

# **Power and Telecom Management Networks: Real-time Anomaly Detection and Correlation**

**Simin Nadjm-Tehrani**

[www.ida.liu.se/~rtslab](http://www.ida.liu.se/~rtslab)

Department of Computer & Information Science

Linköping University, Sweden

and

University of Luxembourg

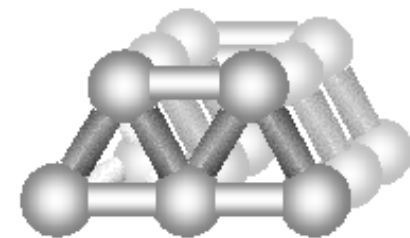


- First EU (IST) Critical Infrastructure Project
  - Outcomes and lessons learnt



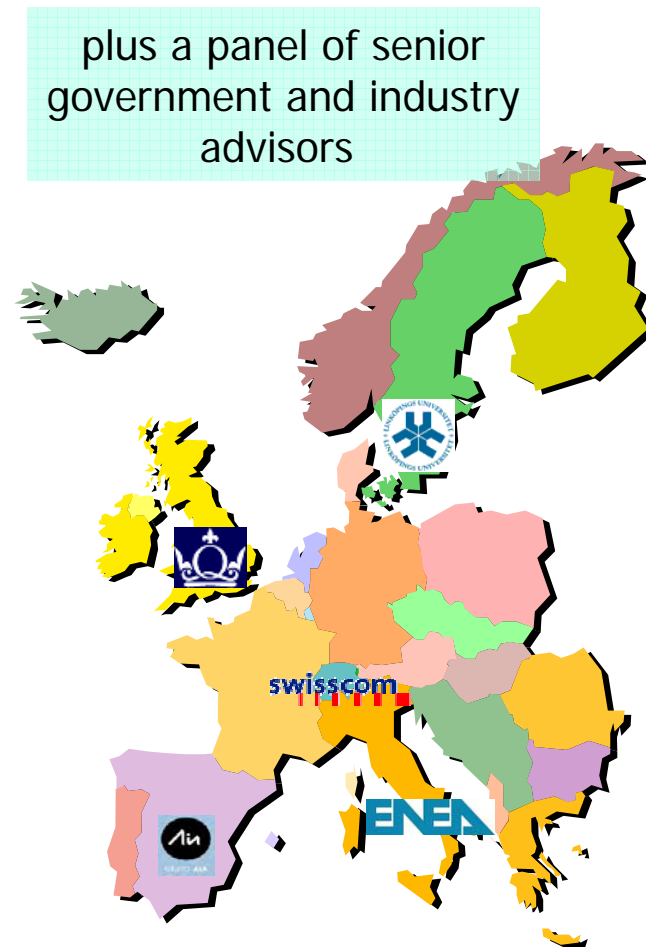
- CRIS: International Institute for Research on Critical Infrastructures

(will return to this tomorrow)



# Safeguard: 2001-2004

- Goal: to enhance survivability of Large Complex Critical Infrastructures (LCCIs)
- Electricity and telecommunications networks as practical examples
- Pre 9/11!



# Where to start?

---

- Power grid of today or tomorrow?
- Telecom of today or tomorrow?

# Restructuring of Power Grid

---

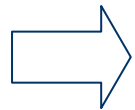
- Deregulation: organisations can enter into bilateral or multilateral power generation contracts
  - Large scale operation: from centralised to distributed control
  - Difficulty of coordination among independent service operators
- Approaching grid capacity
- New monitoring and control problems

# Telecom challenges

---

- Convergence of technologies
  - Everything is changing: services, business models, enabling technologies
- Internet dependability and security paramount to telecoms

## General:

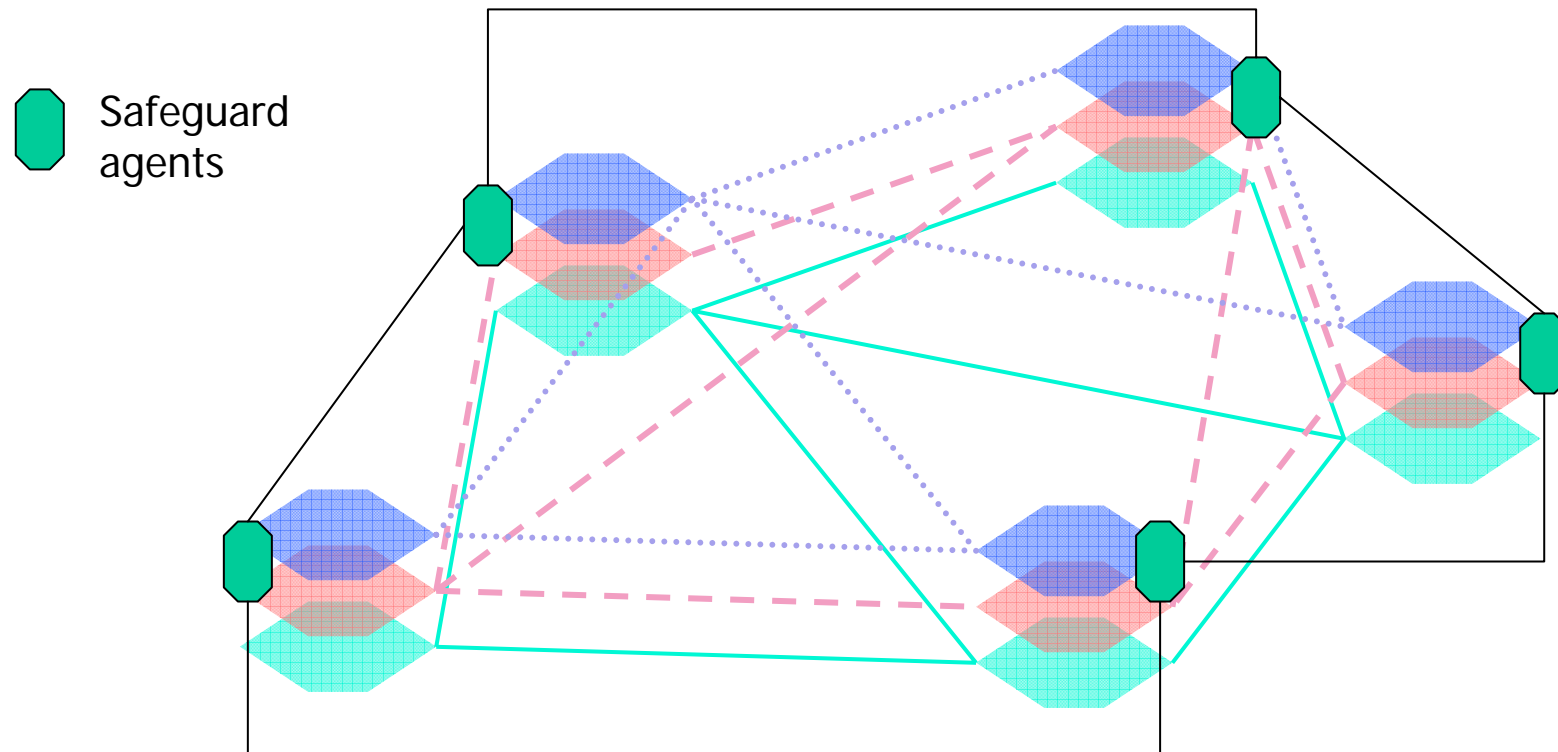


- Increase information quality for administrator
- Recognise unknown attacks
- Predict future overloads

## Telecom specific:

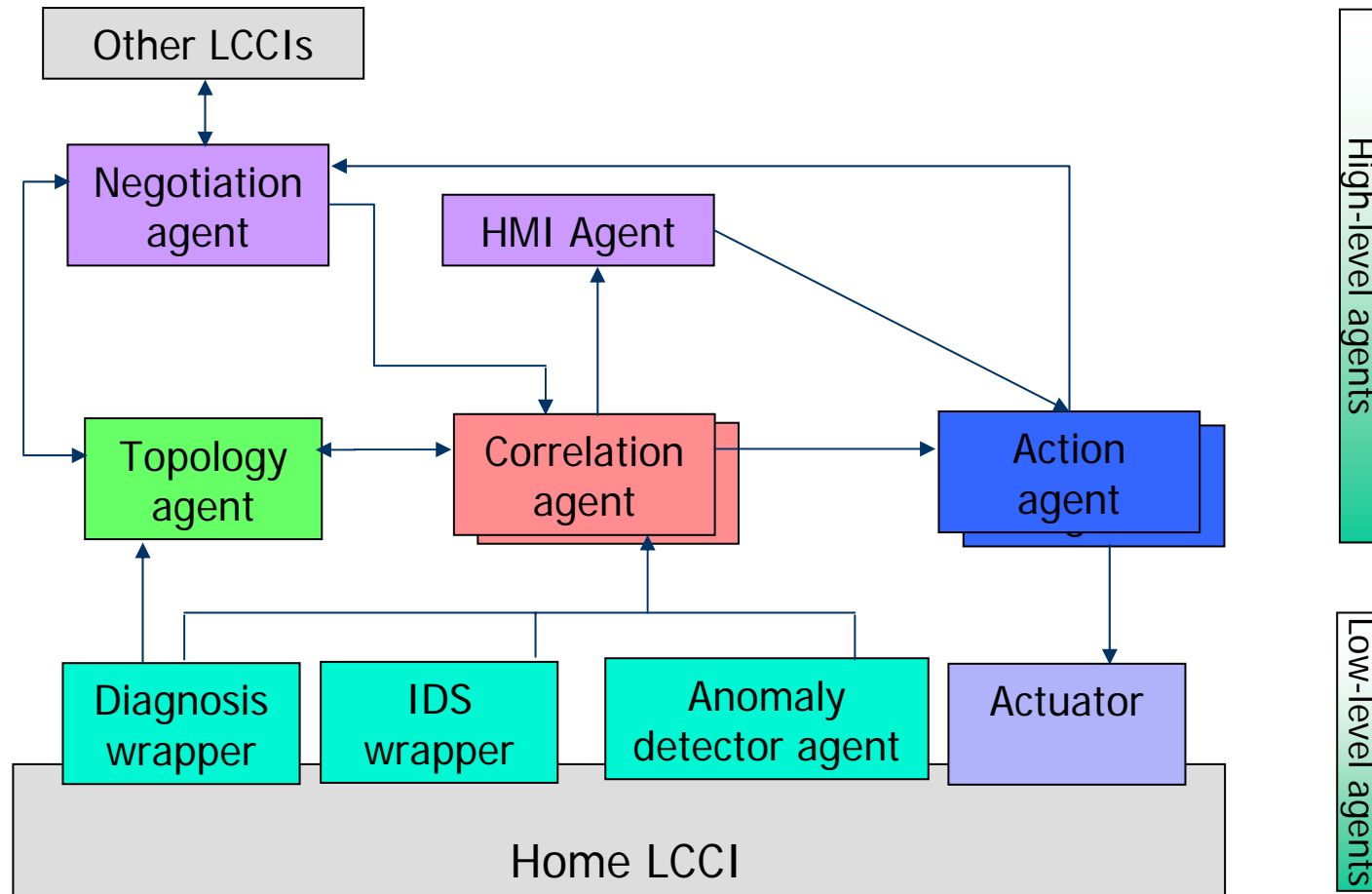
- Decrease no. of alarms
- Decrease false positives (higher availability)

