

# Big Data: Inside the Managing Process

**Petre DINI**

**Concordia University, Canada**

**China Space Agency Center, China**

[petre@iaria.org](mailto:petre@iaria.org)

# Plan

- **Big/small/linked/open data**
- **Data collection/filtering/storing**
- **Data retrieval/selection/interpretation**
- **Domains**
  
- **Case Study**
- **To come**

# Where to start from?

- **For existing systems; mostly, non-standard approaches**
  - - forensic
  - - facial recognition
  - - atmospheric models
  - - universe models
  - - health records
  - - human behavioral models
  - - etc.
  
- **For new systems; must be standard**
  - - syntax
  - - semantic
  - - taxonomy

# Data/Support types

- **Data representation**
  - A la SQL
  - Non-SQL
  - Pictures
  - Voice
- **Support digital (memory)**
- **Support digital (tapes)**
  - Pros/cons
  - (Google story!)

# Data correlation

- Hierarchical data
- Distributed/Isolated, Linked data
- Web / Deep Web [ :) internet profunda!  
Invisible Web, Dark Web, Hidden Web [not indexed]

- Data features

Primary [P]

Secondary [S]

$D ::= \{[P] [S] [context]\}$

# BIG

- Volume
- Complexity, Security, Risks to privacy
- Complex links, (fuzzy links, context-based links)
- Mixed Formal/Informal features
  - ⇔ defined fields (syntax) / text-like information

## Note:

- 90% of world's data was generated in the last two years
- > 90% is unstructured
- + Web ad Cloud offer new possibilities for discovery

## ⇒ New technologies:

For extracting/transforming/loading (ETL) and processing

For cleaning and organizing unstructured data in big-data applications; e.g., **Hadoop**

## Mixed media support

# Source

- **Sensors**
- **System [any] reports**
- **Neural/body systems**
- **Atmospheric measurements [short, medium, long terms]**
- **Universe observations [long term]**
- **Health measurements [small + big,...]**
- **Social measurements [migrations, resources, etc.]**

- **collection**

**Raw**

**Partially processed**

# Retrieval

- **Clustering**
- **Partitioning**
- **Summarizing**
- **Fusion**
- **Compressing**
  
- **? selecting the right features**



# Bumps

- Noise
- Probabilistic data
- Fuzzy-data sets
- Incompleteness
  
- Time-sensitive
- Time-free
  
- Timestamps
- Hierarchical timestamps
- Timestamps: [source][storage][processor][console]  
? Clock synchronization

# Storage

- **Distributed**
- **Access**
  - Internet Neutrality**
  - Accessibility**
- **Transparency Degree**
  - Open [e-government]**
  - OpenData Government**
  - OpenData Forum**
  - Private [financial, health]**

# Status

- Yesterday
- Today
- Tomorrow
  
- #1: big data exist
- #2: big data was dealt with
  
- ? → classical “hype” case

# Applications i [SMALL data]

- **Using Patient Data for Personalized Cancer Treatments**

- improve health outcomes
- support development of new therapies

- **Small Data**

**Seeking personalized data-derived insights form analysis of our digital traces**

**Personal devices Internet services for self-tracking**

**Fitbit**

**Patients like me**

**<http://quatifiedself.com>**

**Digital traces accumulated by social networks, search engines, mobile operators, online games, e-commerce**

# Applications ii [SMALL data]

- ? → regulatory challenges /FDA, HIPAA, privacy policies

Health Insurance Portability and Accountability Act

- Open mHealth <http://openmhealth.org>
- <http://smalldata.tech.cronell.edu>

# Application i [BIG data, bug traces]

- Large-scale bug traces
- Testing network devices before releasing them
- Binary/Linear Downsizing -> Downsizing Ratio
- Reproducing failures to facilitate the debugging process, real-world traffic needs to be captured and later replayed
- → high volume [peak-hour, at Beta Site, 20Gbytes, 30 minutes]
- Remove data redundancy in large a trace
- Note
- Linear Downsizing: rollback-and-reply; whenever a failure is triggered, the failure would be logged and the sequential traces triggering the failure are regarded as a whole and divided into equal-sized pieces of traces from the beginning based on a predefined size, rollback size.
- Binary Downsizing: BD locates the sequential traces triggering the failure by recursively splitting the traces on halves and replaying the smaller ones in turn, until the failure is missed....

# Applications ii [BIG Data]

- **Government Sector**
- **Ref: [Communications ACM, 03/2104, vol. 57, no. 03](#)**
- **BIG Data initiatives**
- **Japan: ITS [Intelligent Traffic System], Info-plosion, MEXT/NSF [Education..]/NSF**
- **UK: HSC [Horizon Scanning Center]**
- **Singapore: RAHS [Risk Assessment and Horizon Scanning]**
- **Korea: KOBIC, MFAFF, MOPAS**
- **EU: DOME [The Netherlands, Switzerland, UK, + 17 countries] + IBM /supercomputing center to handle a data set in excess of one exabyte per day derived from SKA radio telescope**
  - Exascale computing, transport and storage
  - Analyze all raw data collected daily (observable universe)
  - (One exaflops is a thousand petaflops or a [quintillion](#), 10<sup>18</sup>, floating point operations per second.) At a [supercomputing](#) conference in 2009, [Computerworld](#) projected exascale implementation by 2018
- **US: Genome Data on AWS [Amazon Web Services], CDC, NSF/NIH: BIGDATA, US Michigan [Statewide Data Warehouse]**

# Applications iii [BIG data]

- **Local Governments**
- **2011, Syracuse (NY) + IBM → Smarter City**
  - Bid data to help predict and prevent vacant residential properties**
- **Michigan's Department of Information → data warehouse**
  - To provide a single source of information about citizens of Michigan to multiple government agencies and organizations to help provide better services**
- **Facts**
  - ? E-Coli story**
  - ? Driver License story**



