

A photograph of an MRI scanner room. A patient is lying on the scanner table, which is positioned inside the large circular gantry. The room is dimly lit, with light coming from a window with blinds in the background. A monitor is visible on the right side of the room.

**iMinds Dept.  
MEDICAL IT**

**KU Leuven ESAT-STADIUS**

**Serious Data Mining**

**Prof.Dr. Bart De Moor**

**[Bart.DeMoor@iminds.be](mailto:Bart.DeMoor@iminds.be)**

# Content

## **Big Data**

What

Who

## **Six issues**

Data

Compute Infrastructure

Storage Infrastructure

Analytics

Visualization

Security & Privacy

## **Machine learning as a commodity**

## **Expertise of ESAT-STADIUS, KU Leuven**

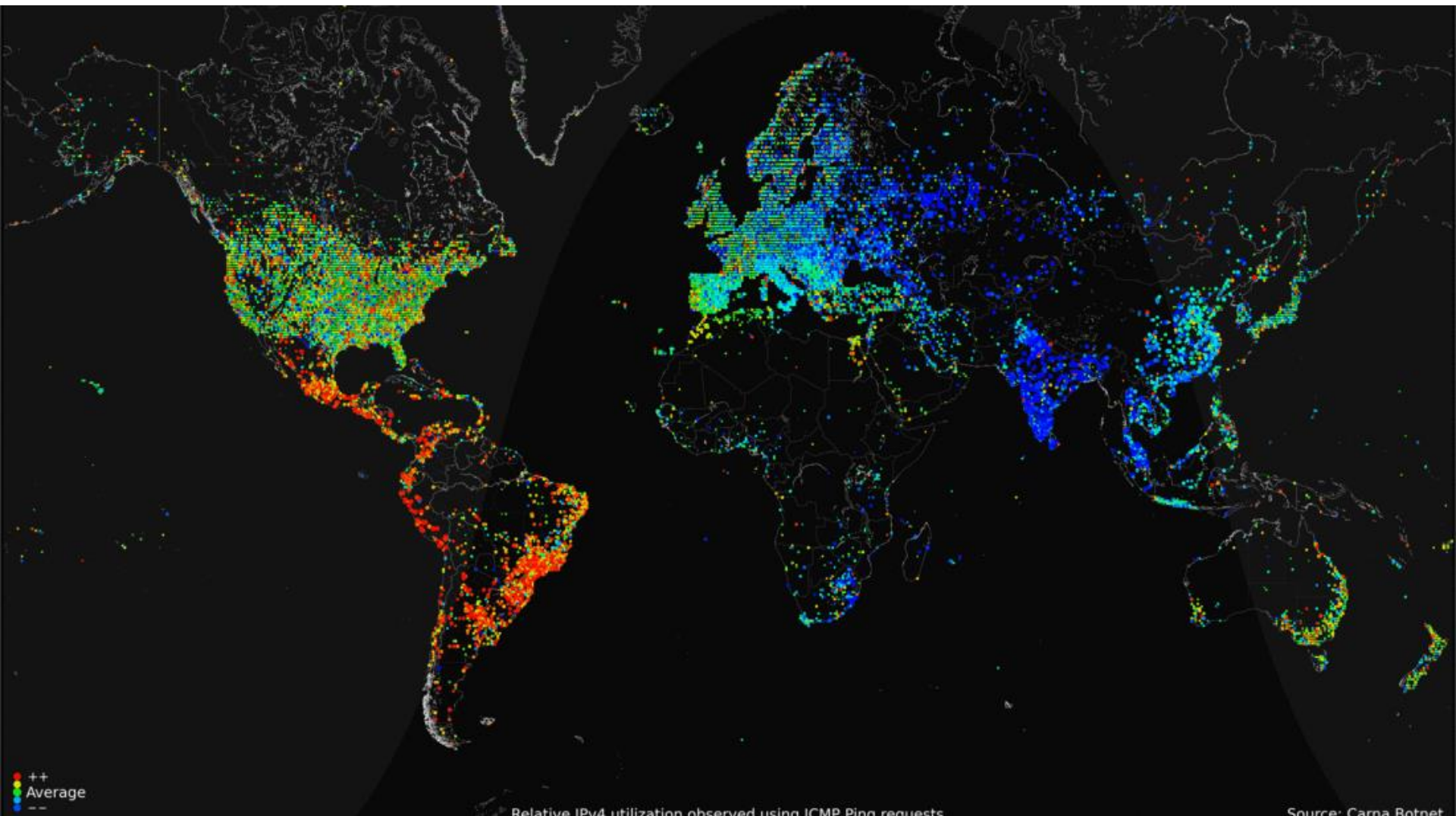
Books & Spin-offs

Algorithms

Applications



# WWW







Grains of rice the world consumes annually: **27.5 quadrillion**



Amount of data the world consumes every 30 minutes: **40.4 petabytes**

**We consume more bytes on the internet in 30 minutes than grains of rice in a year.**

1 million = 1 000 000  
1 billion = 1 000 000 000  
1 trillion = 1 000 000 000 000  
1 quadrillion =  
1 000 000 000 000 000