# The Safeguard approach

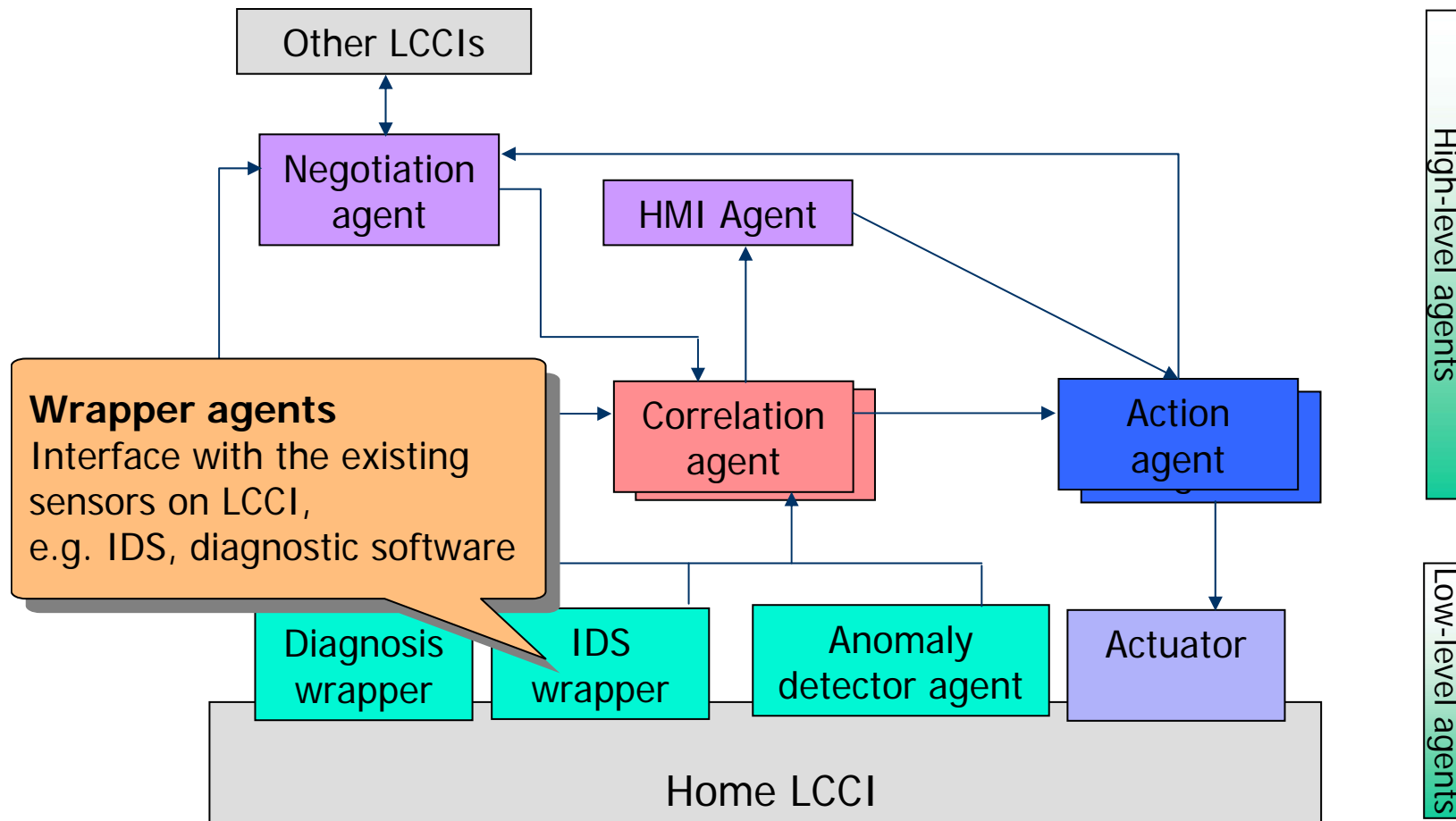




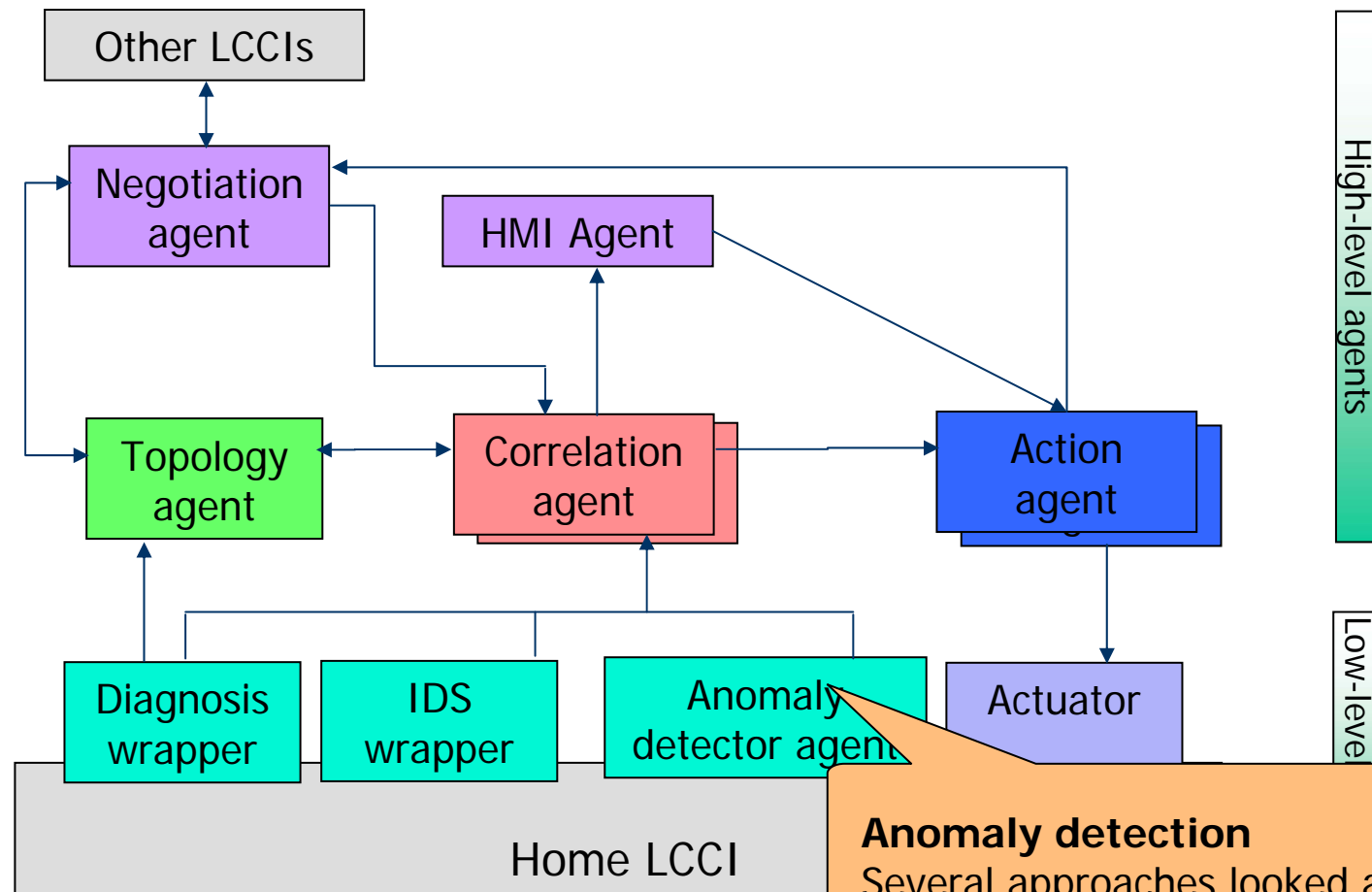
# Context: Safeguard architecture



# The Safeguard architecture



# The Safeguard architecture

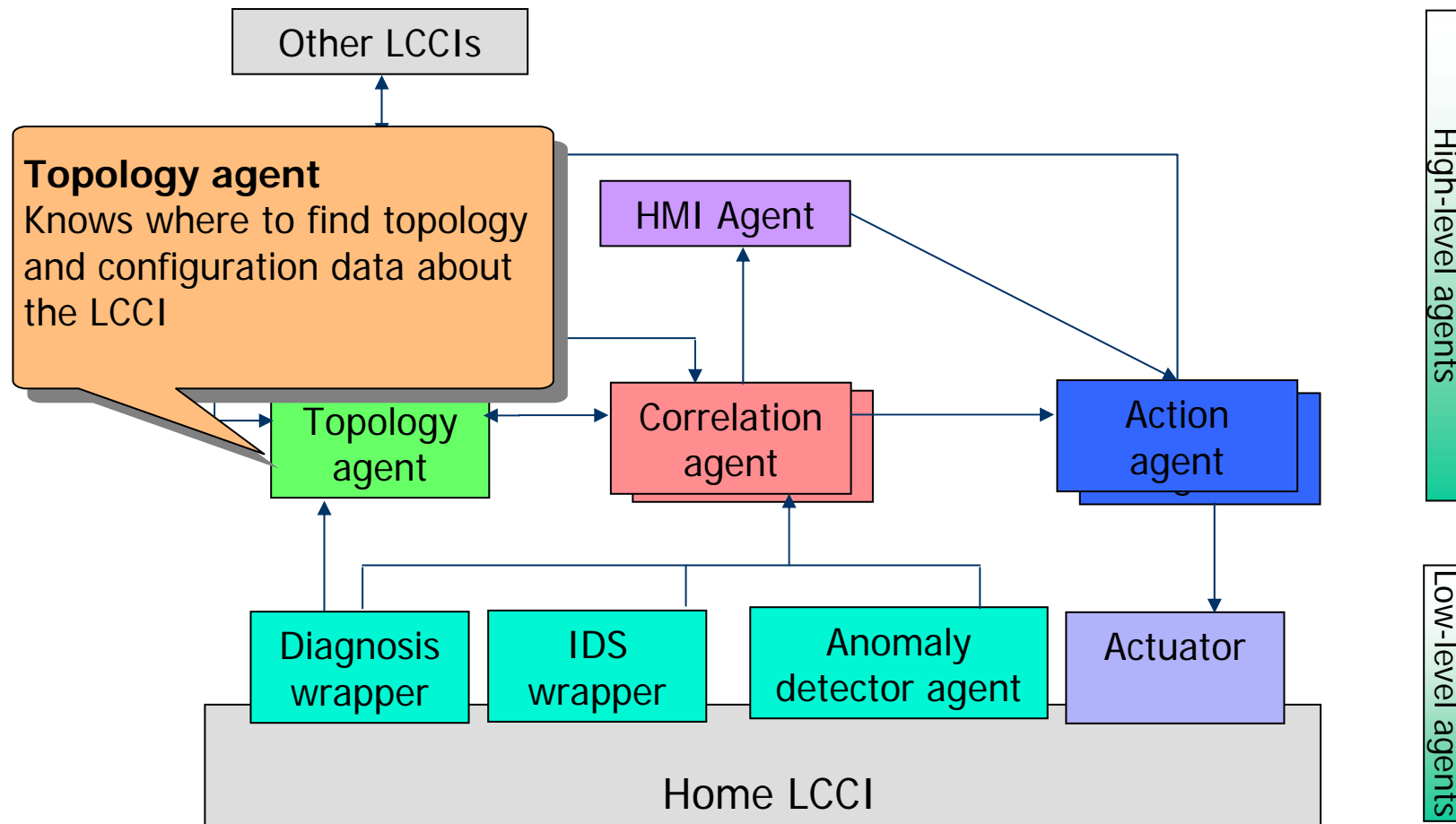


## Anomaly detection

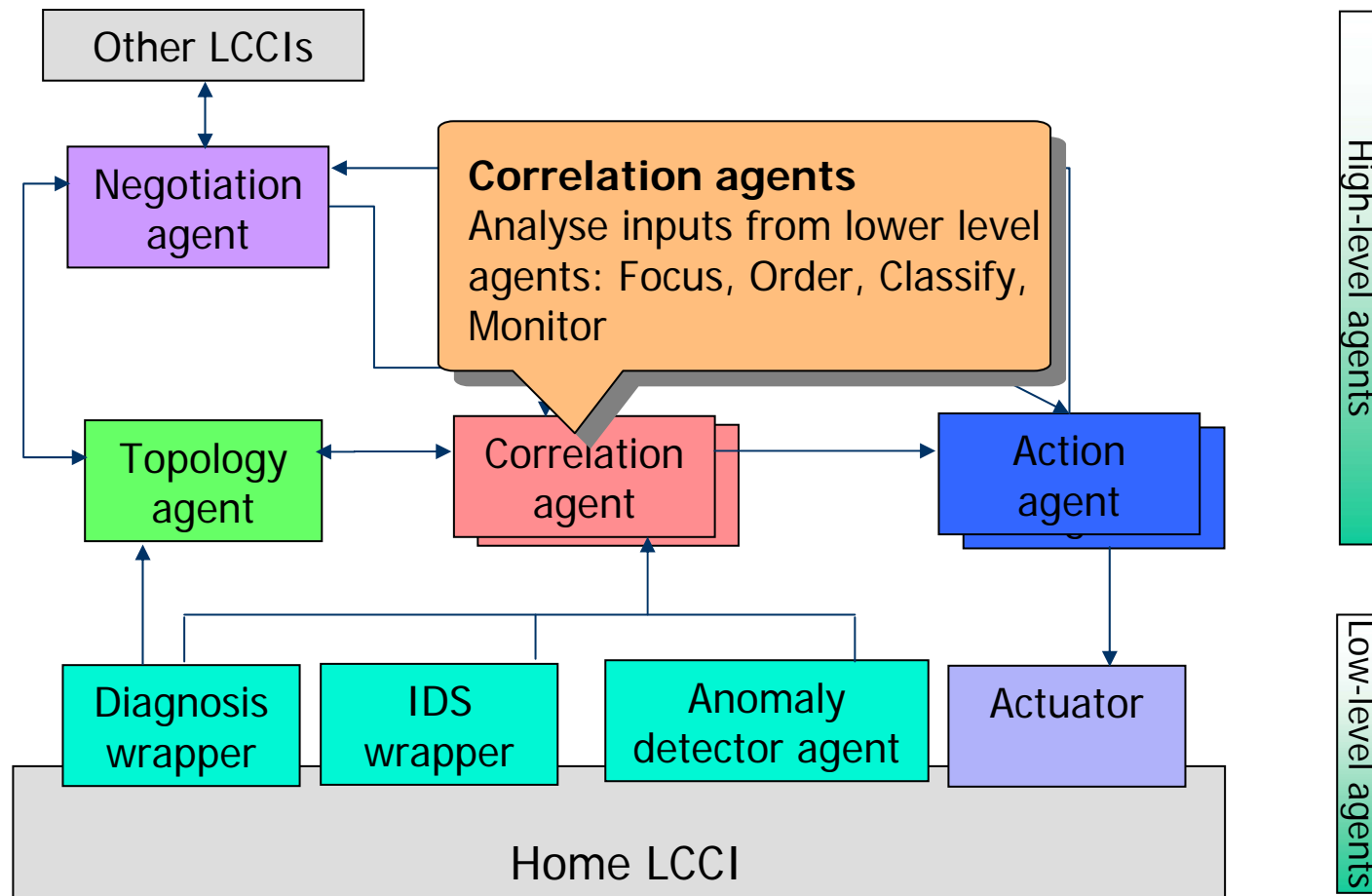
Several approaches looked at:

- Case Based Reasoning
- Clustering (ADWICE)
- Invariant detection

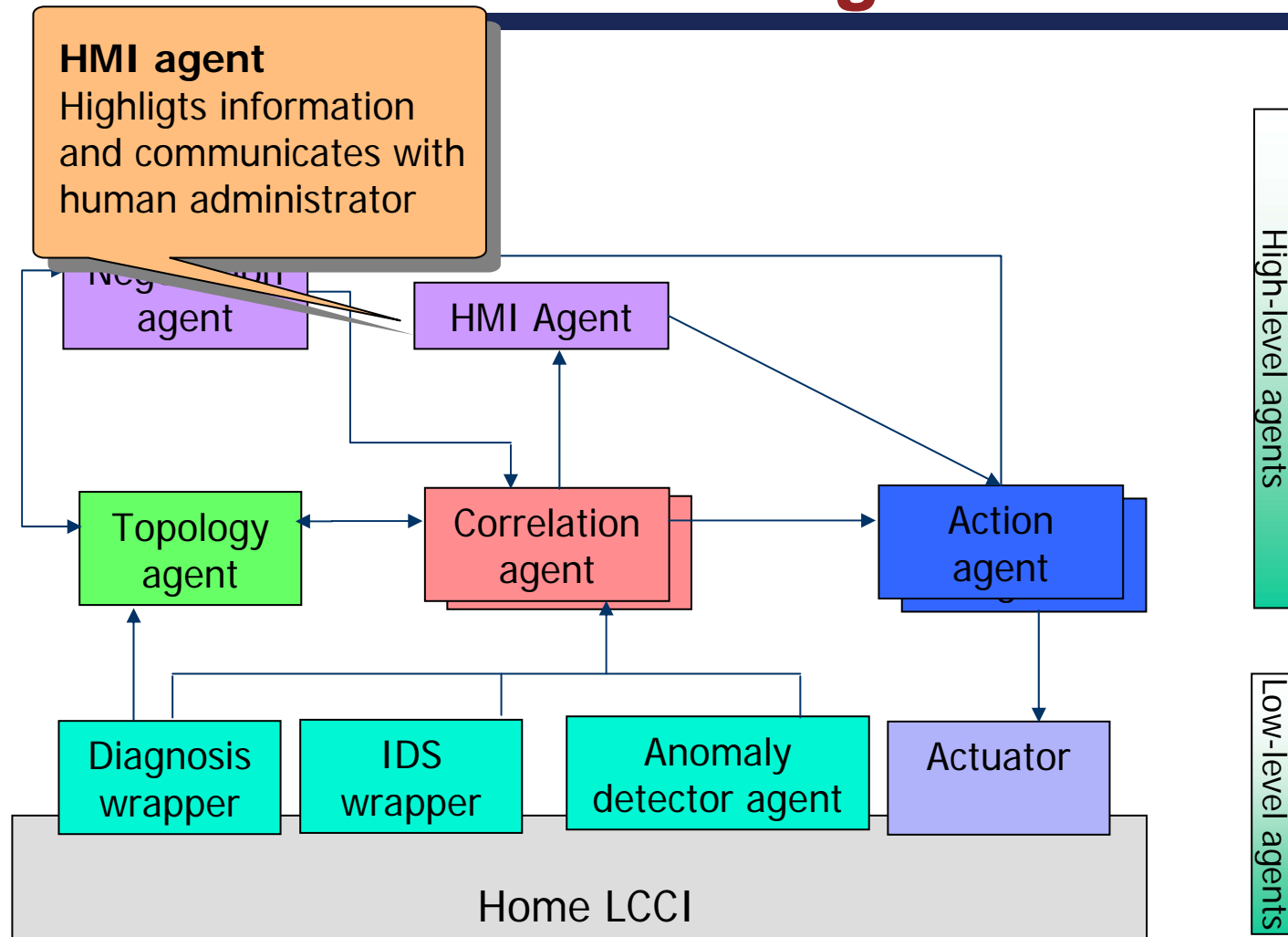
# The Safeguard architecture



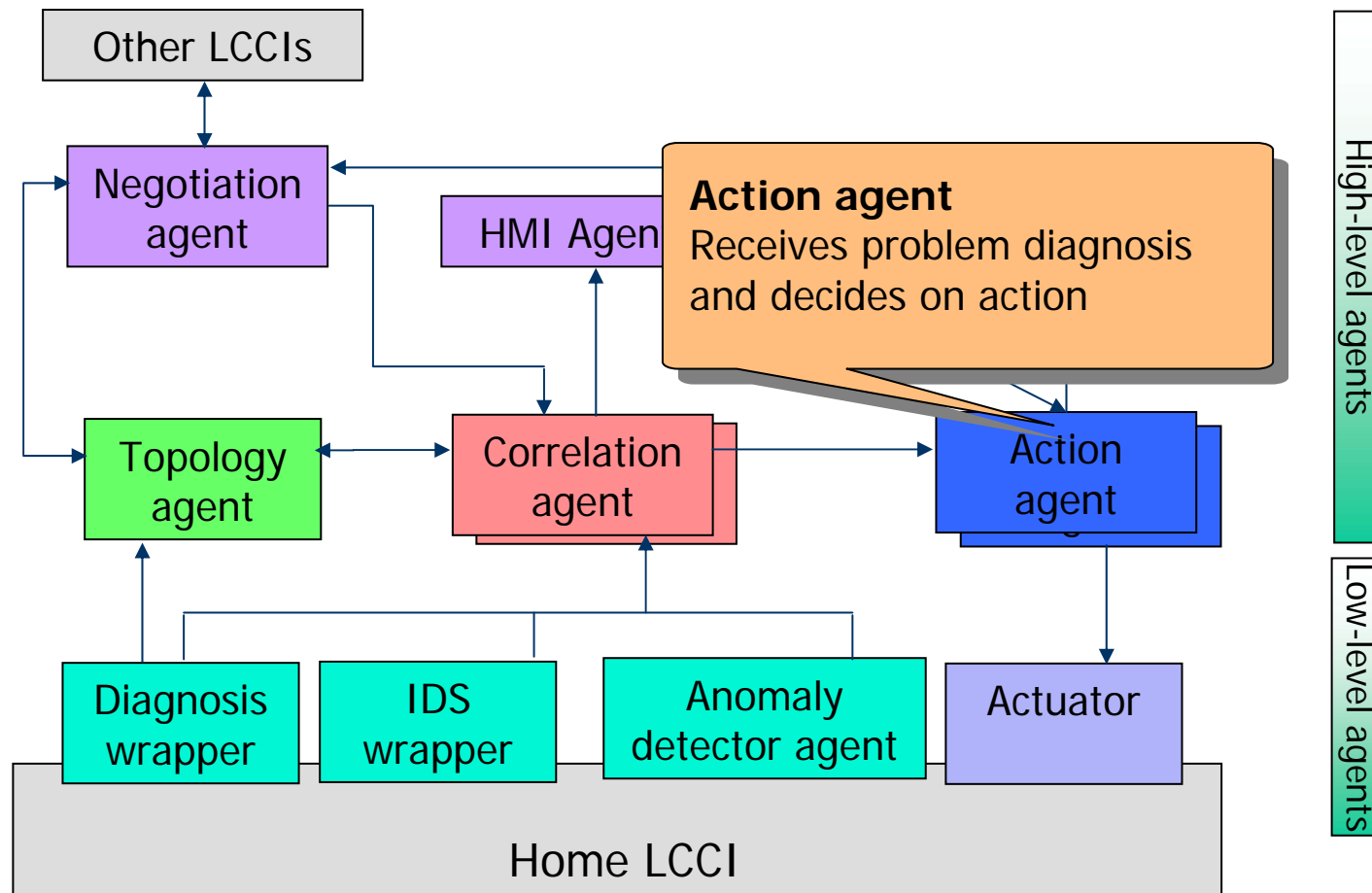
# The Safeguard architecture



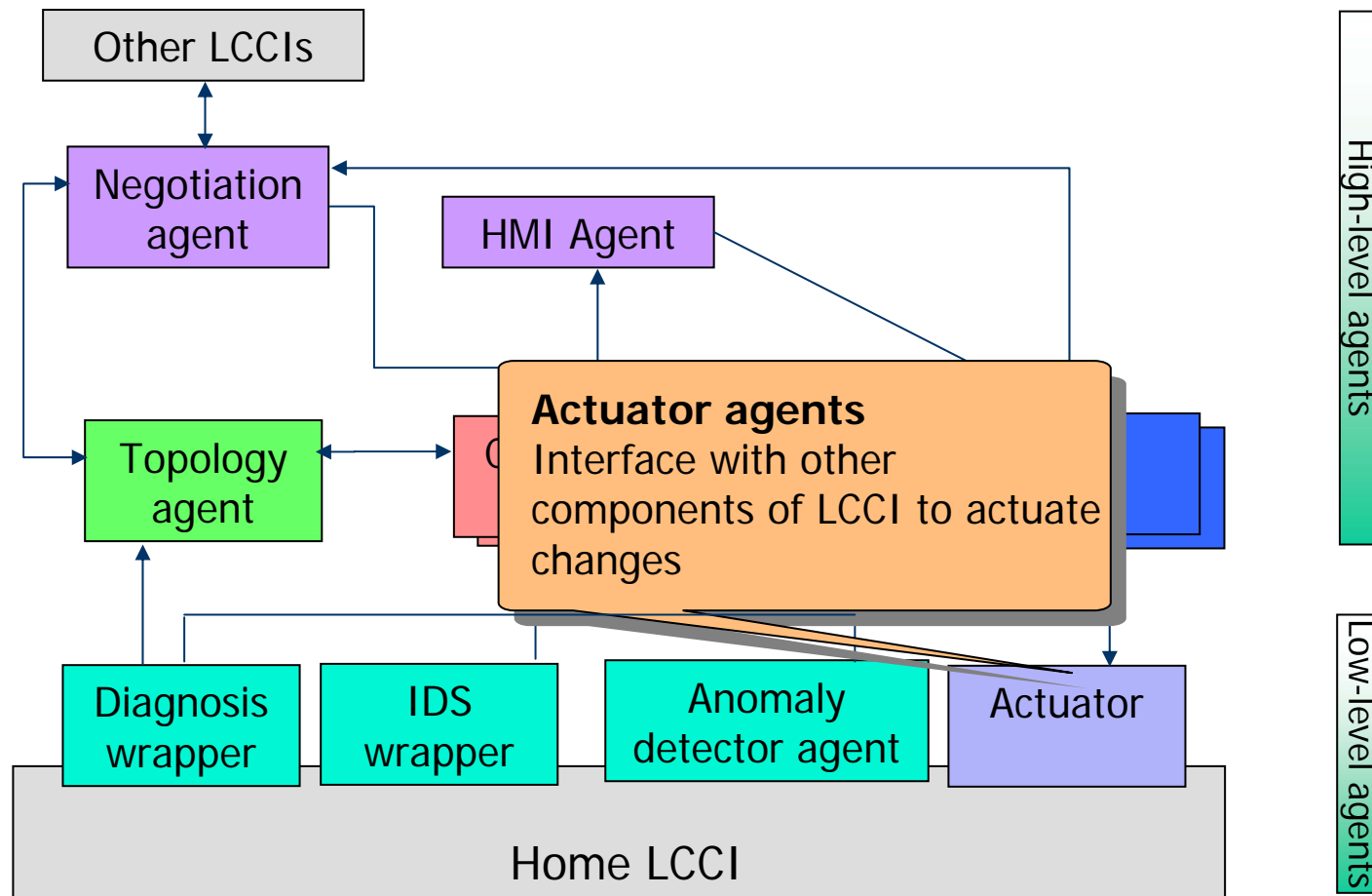
# The Safeguard architecture



# The Safeguard architecture

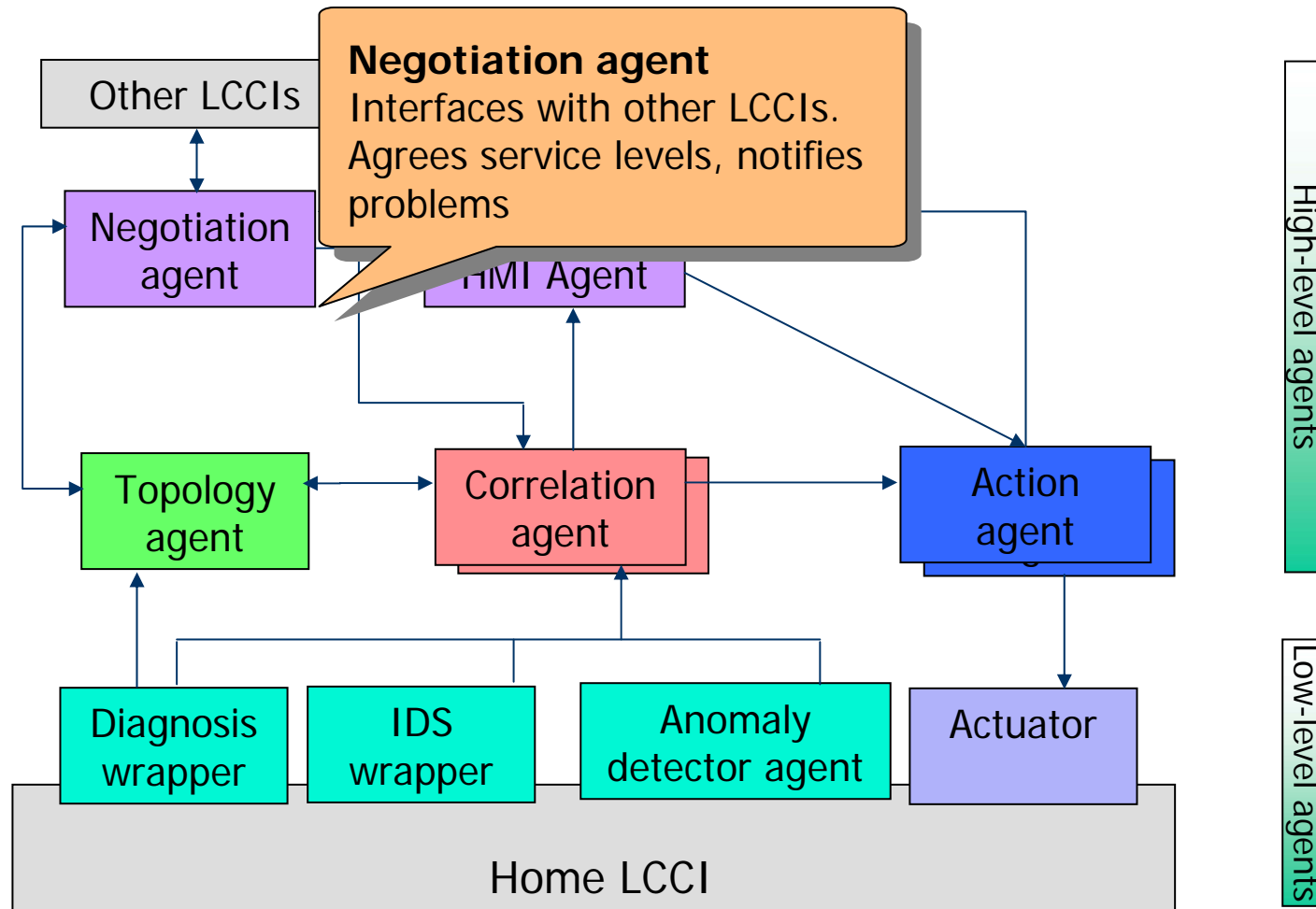


# The Safeguard architecture





# The Safeguard architecture

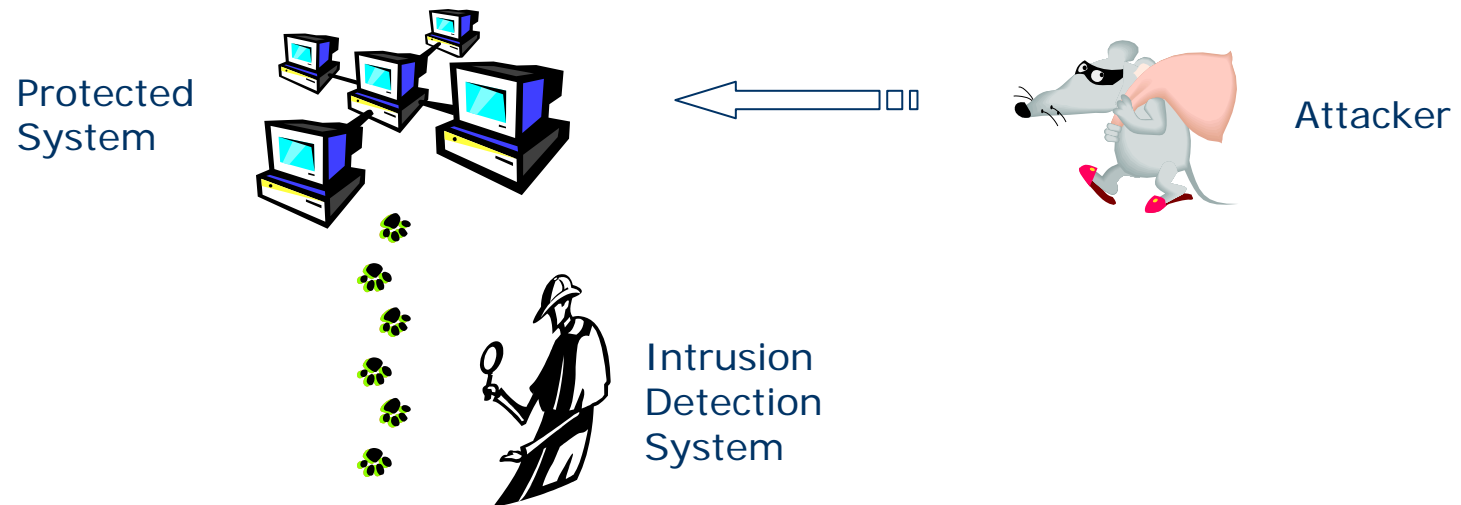


# Anomaly Detection

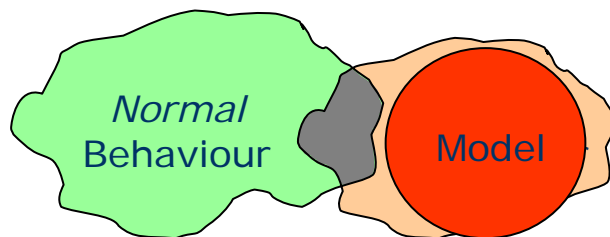
---

- ADWICE: Anomaly Detection With fast Incremental ClustEring
- Joint work with Kalle Burbeck
- Not a silver bullet: part of the larger Safeguard context

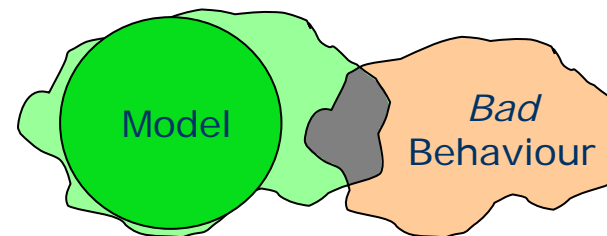
# Intrusion detection



## *Misuse Detection*

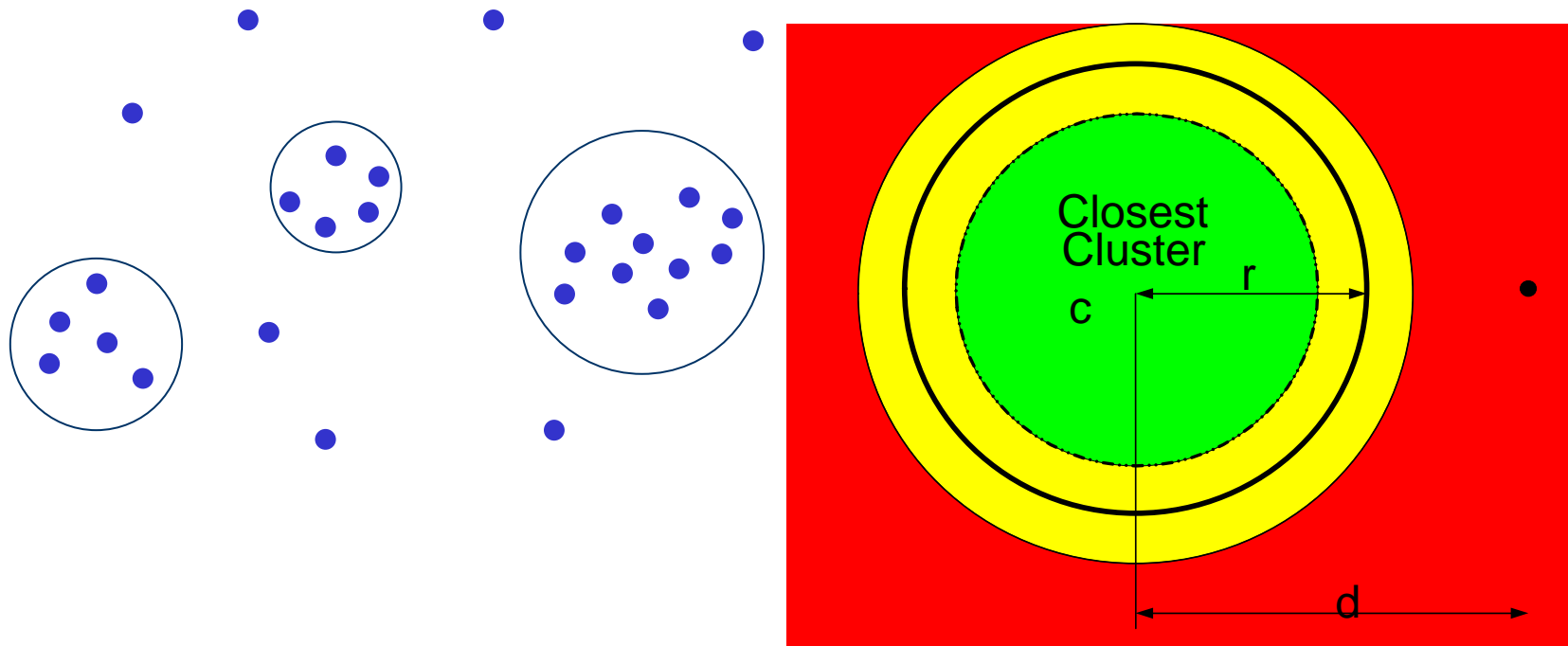


## *Anomaly Detection*



# Clustering

- ADWICE uses clusters to represent normality
- Adaptation of an existing data mining algorithm (BIRCH)



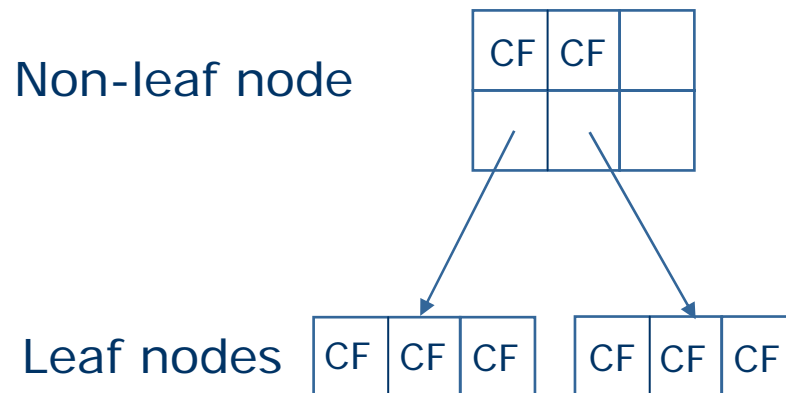
# What is a data point?

---

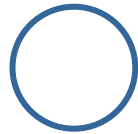