# Case study

## Positioning

### Issues

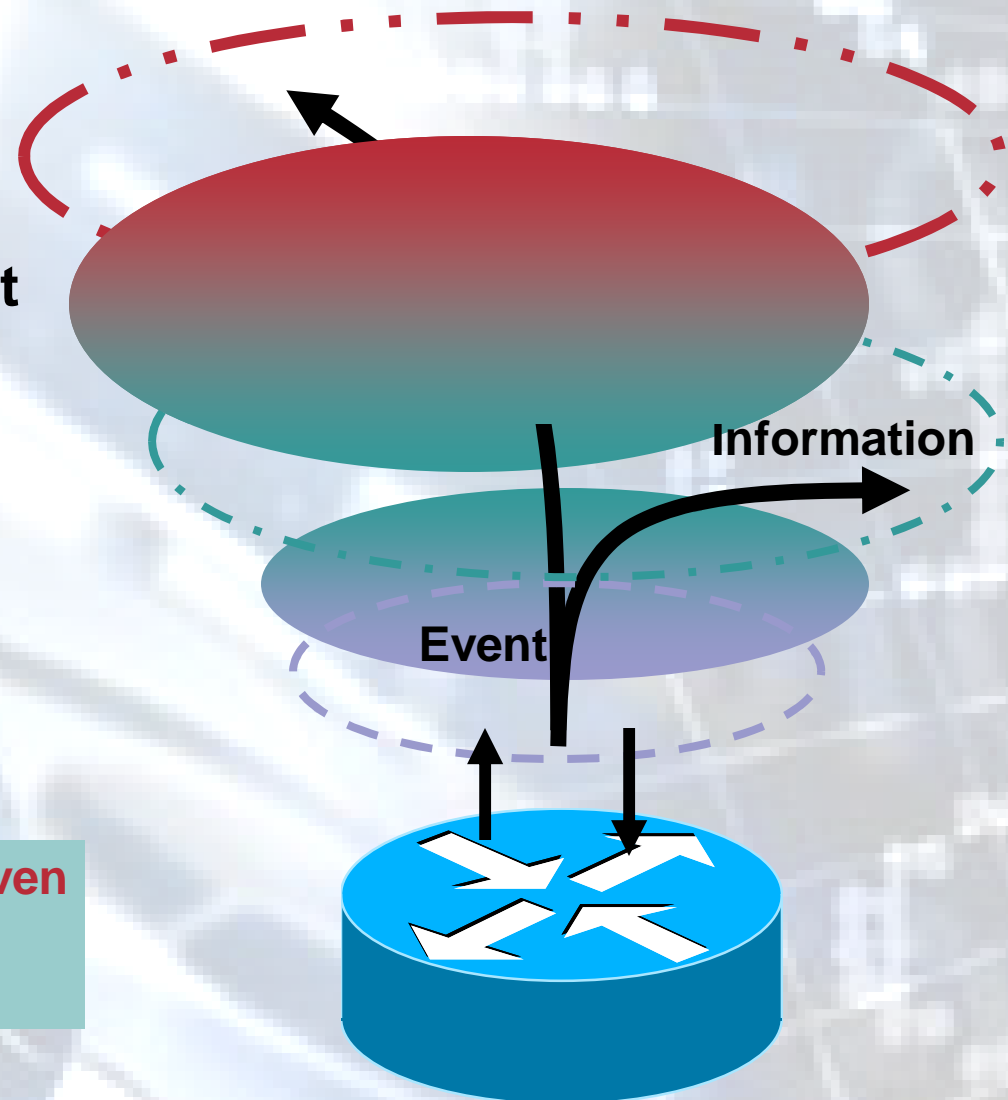
- **Event definition**
- **Event transport**
- **Event processing**
- **Business-driven events**

# Positioning

- **Layered event process architecture**
  - Issuing events
  - Processing events
    - ? Performance
- **Information bus**
  - Publishing events
  - Subscribing to events
    - ? Access/ transport
- **Towards autonomic event processing**
  - Network smartness vs. network management

# Get the infrastructure behavior

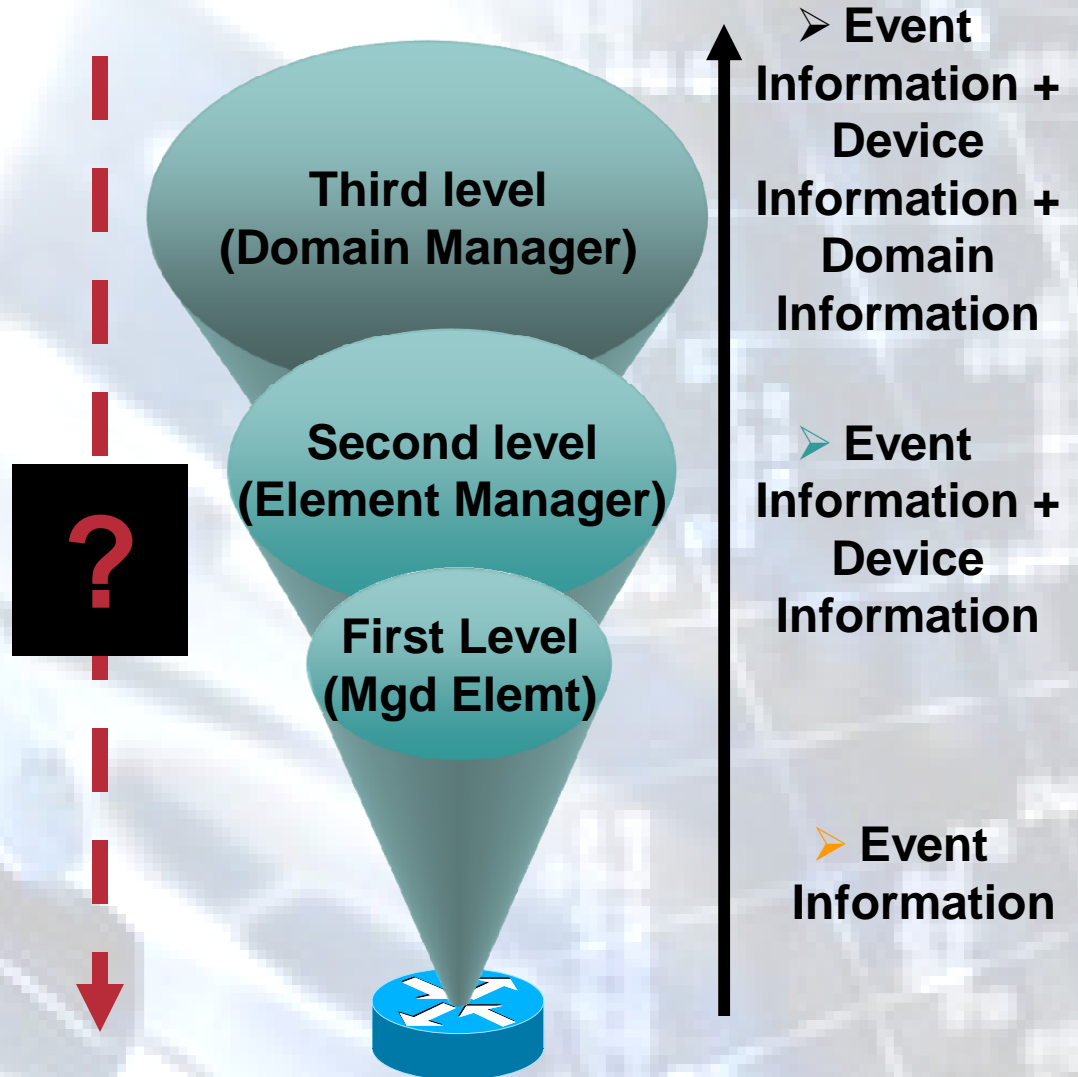
- Act (pre-emptive, proactive, reactive,...)
- Correlate (diagnostic, troubleshooting, impact, root cause, ...)
- Get status (push/poll)



