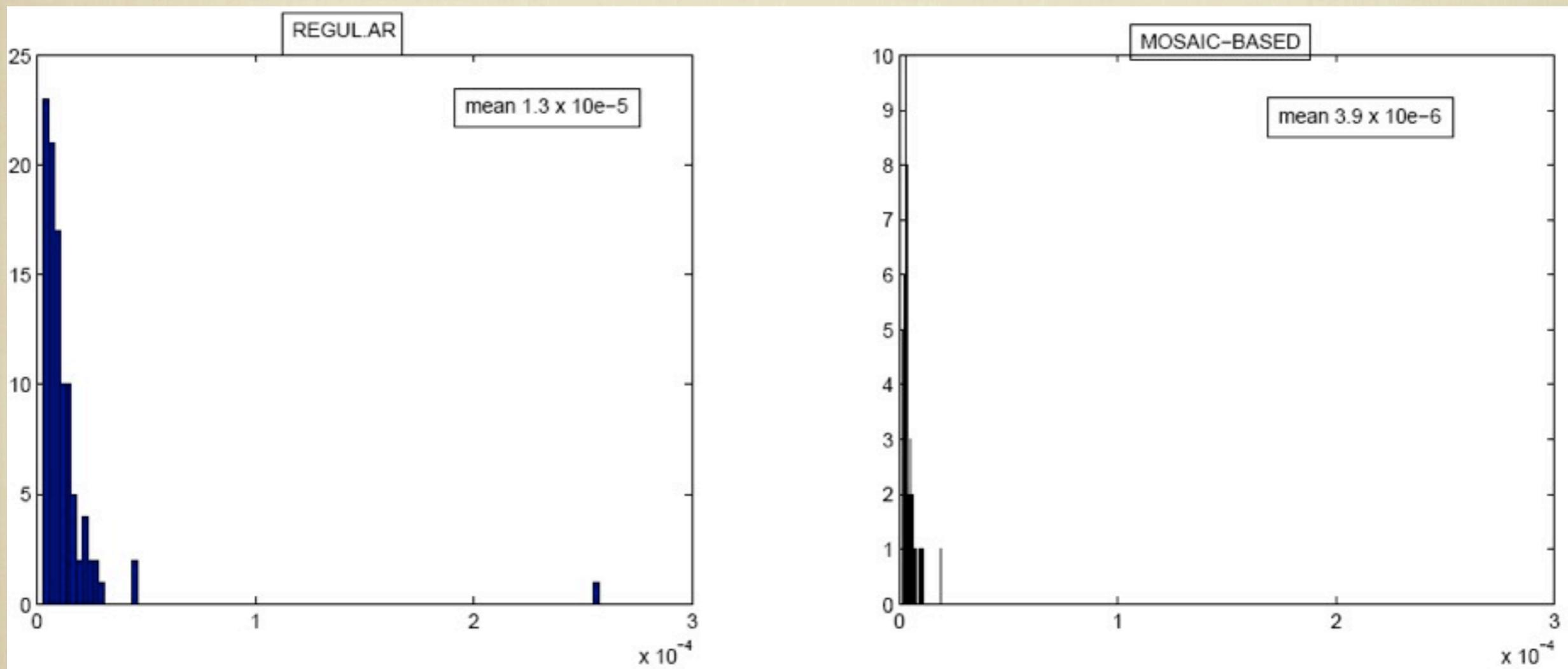
- BUT THIS COULD BE RATHER COSTLY TO EVALUATE

# ANOMALY DETECTION

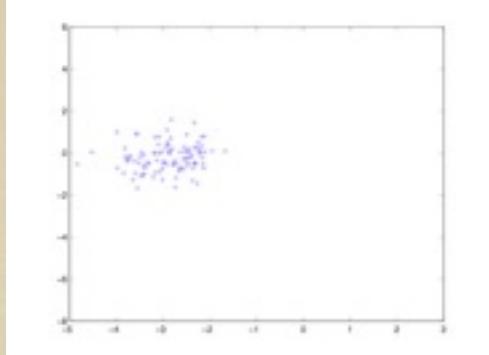
- KULLBACK-LEIBLER DISTANCE IS EASIER TO EVALUATE

$$D(f_0 \parallel f_t) = \sum_{\text{cells } C_j} f_0(C_j) \cdot \log \left( \frac{f_0(C_j)}{f_t(C_j)} \right)$$