- General: A set of numeric values
  - E.g. measurements from sensors
- What about IP packets?
  - A vector of alphanumeric values in header of an IP packet
  - Transformed into vector of numeric values
  - In our tests: 41 dimensions
- Need efficient storage and search among summaries of collections of data points

# Basic ADWICE concepts

- CF (Cluster Feature)
  - Summary of cluster
  - [No, Sum, Sum of sq]
- Index: CF Tree
  - Maximal number of clusters (M)
  - Threshold requirement (TR)
  - Branching factor (B)



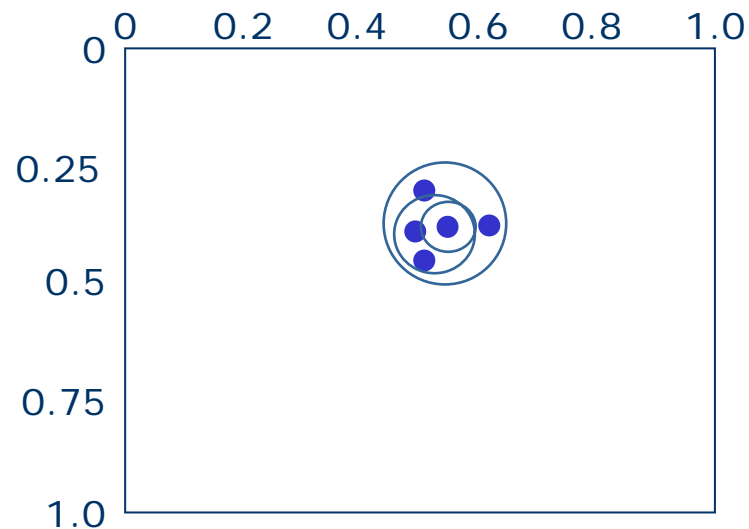
Threshold:



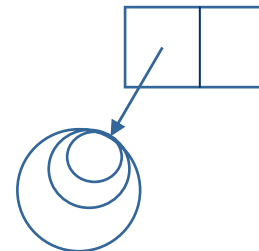
Max Number of Clusters: 3

Branching factor: 2

Data Space

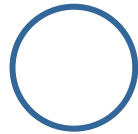


CF Tree



# ADWICE training

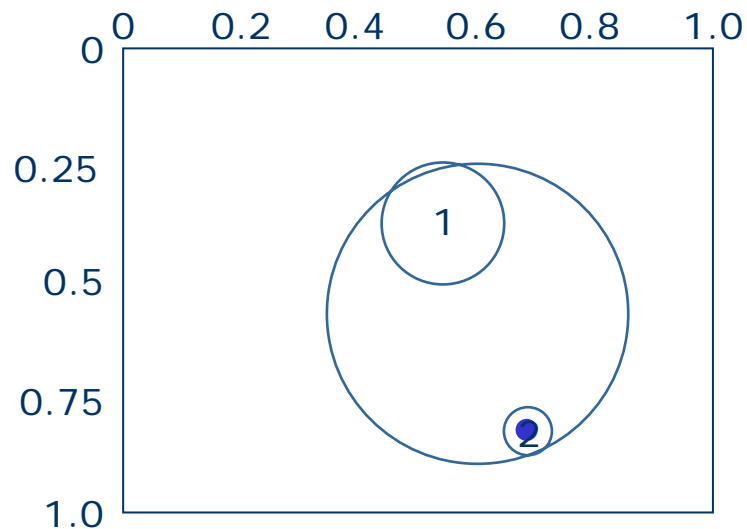
Threshold:



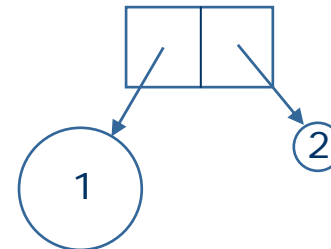
Max Number of Clusters: 3

Branching factor: 2

Data Space



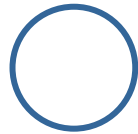
CF Tree





# ADWICE training

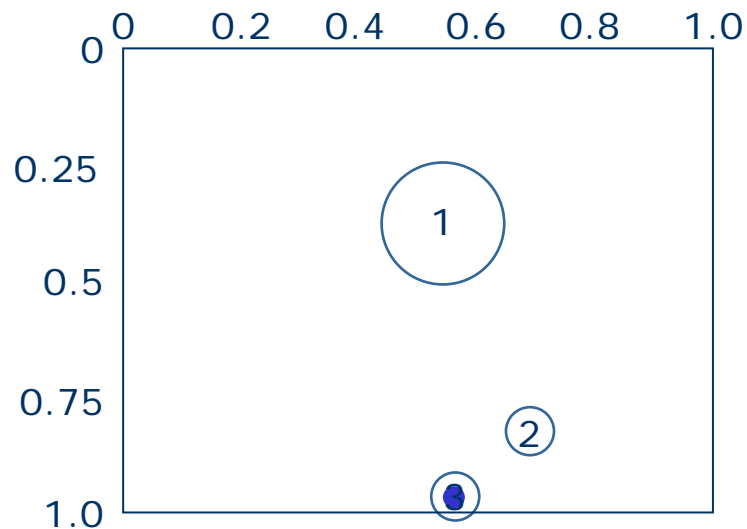
Threshold:



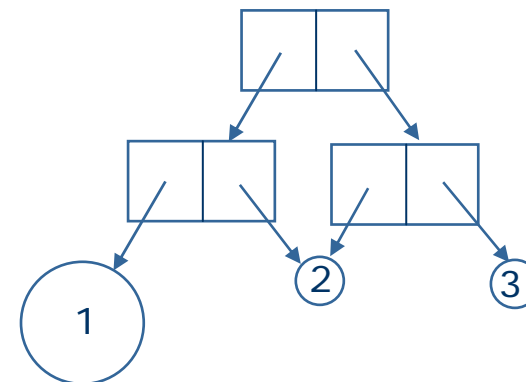
Max Number of Clusters: 3

Branching factor: 2

Data Space

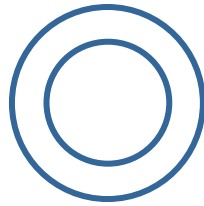


CF Tree



# ADWICE training

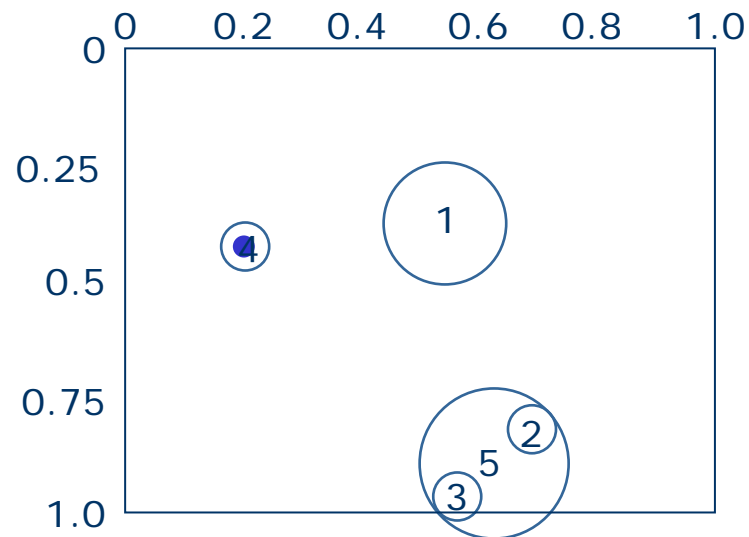
Threshold:



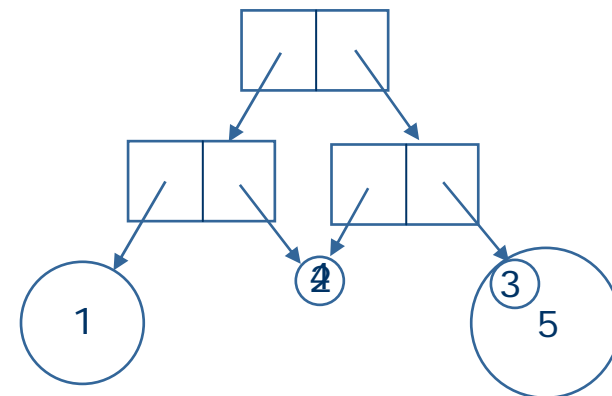
Max Number of Clusters: 3

Branching factor: 2

Data Space

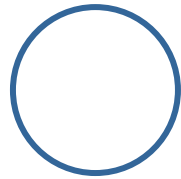


CF Tree



# ADWICE detection

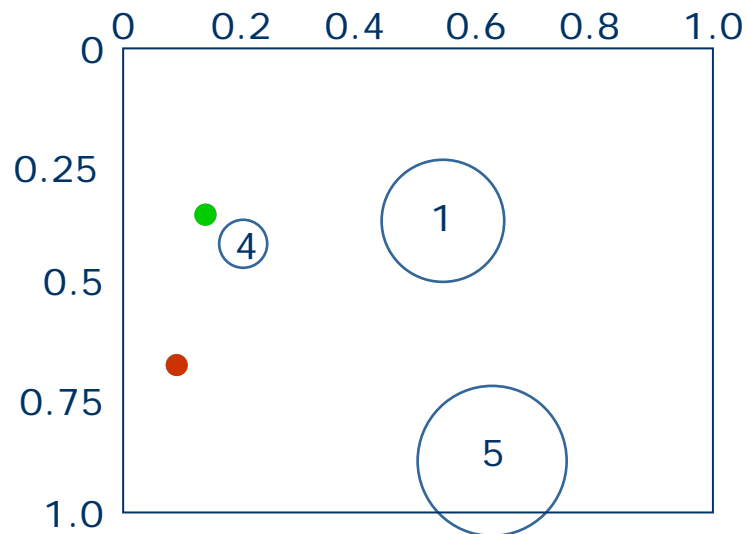
Threshold:



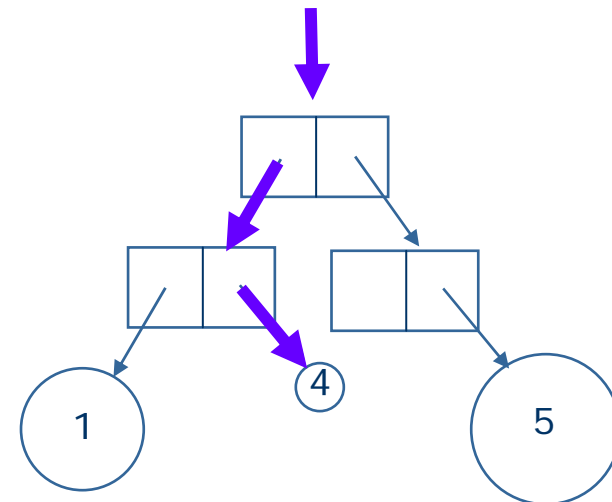
Max Number of Clusters: 3

Branching factor: 2

Data Space

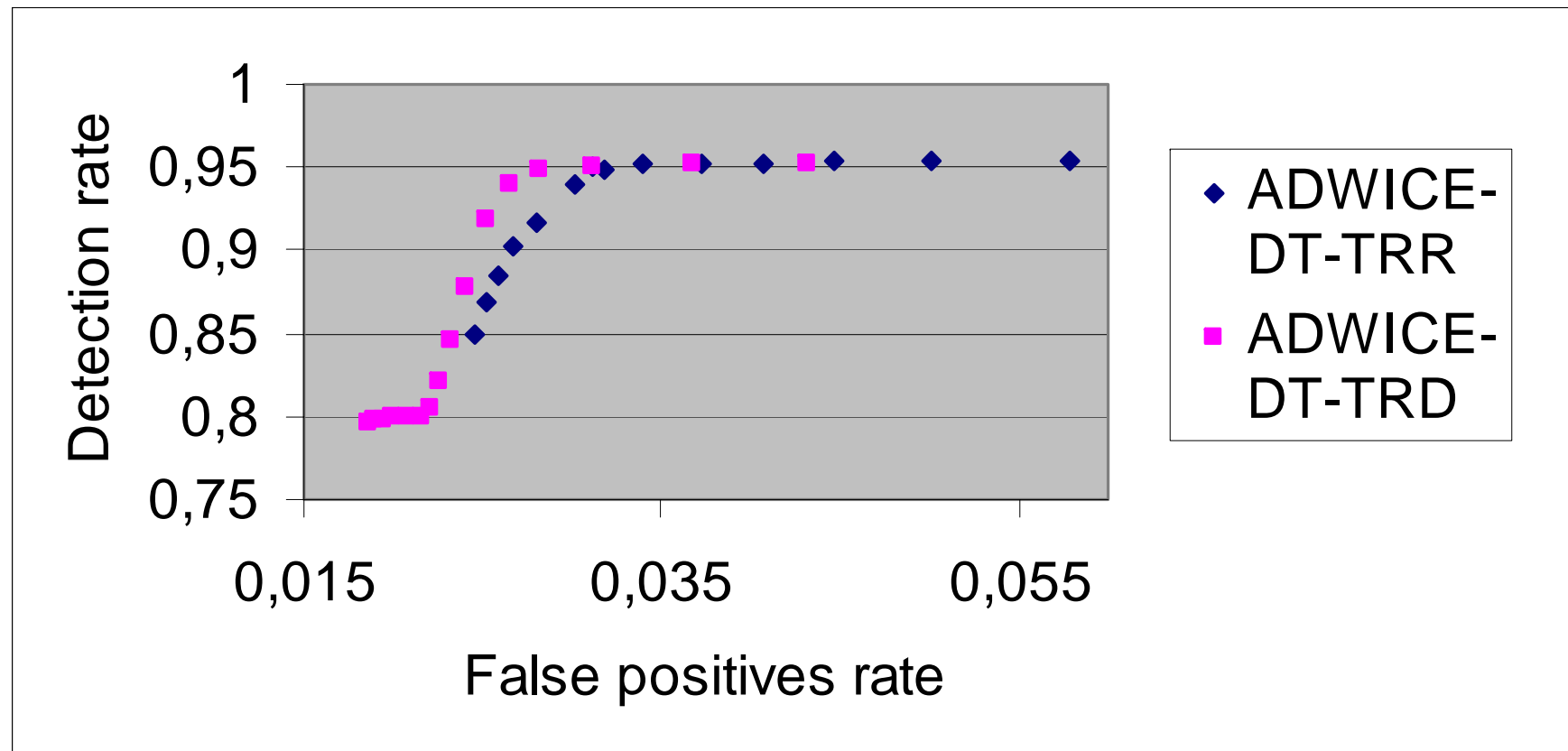


CF Tree



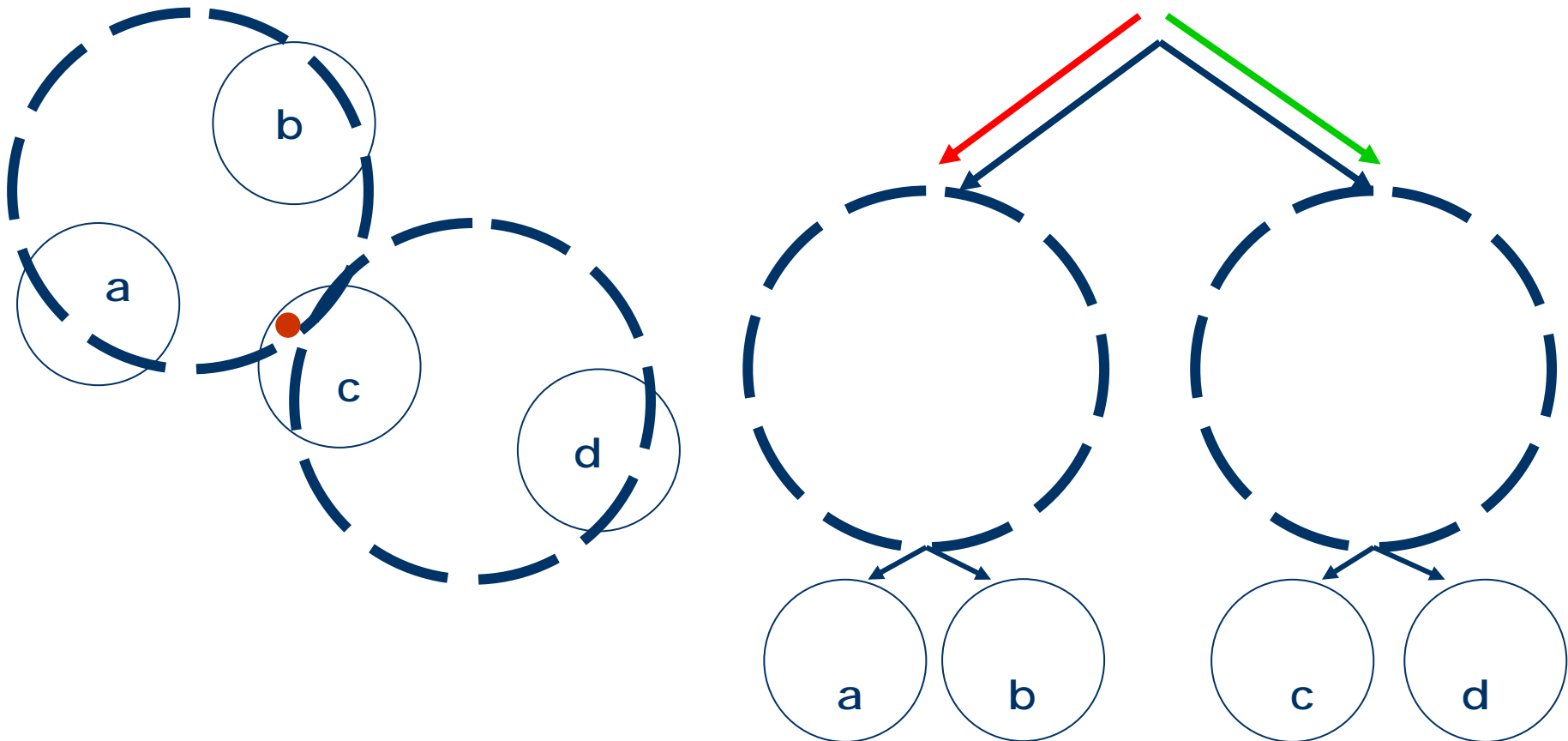
- KDD99 Data
- General properties
  - Session records (TCP/UDP summaries)
  - 41 features (flags, service, traffic stats ...)
- Training data
  - 4 898 431 session records
  - 972 781 normal, the rest (attacks) not used
- Testing data
  - 311029 session records
  - normal data and 37 different attack types

# Detection rate vs. false positives



# Index errors

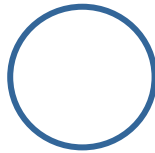
- Some false positives are due to index errors



- A new version of the algorithm: separates cluster formation and index updates
- How does ADWICE- Grid work?