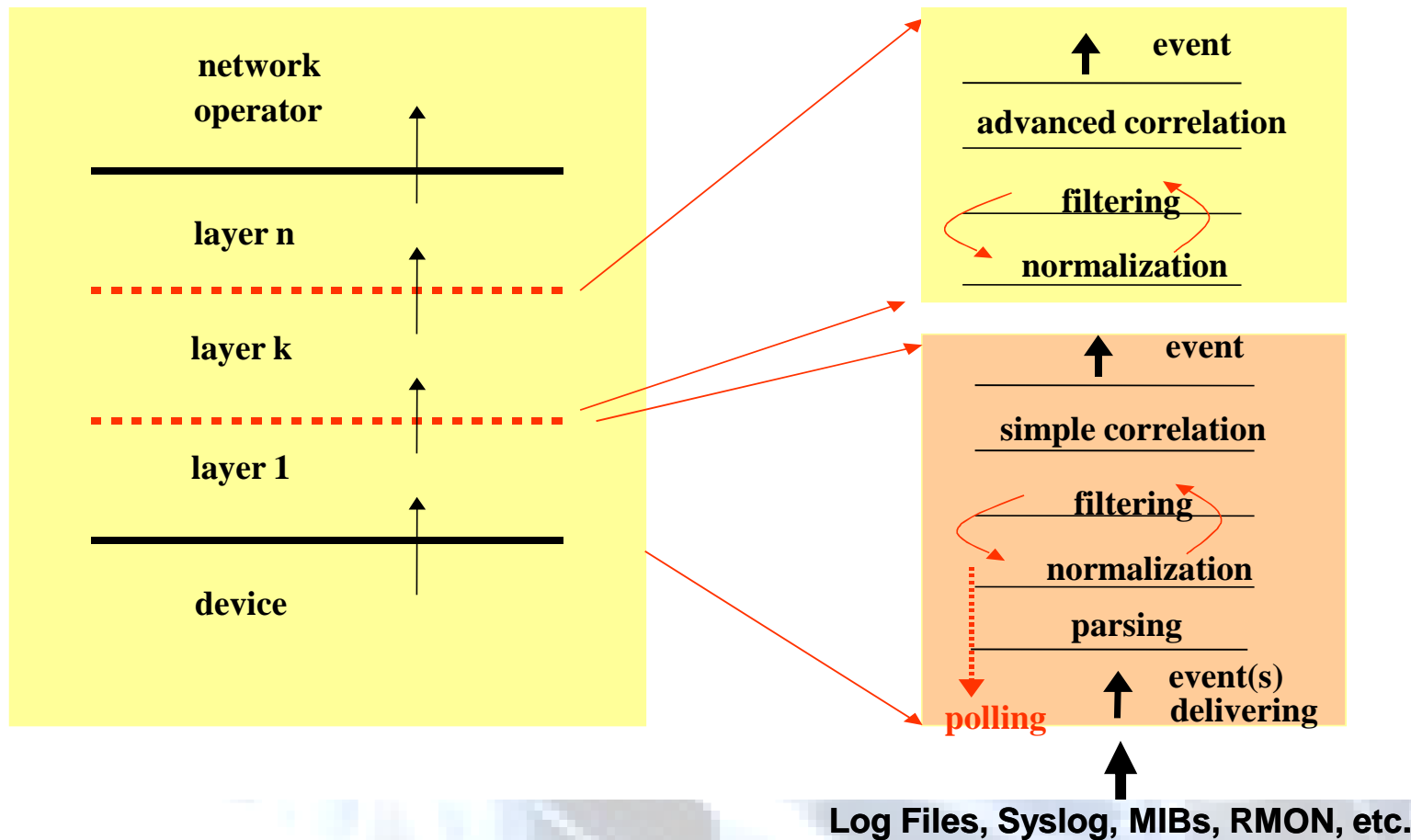
All operations can be policy-driven  
- top-down  
- bottom-up