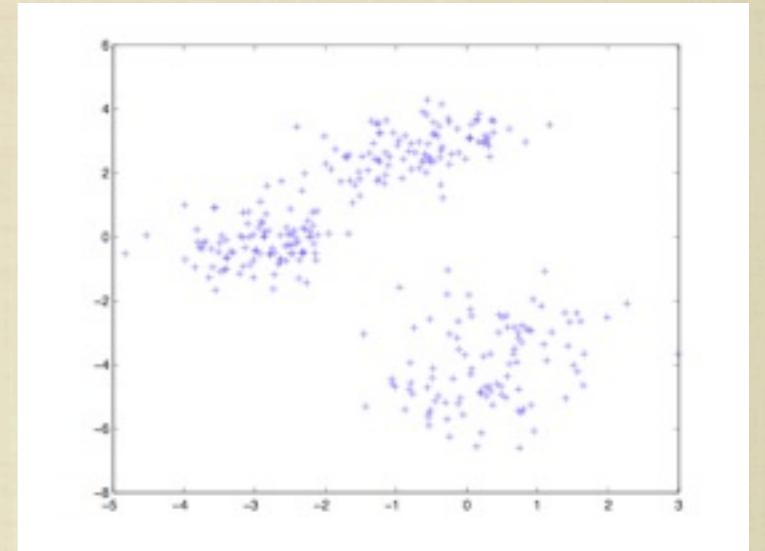
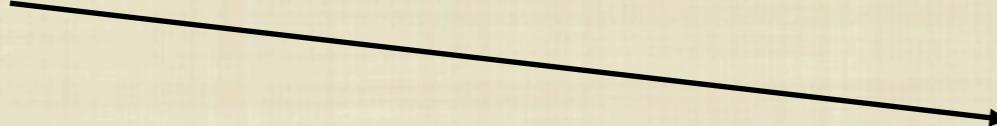
SINCE WE CAN TAKE IT AS A DISCRETE DISTRIBUTION  
OVER THE BINS.



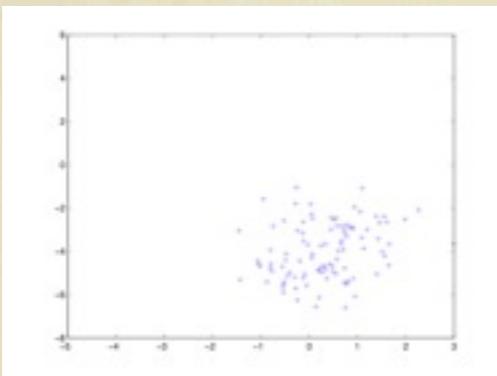
# GLOBAL DISTRIBUTION



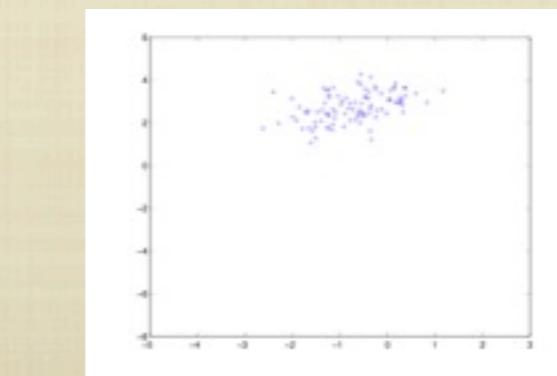
$f_1(x), \frac{N_1}{N}$  points



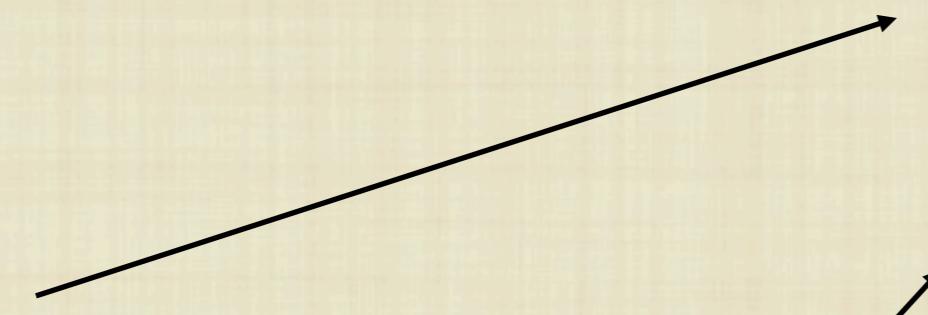
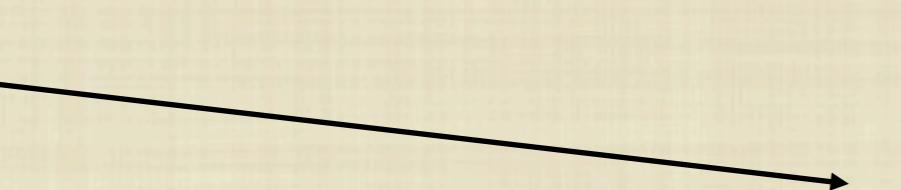
$$f(x) = \sum_i \frac{N_i}{N} f_i(x)$$