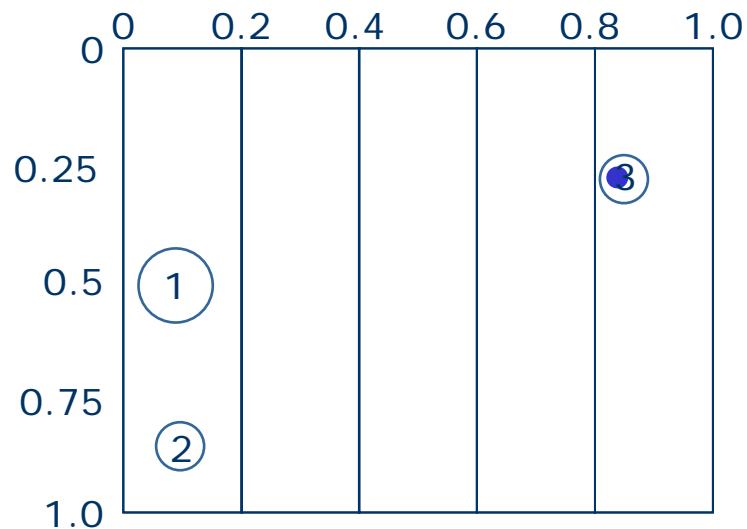
# ADWICE-Grid: Training

Threshold:

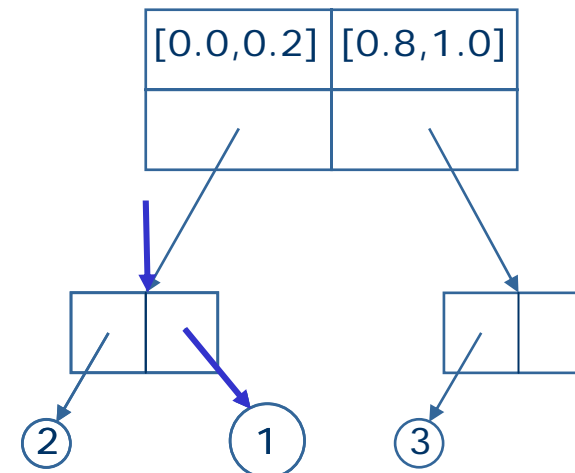


Max clusters in Leaf: 2

Data Space



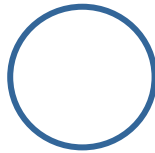
CF Tree





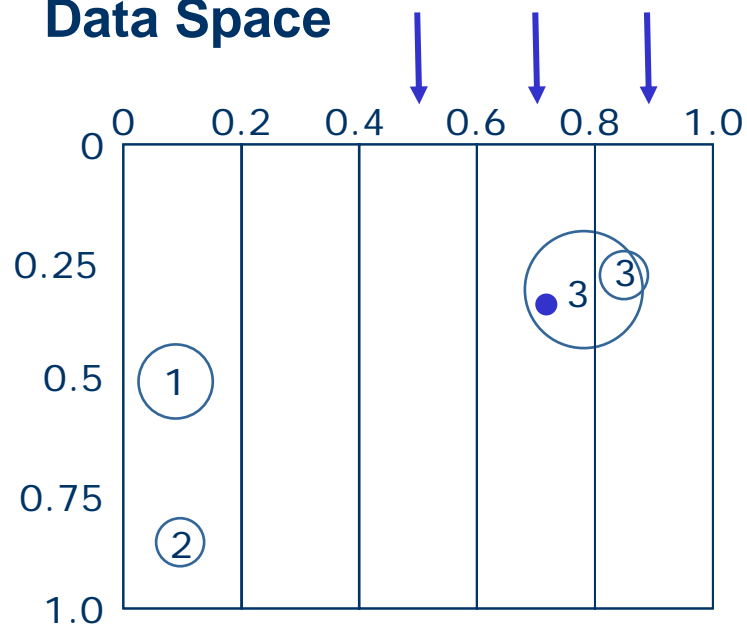
# ADWICE-Grid: Training

Threshold:  
(Search width)

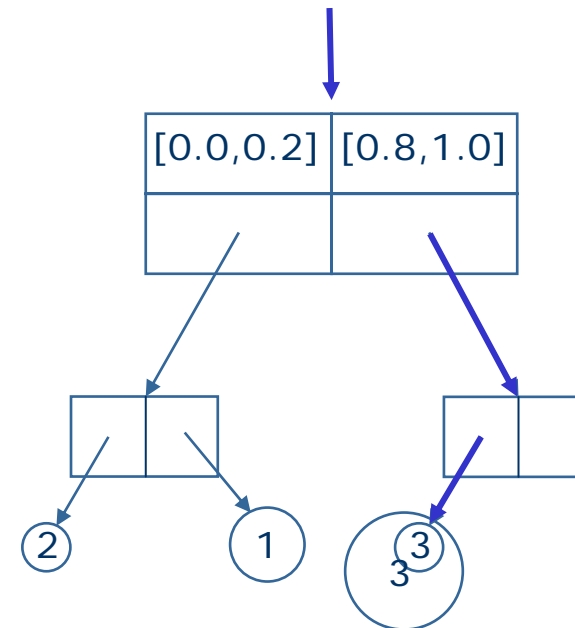


Max clusters in Leaf: 2

Data Space

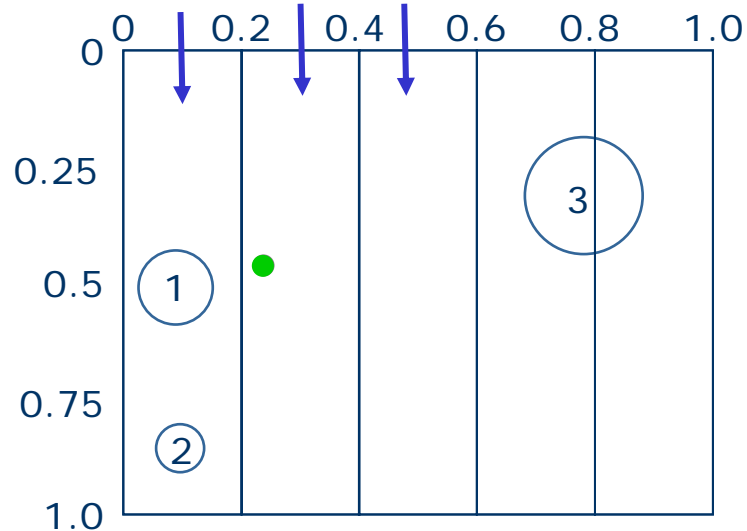


CF Tree

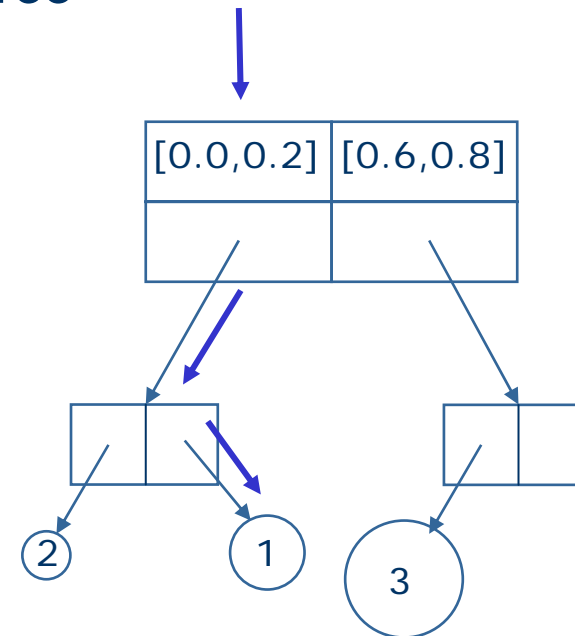


# ADWICE-Grid: Detection (1)

Data Space

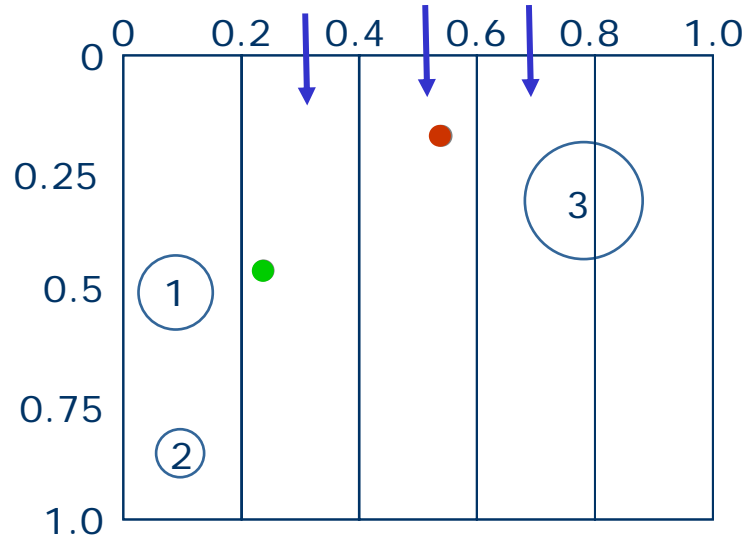


CF Tree

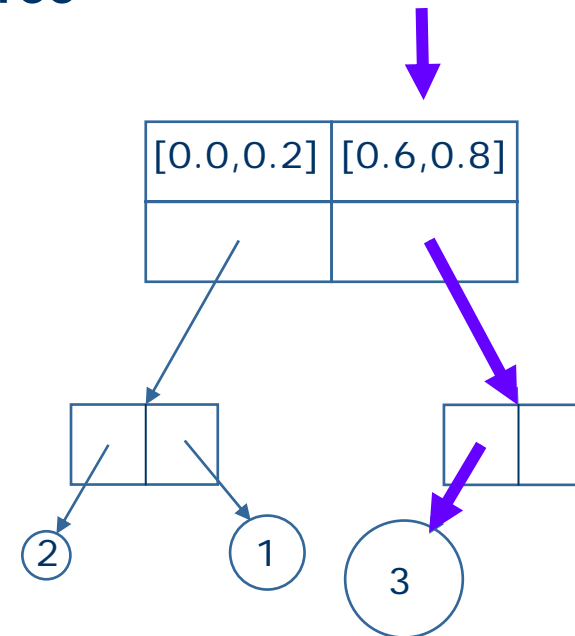


# ADWICE-Grid: Detection (2)

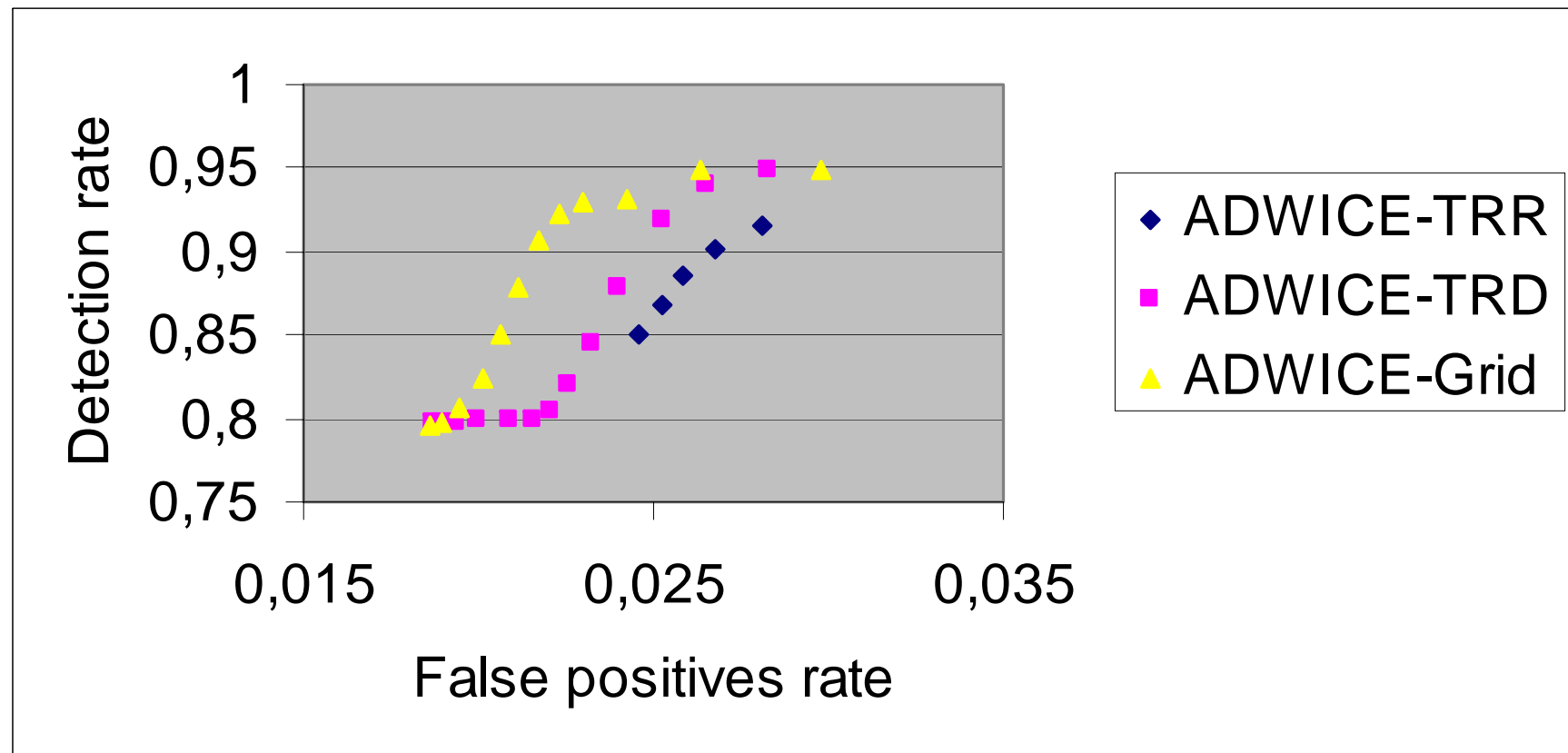
Data Space



CF Tree



# Detection rate vs. false positives



# Alarm aggregation

- Anomaly detection may produce many similar alarms (e.g. DoS, Probes, False positives)
- Similar alarms can be aggregated without losing accuracy