# Bottom-up vs. Top-down

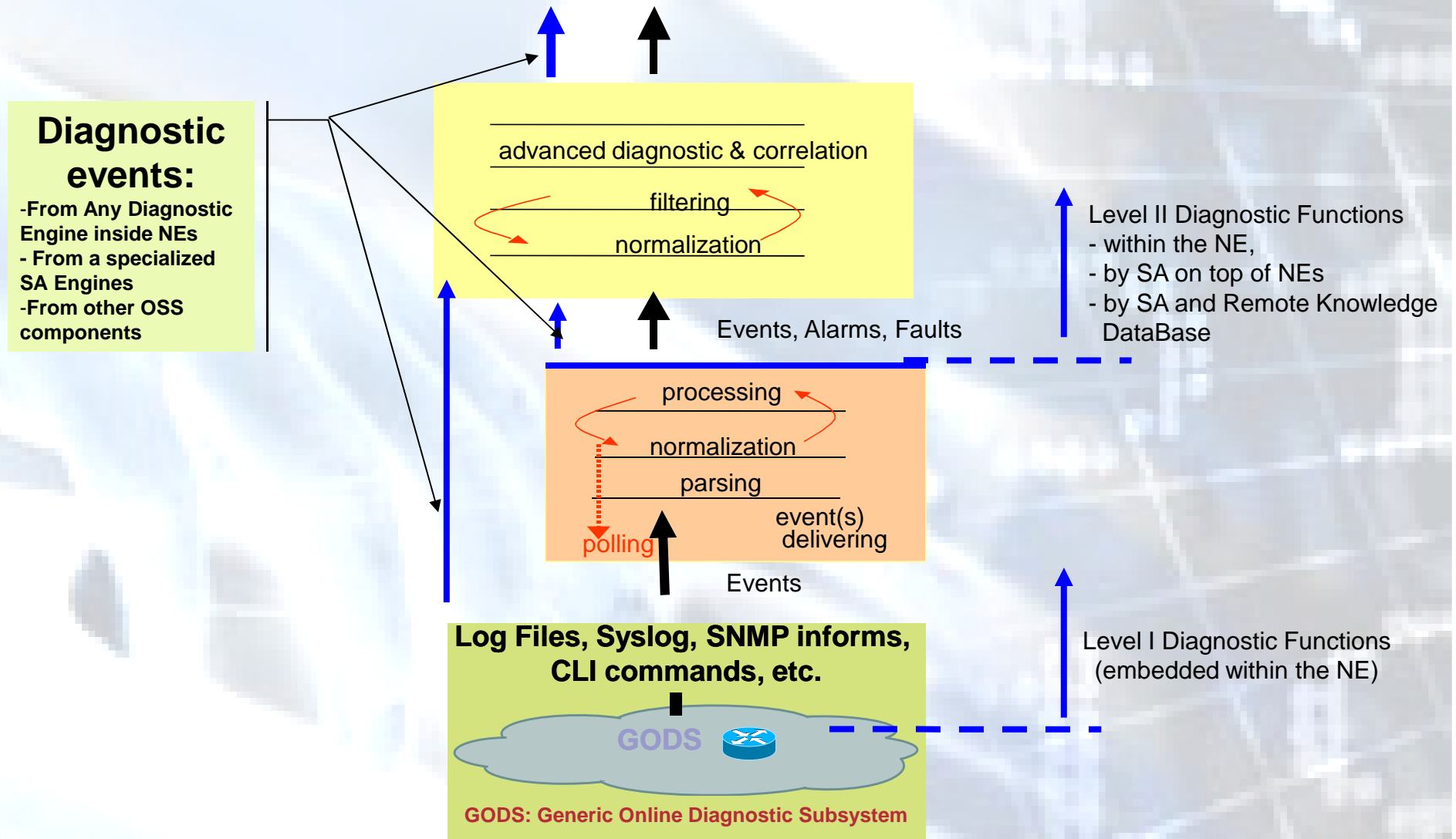
- Domain Manager enriches with domain information
- EMS enriches with multi-device information
- Notification Engine collects OS notifications