$f_2(x), \frac{N_2}{N}$  points

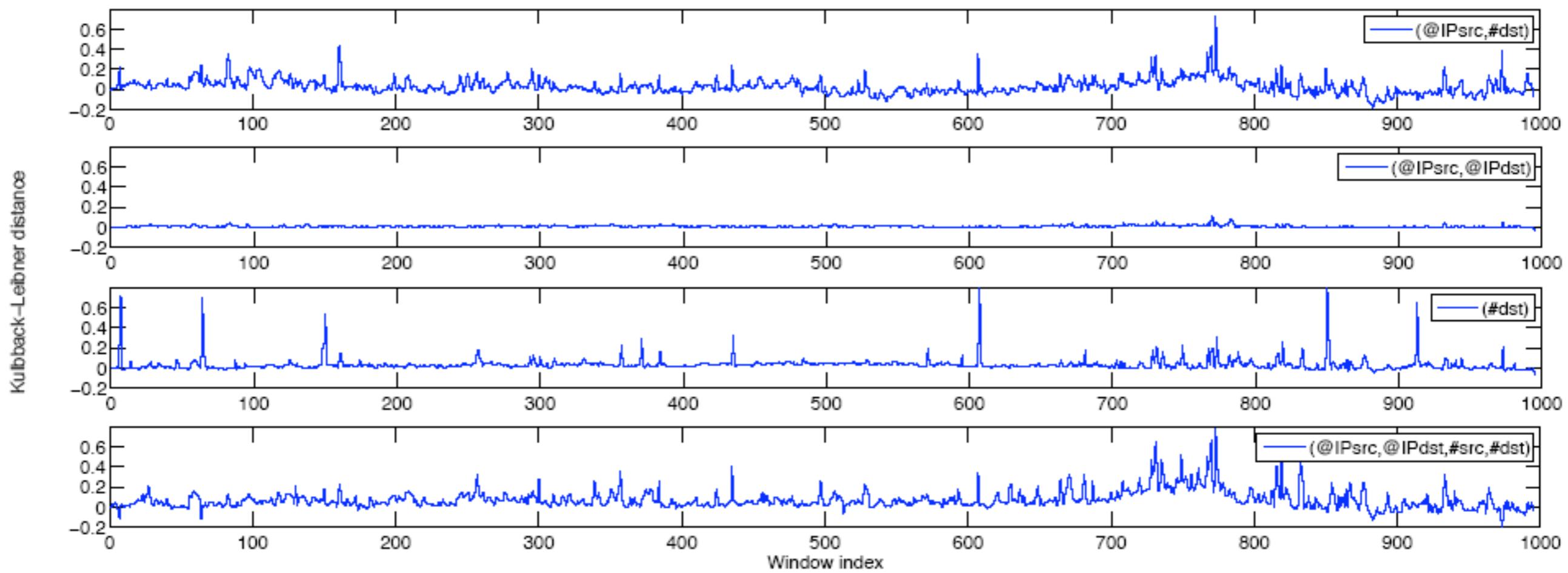


$f_3(x), \frac{N_3}{N}$  points



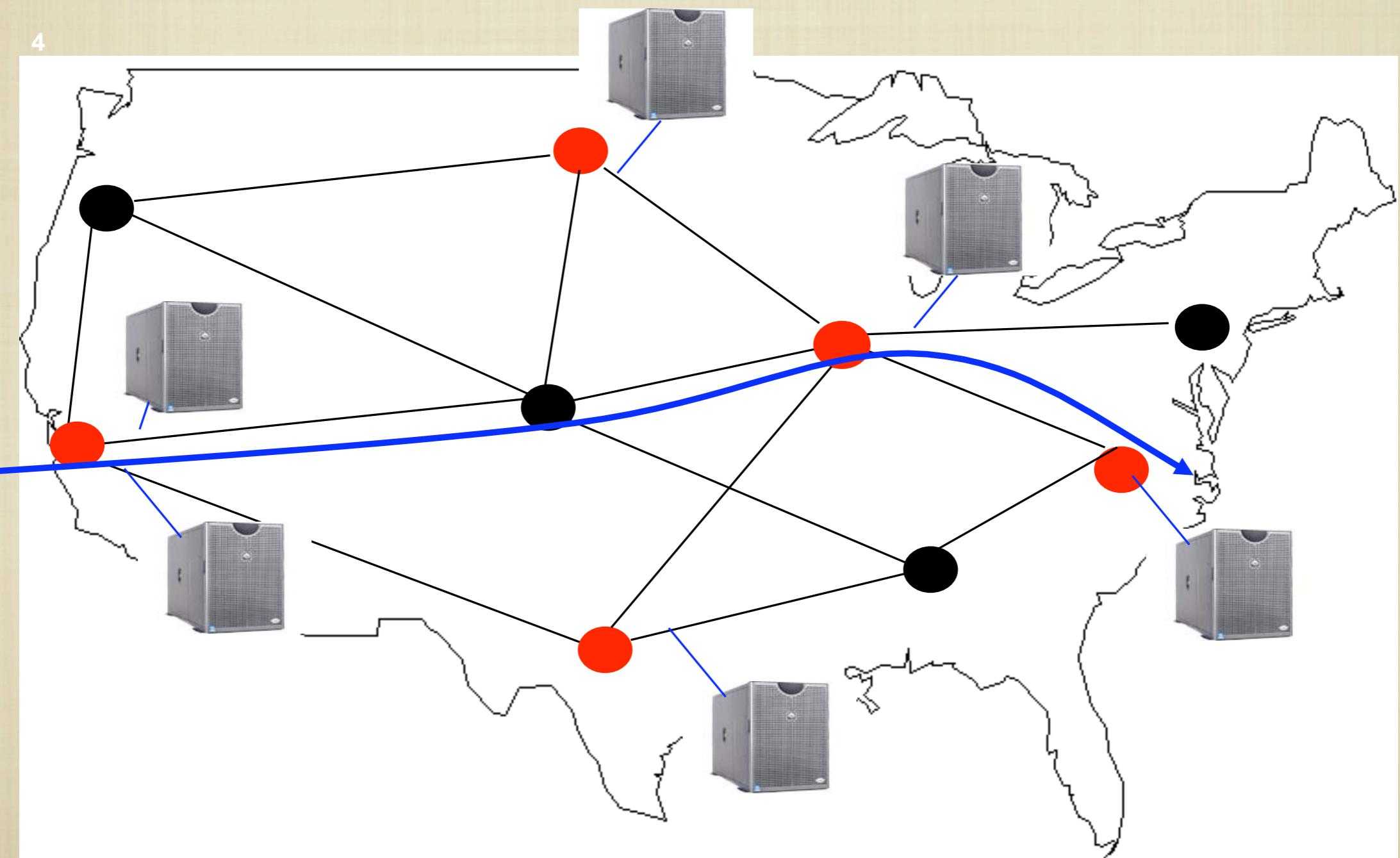
# NON PARAMETRIC AD

Variation of Kullback–Leibnner distance compared to the first window for a full day



# DISTRIBUTED ANOMALY DETECTION

# DISTRIBUTED ANOMALY DETECTION



# DISTRIBUTED ANOMALY DETECTION !

- SAME AS CENTRALIZED
  - ALL THE STEPS SHOULD BE RUN IN A DISTRIBUTED WAY
    - DISTRIBUTED MODEL IDENTIFICATION
    - DISTRIBUTED FILTERING
    - DISTRIBUTED DECISION
  - FOR PARAMETRIC AND NON PARAMETRIC APPROACHES !
    - PARAMETRIC : STEPS 1 AND 2 DISTRIBUTED, STEP 3 LOCAL
    - NON PARAMETRIC : STEPS 1 AND 2 LOCAL, STEP 3 DISTRIBUTED

# PARAMETRIC VS. NON PARAMETRIC

- PARAMETRIC CASE
  - KNOWING THE AUTOCOVARIANCE WHAT ARE THE BEST PROJECTIONS
    - LOCAL KLT GIVES BEST LOCAL PROJECTIONS
    - NOT OPTIMAL
  - DISTRIBUTED KLT
    - SOLVE GLOBAL OPTIMIZATION BY AN ITERATIVE ALGORITHM
- NON PARAMETRIC CASE
  - USE RANDOM PROJECTION
    - IF ENOUGH USED CONVERGES TO KLT
      - ASYMPTOTICALLY OPTIMAL
    - COMPRESSED SENSING IDEA

# DISTRIBUTED KLT

- 4 NODES WITH 10 OBSERVATIONS EACH
- EXCHANGING 2 PROJECTIONS EACH
- JOINT KLT MSE = 8.5
- DISTRIBUTED KLT MSE = 15.2
- 90% OF AD COINCIDENCE