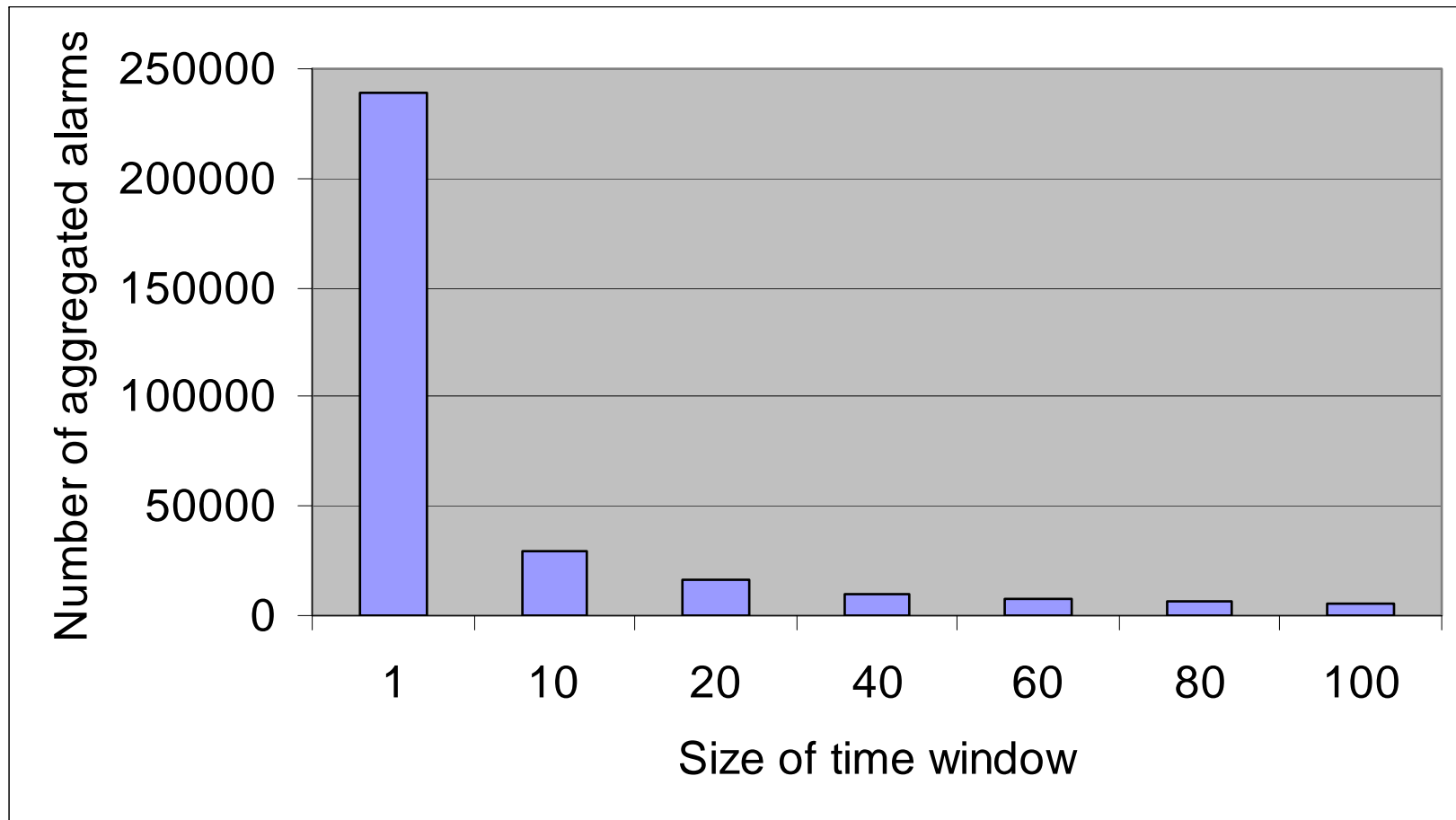
*Normal alarms:*

$\underbrace{\langle t1, \text{HTTP}, \dots \rangle \quad \langle t2, \text{HTTP}, \dots \rangle}$

*Aggregated alarm:*

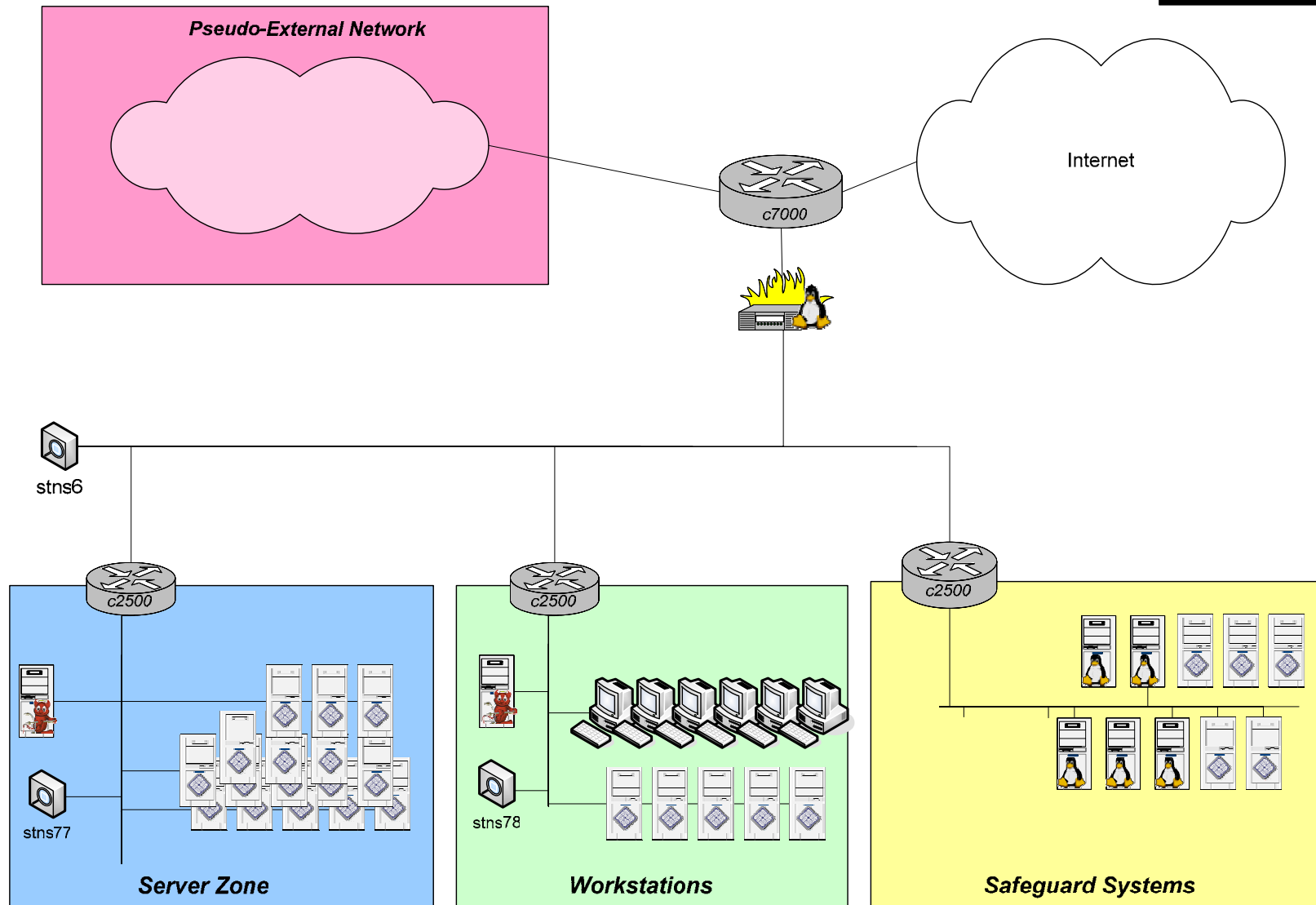
$\langle \text{Start}, \text{End}, \text{Count} = 2, \text{HTTP}, \dots \rangle$