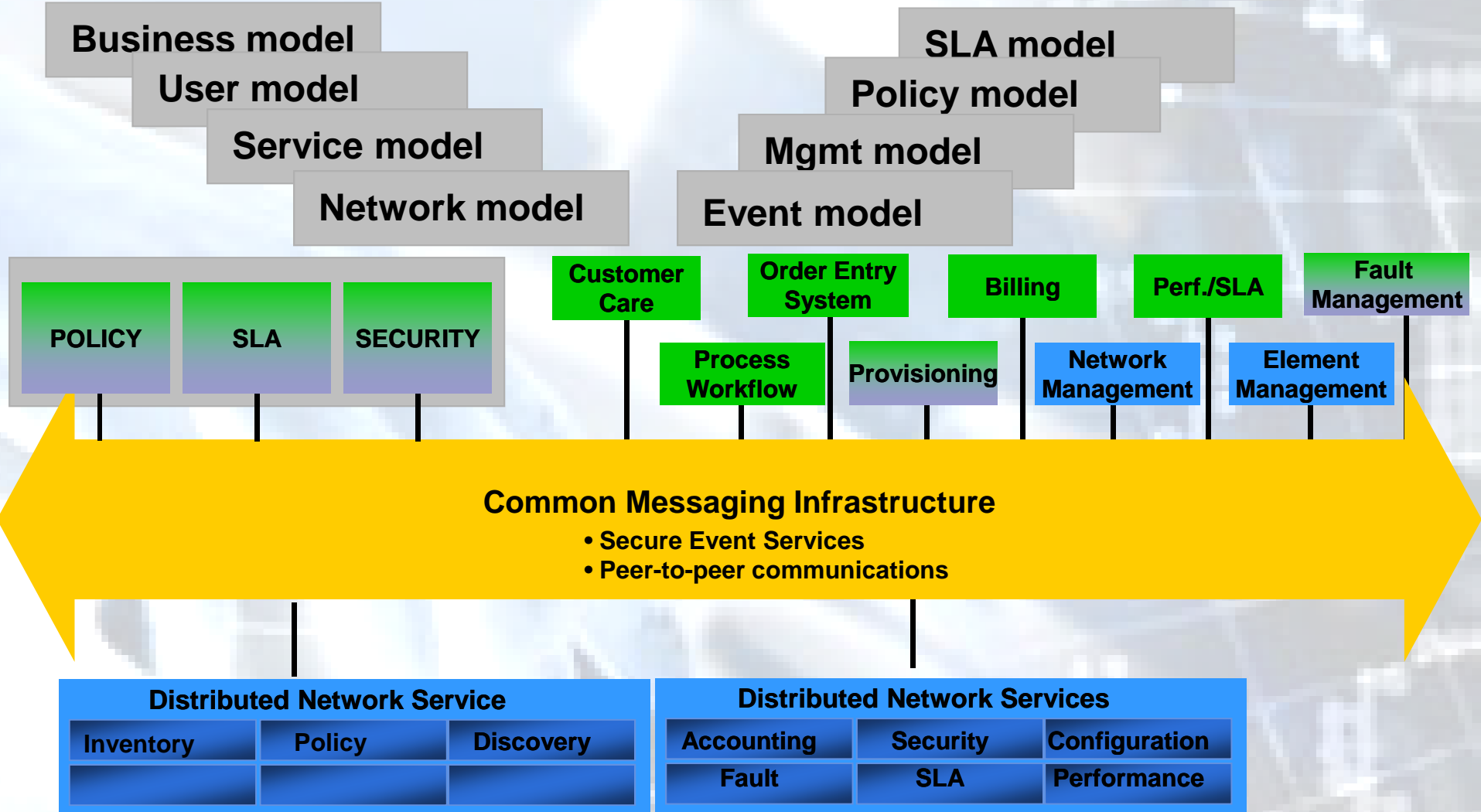
# A Layered Processing View



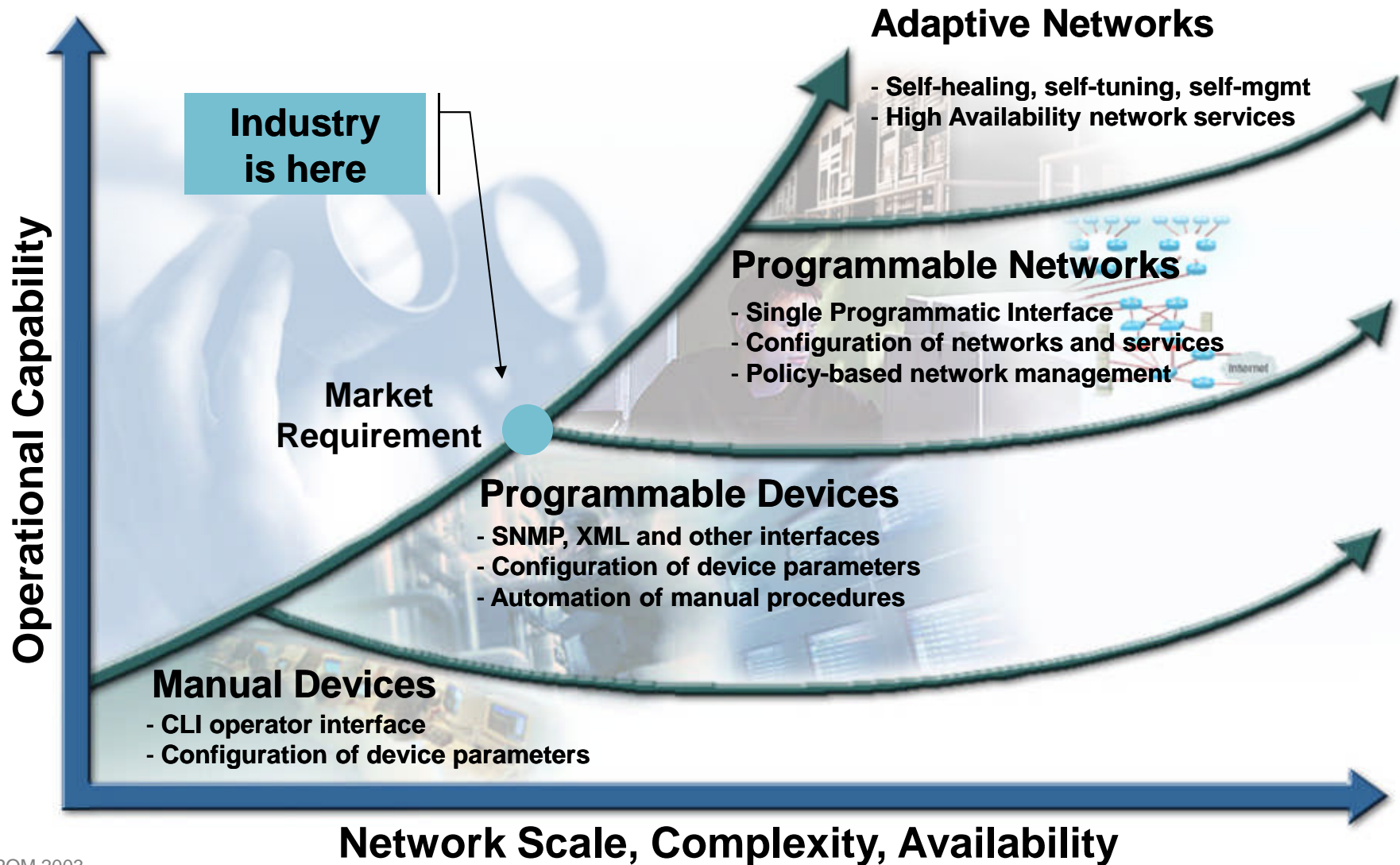
# Multi-level diagnostic



# Communication Bus

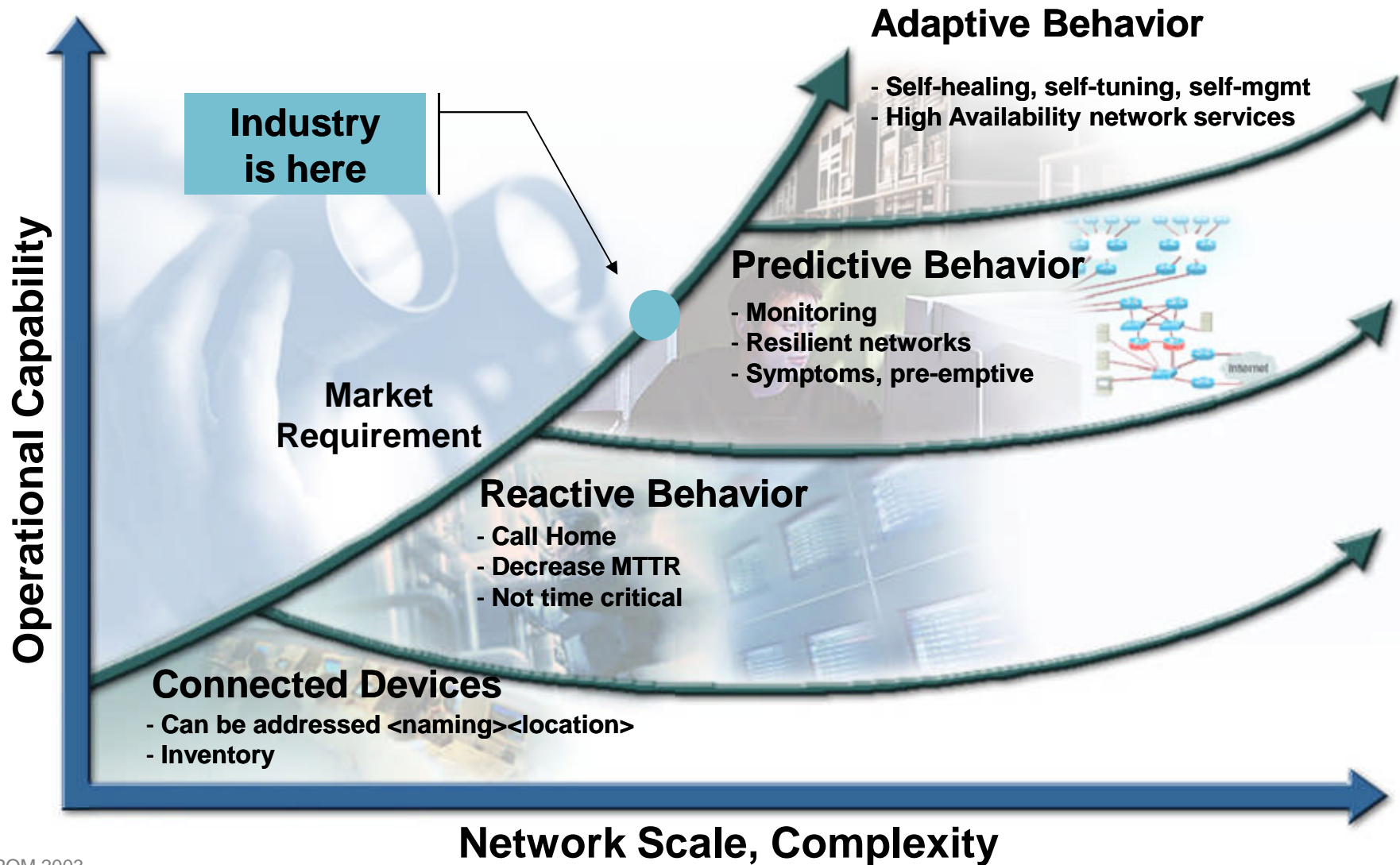


# BIG Data and Evolution of Network Manageability

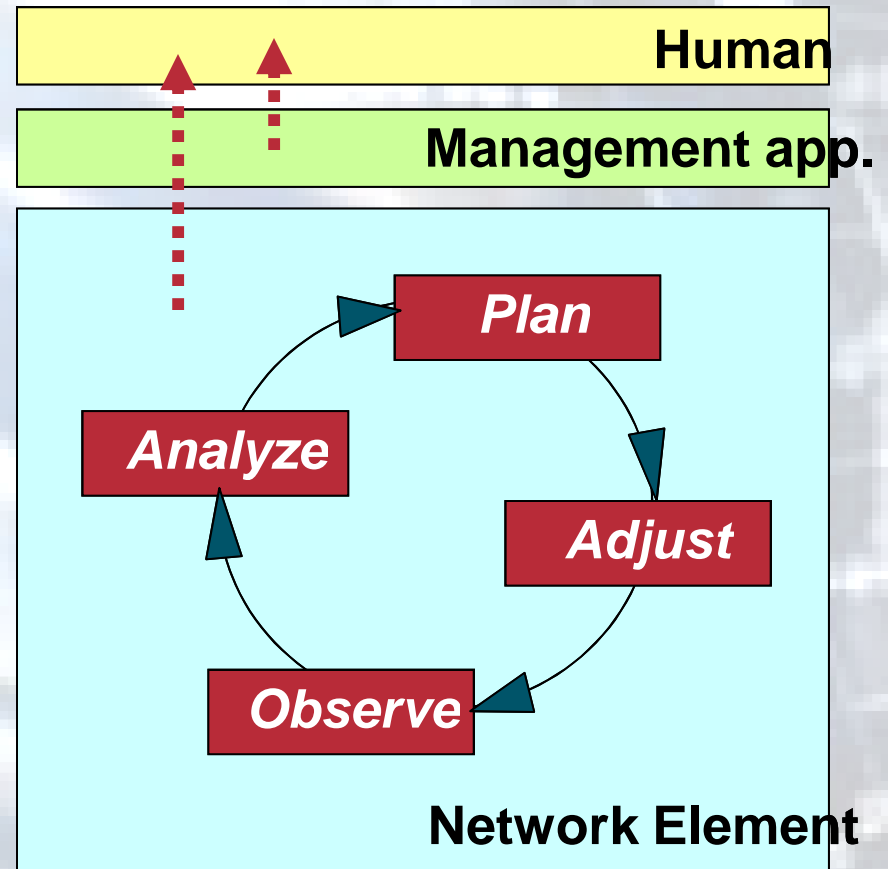
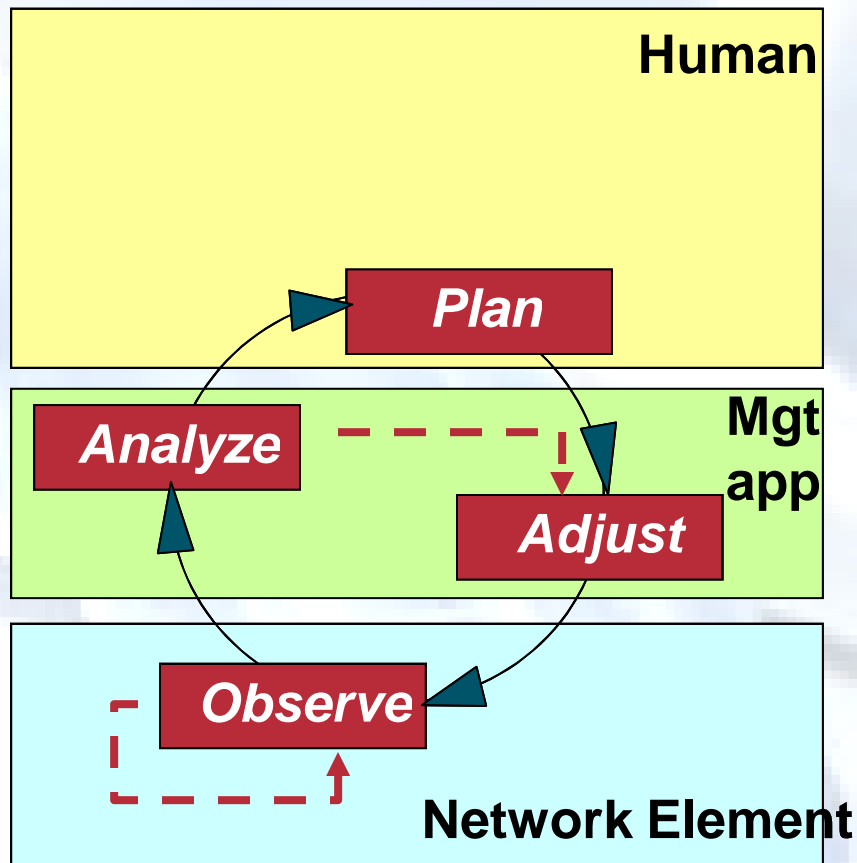