# Alarm aggregation results



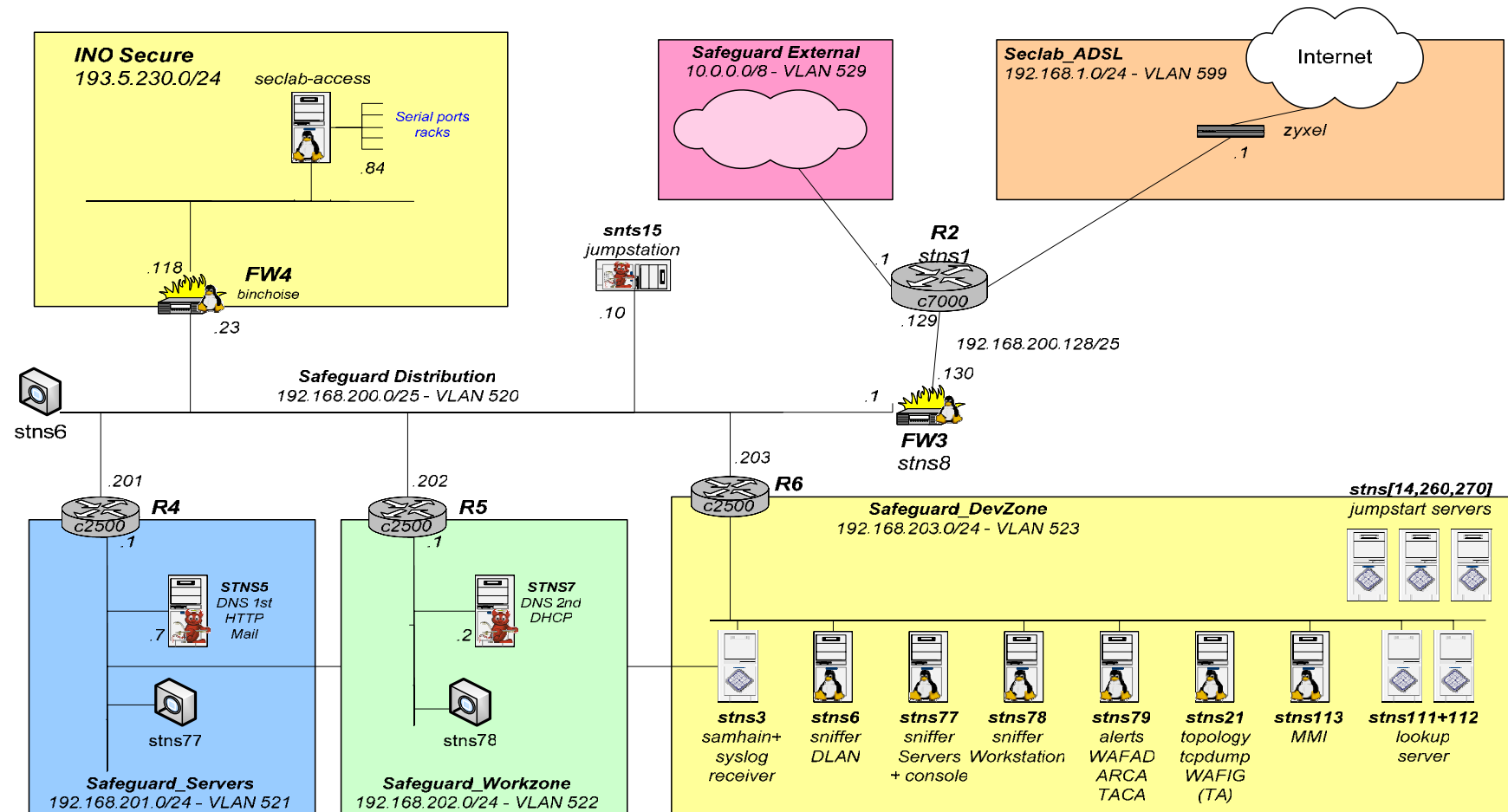
# Safeguard 100+ test network

Safeguard



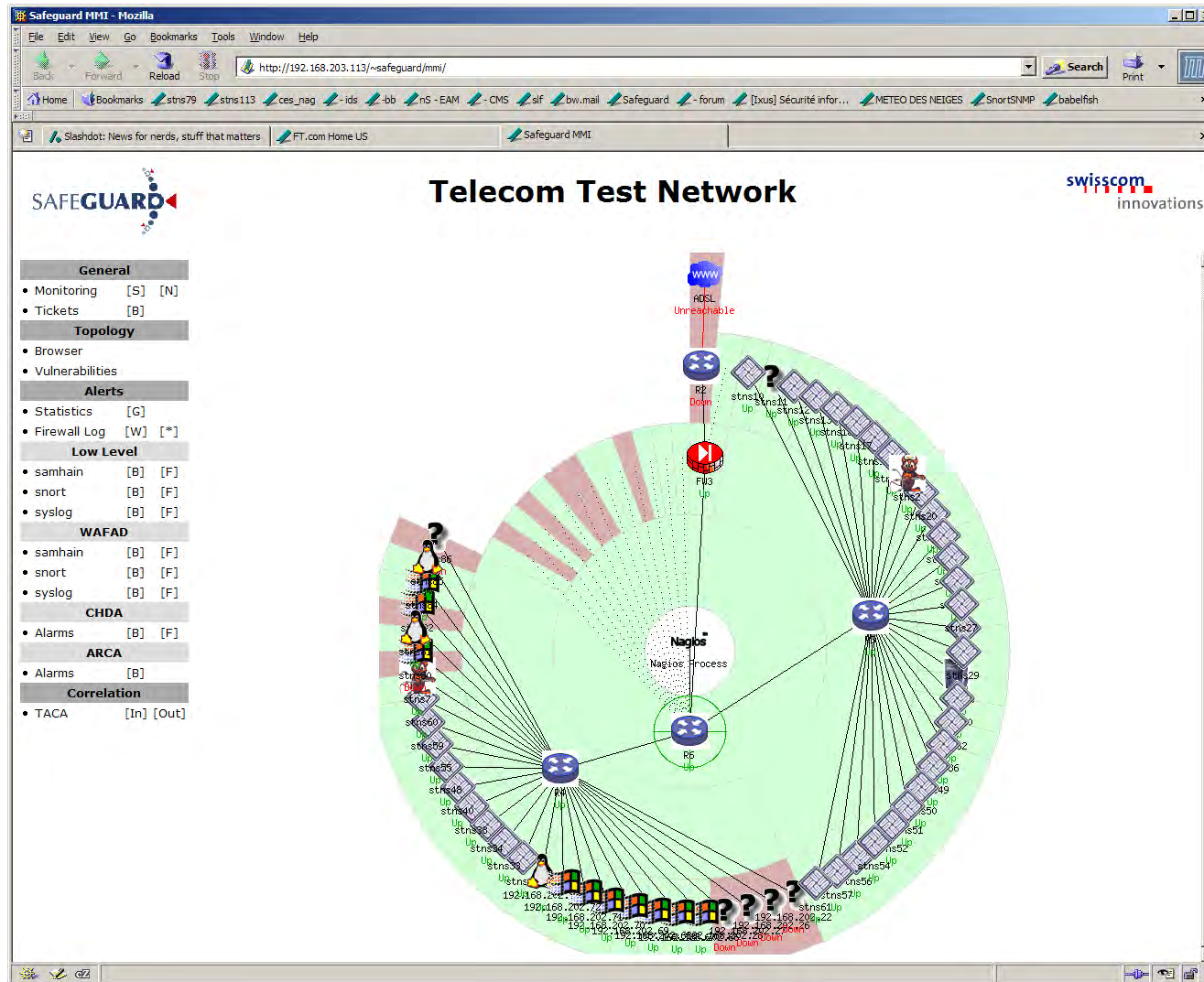
# Agents deployment

## Safeguard

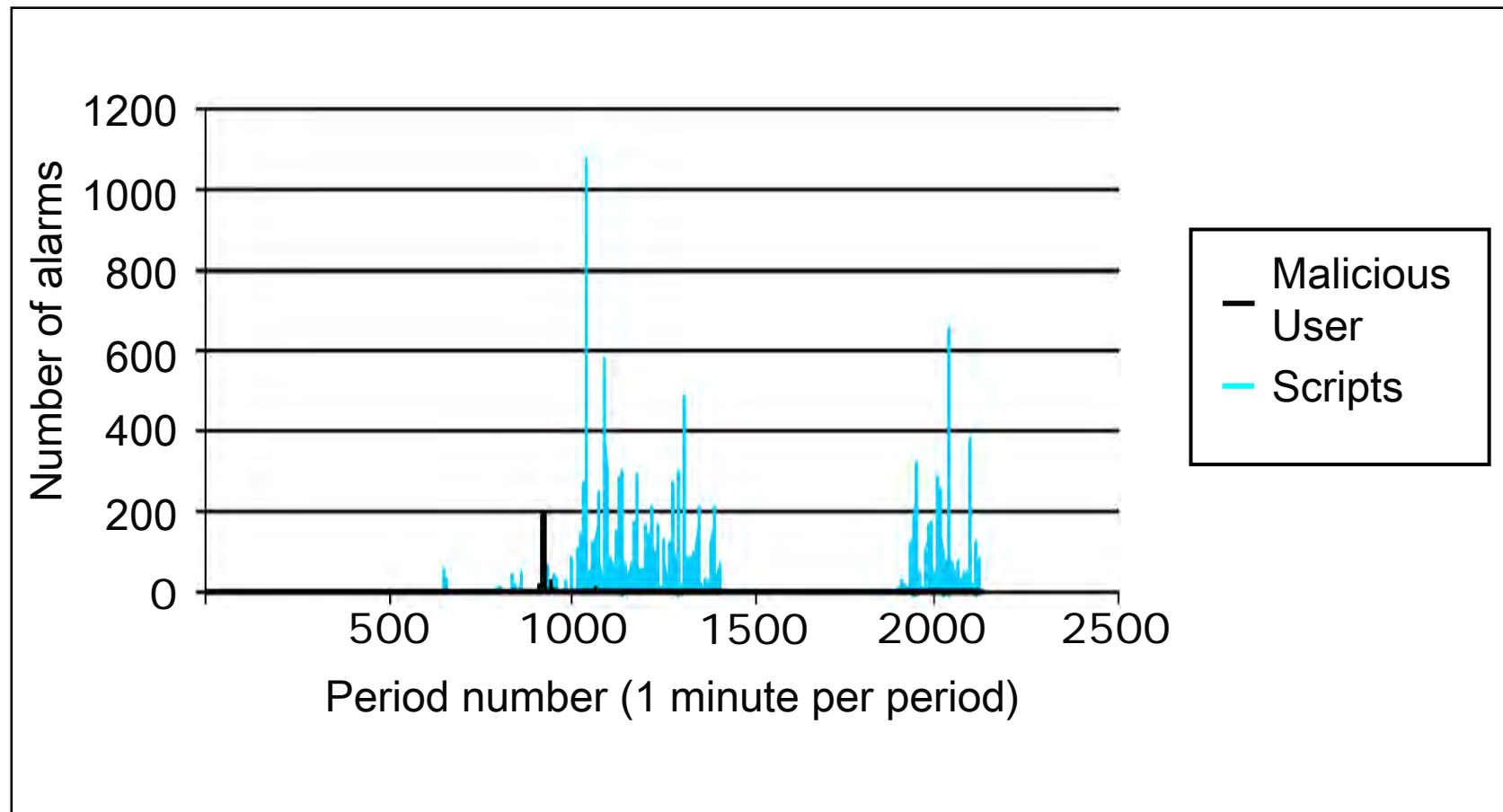




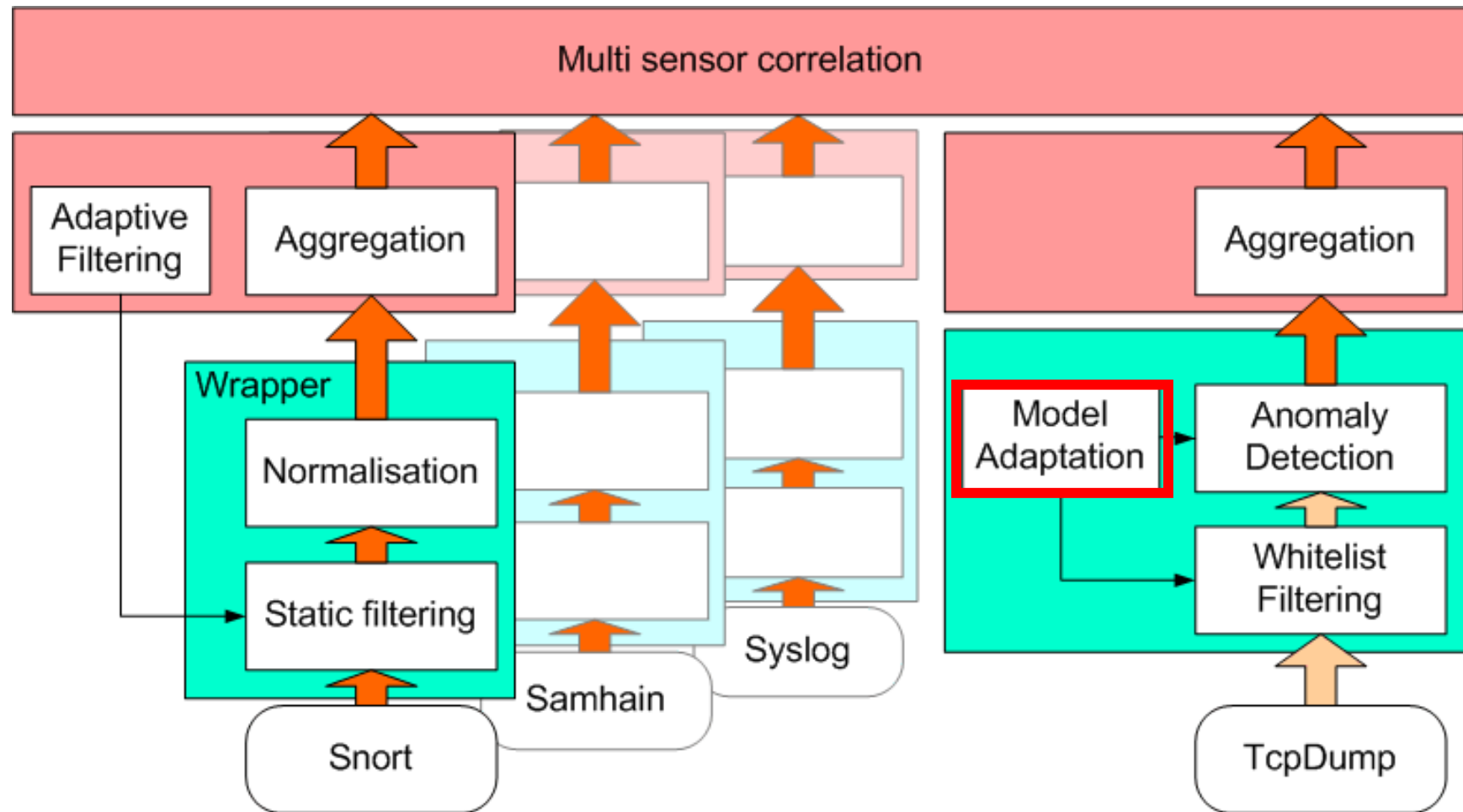
# One HMI agent interface



# A Safeguard scenario

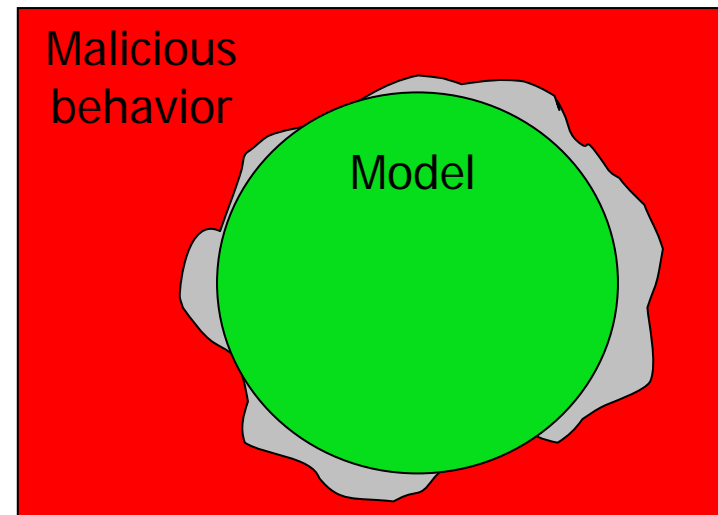


# Correlating alarms



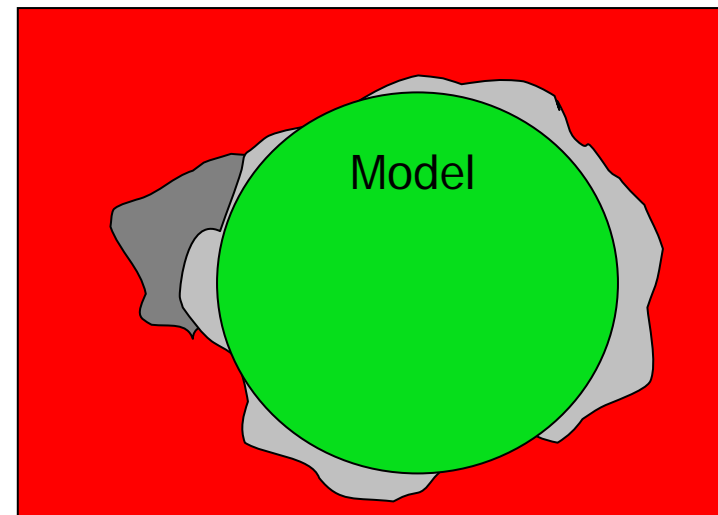
# Need for normality adaptation

- Normality is not static!



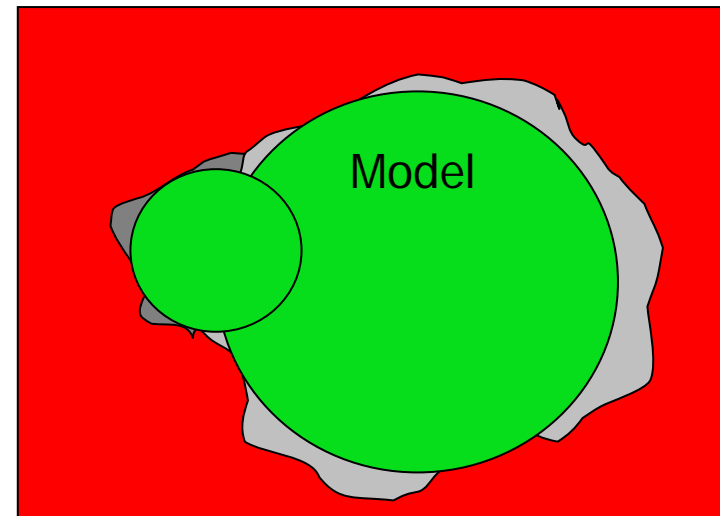
# New cases of normality

- Normality changes
  - New type of normal behaviour
- Old model incomplete
  - Evaluation using KDD data gives ~300 false positives for new normality

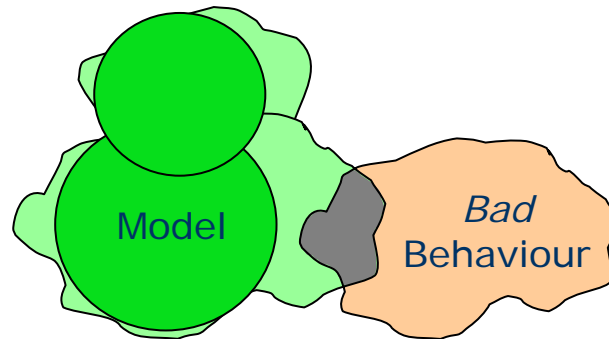


# Evaluation of normality adaptation

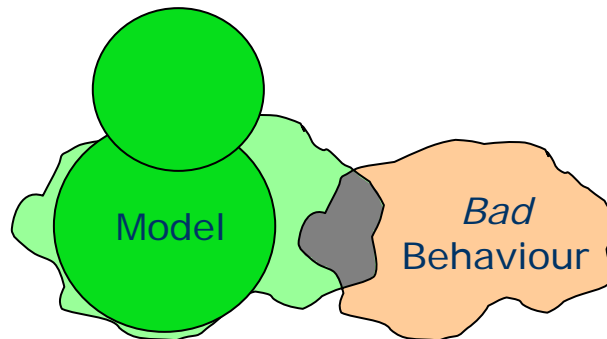
- Admin or system reacts
  - Recognize new false positives
  - Tells ADWICE to learn this behaviour
- Normality model adapted
  - From 300 to 3 false positives!



# Forgetting



- System keeps track of model usage
  - If time since last usage is very long for subset of clusters
  - Decrease size (influence) of those clusters and finally remove them if not used



Safeguarding critical infrastructures needs:

- Adaptive elements
- Incremental and scalable algorithms
- High performance for large volume of data
- Demonstration on realistic test beds
  - Research on open data sets :-)
- Understanding and mitigating interdependencies



- Application of ADWICE in anomaly detection for water management systems
  - Cooperation with Environment Protection Agency (EPA), USA
  - 50 scenarios: Time series data from simulated water system over an interval of one week
- Banking applications (Luxembourg 😊)