# BIG Data and Evolution of Network Smartness



# BID Data and Autonomic Computing



(a) Typical management control loop (b) Closed management control loop in autonomous network

# Challenging Issues

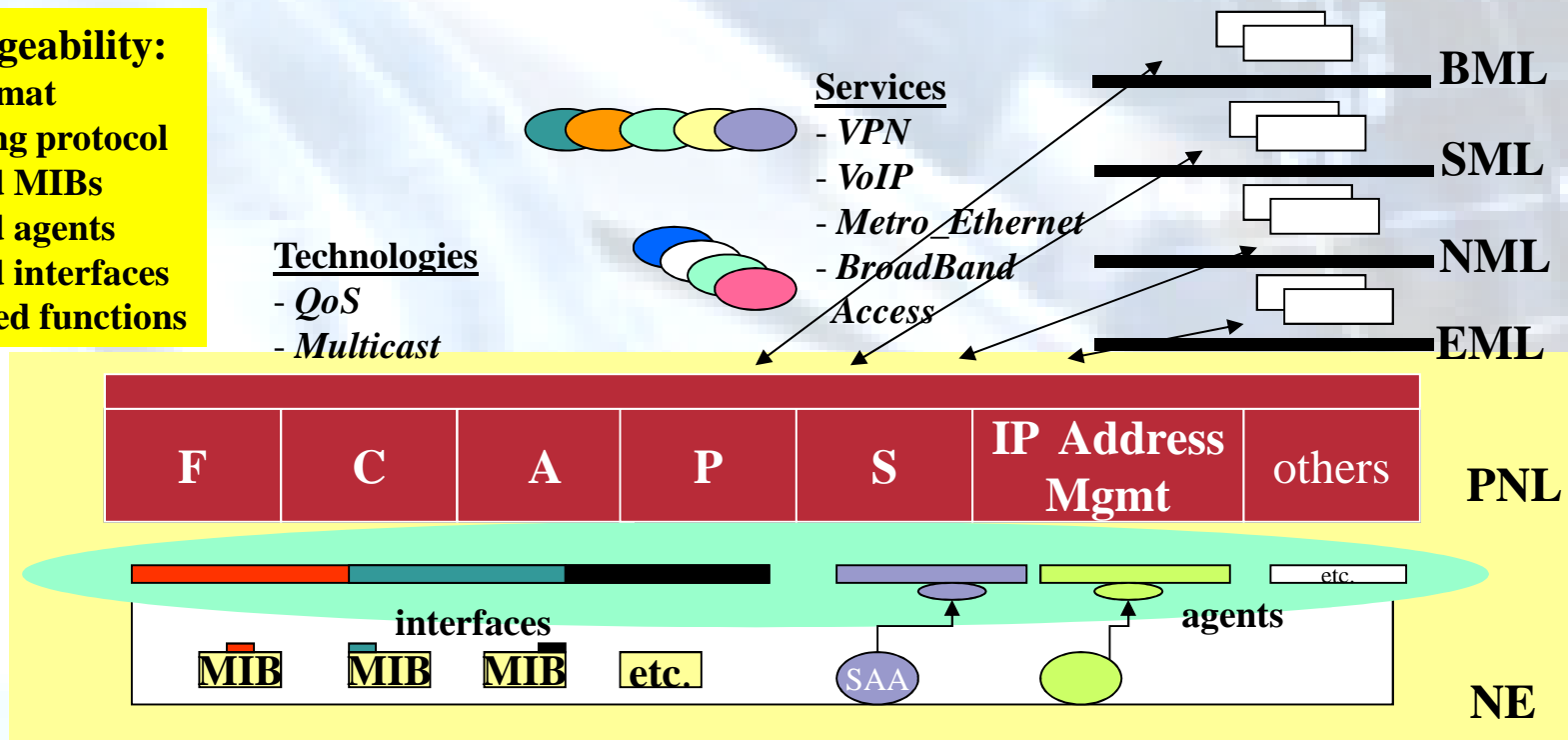
**Too Many**



# Syntax Issues

- Various formats
- Myriad of conversions needed
- Lack of syntax control

**NE Manageability:**  
 ? data format  
 ? conveying protocol  
 ? required MIBs  
 ? required agents  
 ? required interfaces  
 ? embedded functions



# Syslog Message "Body" Format in the IOS

CONSOLE

\* Sep 20 01:12:31: %SYS-5-CONFIG\_I: Configured from console by vty1 (144.254.9.79)

Timestamp	IOS Component	Severity	Mnemonic	Message-text
-----------	---------------	----------	----------	--------------

Timestamp from the server

SERVER

Sep 20 01:07:00 router.cisco.com 571: Sep 20 01:12:31: %SYS-5-CONFIG\_I: Configured from console by vty1 (144.254.9.79)

Router

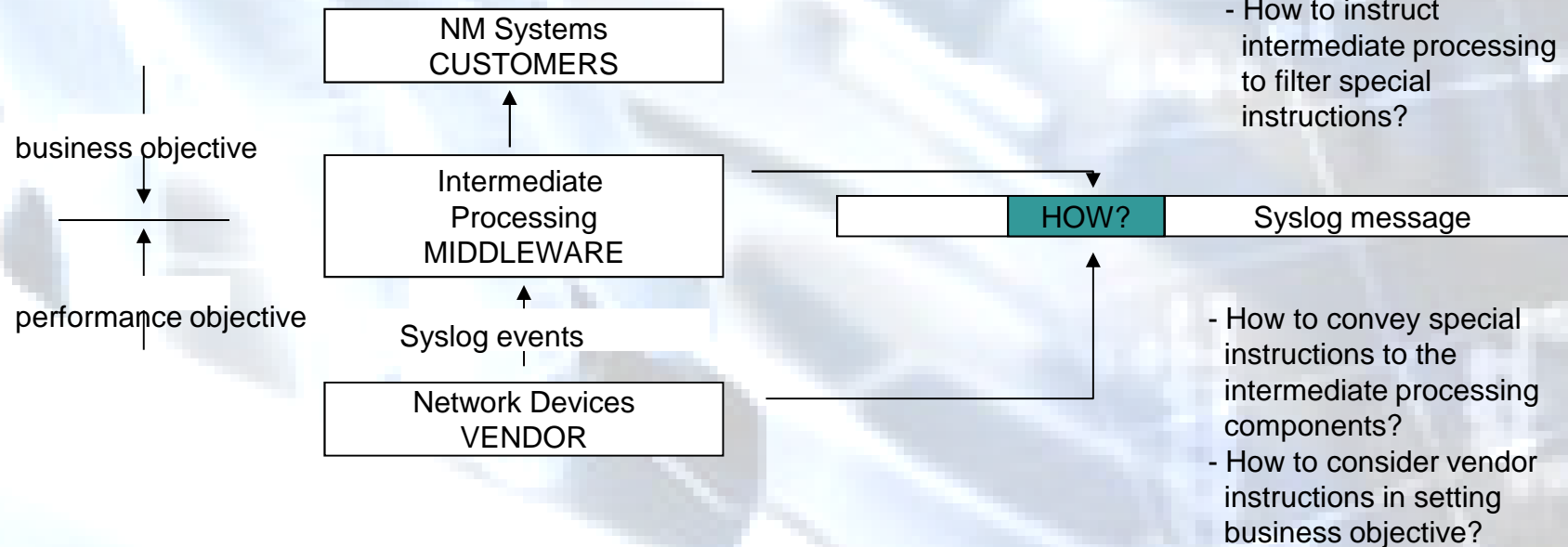
Timestamp from the router

- NTP is needed!
- Header:level can be different than Body:severity

# Semantic Issues



- Naming
- Context-defined
- Smart events



# XML Tagging is Not Enough

% versus <XML>

: <<a><b><c>>  
: (((a)(b)c))

1. <a> <b> <c>  
? ? ?

2. <<a> --- r1 -- <b>> -- r2 -- <c>  
? ?

e.g.,  
<a> -- Interface (? OID)  
<b> -- Port (? OID)  
<c> -- Severity

?  
Tag table (??)  
Tag List:  
<name><semantics>

?  
Tag relationships

?  
Naming service  
required

- Despite the problems caused by its use:
  - – The messages don't have a standardized definition
  - – Priority is geared toward UNIX problems
  - – Priority is not used consistently
  - – Not reliable
  - – Not secure
- some key features, (i) ease of use for developers, (ii) familiarity, and (iii) ubiquity makes it a workable solution.

# Timestamps issues

- **Format**
- **Clock-free event sources**
- **Sources-destination timestamps**
- **Delay tolerant networks**
- **Localizing processing**
  - Local synchronization
  - Wide synchronization
- **Reliable timestamps**



# Example: Syslog

**[ field1 ] % [ field2 ] % [ severity ] % [ priority ] % [ mnemonic ] % [ free form field ]**

## Well identified fields

**[timestamps]  
[facility ]  
[severity ]  
[priority]  
[mnemonic]**

## Free form field (the richest in semantic)

**[..English plain text..]**

## Field separator

**%**

## Issues

- Number of fields varies
- Value space of the fields is not uniform/standardized
- Semantic of timestamps is not uniform/or not defined
- Mnemonic is not modeled
- The English text is only humanly readable/useful
- Automation is difficult due to the “natural language processing” needs

# Things started to get fixed

- **Syslog, SNMP/MIB: IETF**
- **Adaptive message format: IBM/Cisco**
- **Intrusion detection format: IETF**
- **Anomaly report format: OASIS**
- **Incident handling format: IETF**
  
- **NGN management : ITU-T [Focus group]**

# Still to answer...

- **Concepts such utility-based computing, autonomic computing, diagnosis-in-the-box, diagnosis out-of-box, adaptable applications, self-adaptable applications, and reflexive environments require a new approach of formalizing events, architecting event-based systems, and integrating such systems.**
- **Additionally, GRID systems bring into the landscape the concept of intermittent and partial behavior related to resource sharing that may require a special semantic on SLA/QoS violation events.**
- **Events related to traffic patterns and the dynamics of performance and availability changes in such environments requires particular metrics and processing, as well [accounting, outage].**
- **Another hot area quite poorly covered in terms of event-related requirements is MPLS OAM and all aspects related to MPLS VPN.**

# ALL DATA

- **The First International Conference on Big Data, Small Data, Linked Data and Open Data**
- **ALLDATA 2015**
- **April 19 - 24, 2015 - Barcelona, Spain**

<http://www.iaria.org/conferences2015/ALLDATA15.html>

# ALLDATA 2015

- **BIG DATA**  
Big data foundations; Big data understanding (knowledge discovery, learning, consumer intelligence); Big data semantics, search and mining; Big data processing and transformations; Big data handling, simulation, visualization, modeling tools, and algorithms; Managing big data (large-scale, integration, etc.); Unknown in large Data Graphs; Reasoning on Big data; Big data analytics for prediction; Applications of Big data (health, financial, social, weather forecasting, etc.); Business-driven Big data; Big data and cloud technologies; Technologies handling Big data; High performance computing on Big data; Big data persistence and preservation; Big data protection, integrity and privacy ; Big data toolkits
- **SMALL DATA**  
Social networking small data; Relationship between small data and big data; Statistics on Small data; Handling Small data sets; Predictive modeling methods for Small data sets; Small data sets versus Big Data sets; Small and incomplete data sets; Normality in Small data sets; Confidence intervals of small data sets; Causal discovery from Small data sets; Deep Web and Small data sets; Small datasets for benchmarking and testing; Validation and verification of regression in small data sets; Small data toolkits
- **LINKED DATA**  
RDF and Linked data; Deploying Linked data; Linked data and Big data; Linked data and Small data; Evolving the Web into a global data space via Linked data; Practical semantic Web via Linked data; Structured dynamics and Linked data sets; Quantifying the connectivity of a semantic Linked data; Query languages for Linked data; Access control and security for Linked data; Anomaly detection via Linked data; Semantics for Linked data; Enterprise internal data 'silos' and Linked data; Traditional knowledge base and Linked data; Knowledge management applications and Linked data
- **OPEN DATA**  
Open data structures and algorithms; Designing for Open data; Open data and Linked Open data; Open data government initiatives; Big Open data; Small Open data; Challenges in using Open data (maps, genomes, chemical compounds, medical data and practice, bioscience and biodiversity); Linked open data and Clouds; Private and public Open data; Culture for Open data or Open government data; Data access, analysis and manipulation of Open data; Open addressing and Open data; Specification languages for Open data; Legal aspects for Open data

# Q&A

**Thanks!**