1 kB = 1 000  
1 MB = 1 000 000  
1 GB = 1 000 000 000  
1 TB = 1 000 000 000 000  
1 PB = 1 000 000 000 000 000

1 TB  
= large university library  
= 212 DVD discs  
= 1430 CDs  
= 3 year music in CD quality





The Industrial Internet, a connected network of intelligent machines working the way they are intended, will transform business as dramatically as the consumer Internet has changed our lives.



## The Industrial Internet

1 2 5 0 0 0 0 0 0 0 0 0

Connected Devices



Devices Per Capita Worldwide



2000 - 2010: there was rapid growth in connectivity and its transformative impact on the world have laid the foundation for the Industrial Internet.







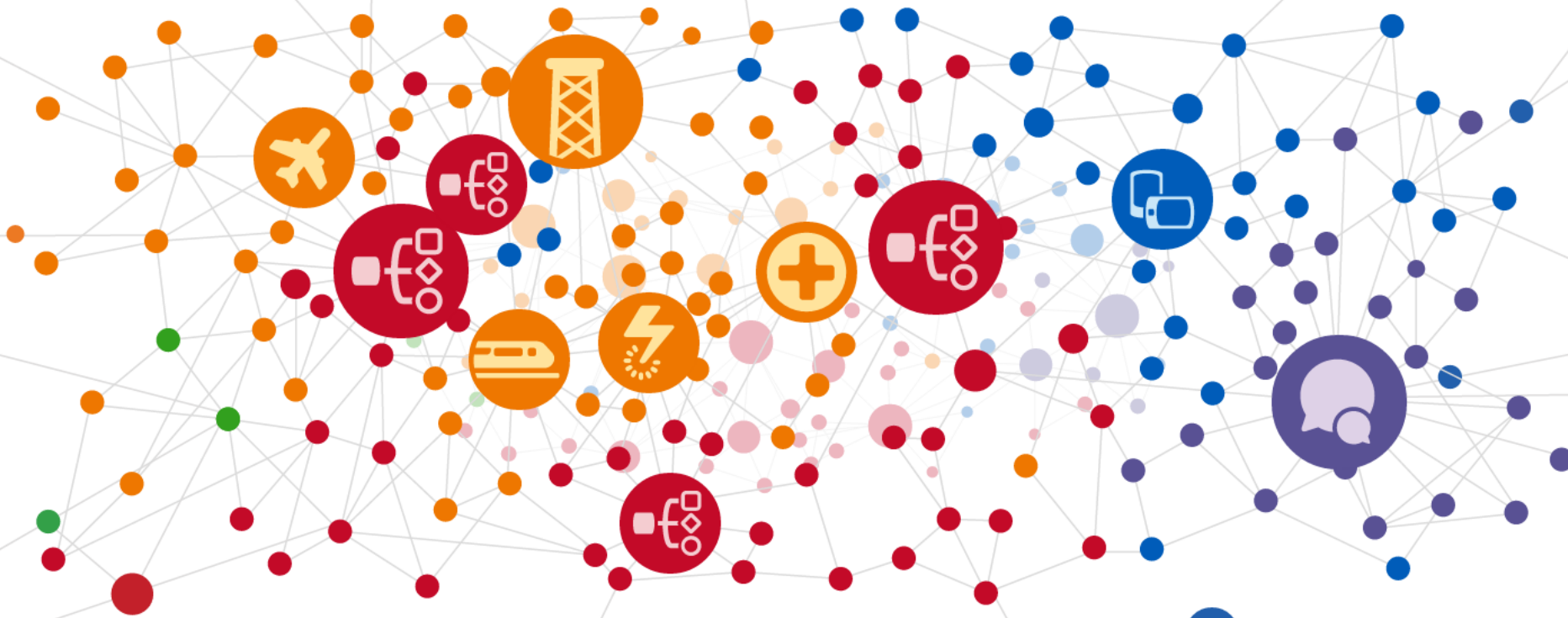
# The Industrial Internet

5 2 0 0 0 0 0 0 0 0

Connected Devices



Devices Per Capita Worldwide

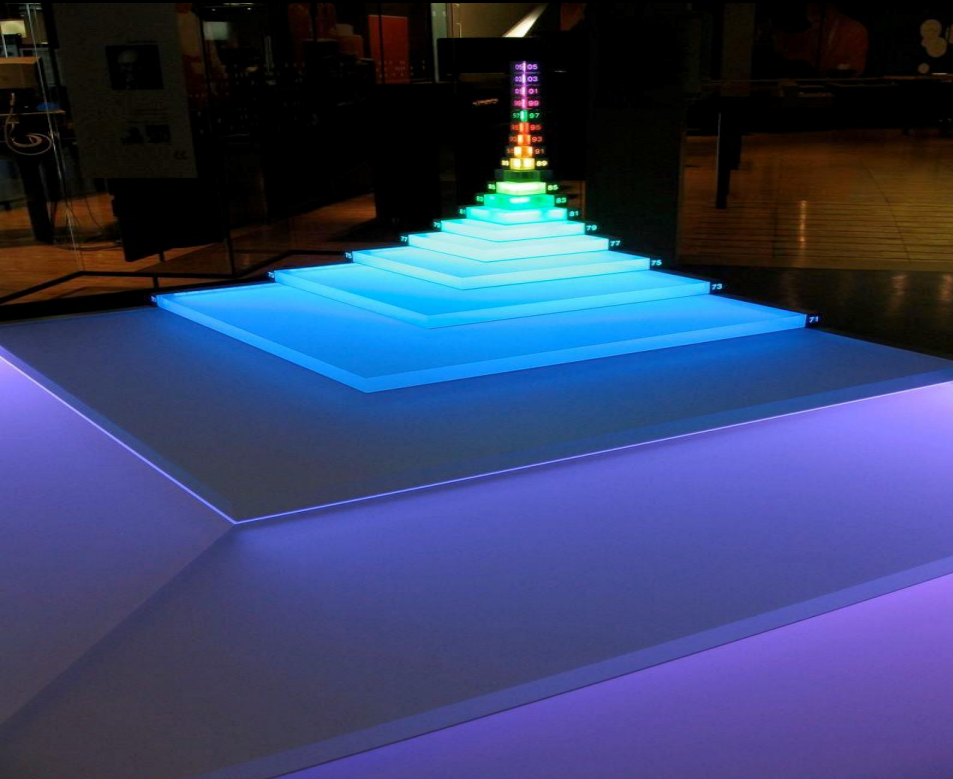


Past



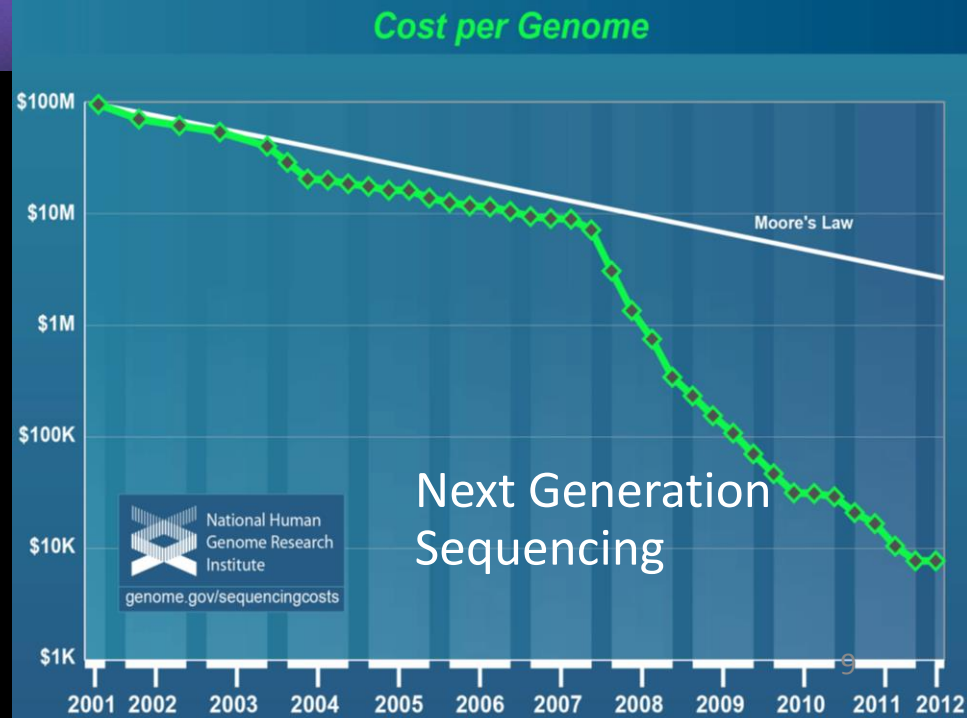
2020+

2015 - 2030: when the consumer and industrial Internet become one, merging minds and machines, we have the potential to connect 7B+ people and 50B assets to make the world work better.



**Moore's law:**  
 computing power  
 doubles  
 every 18 months

**Carlson's law:**  
 complexity/cost  
 evolves  
 exponentially



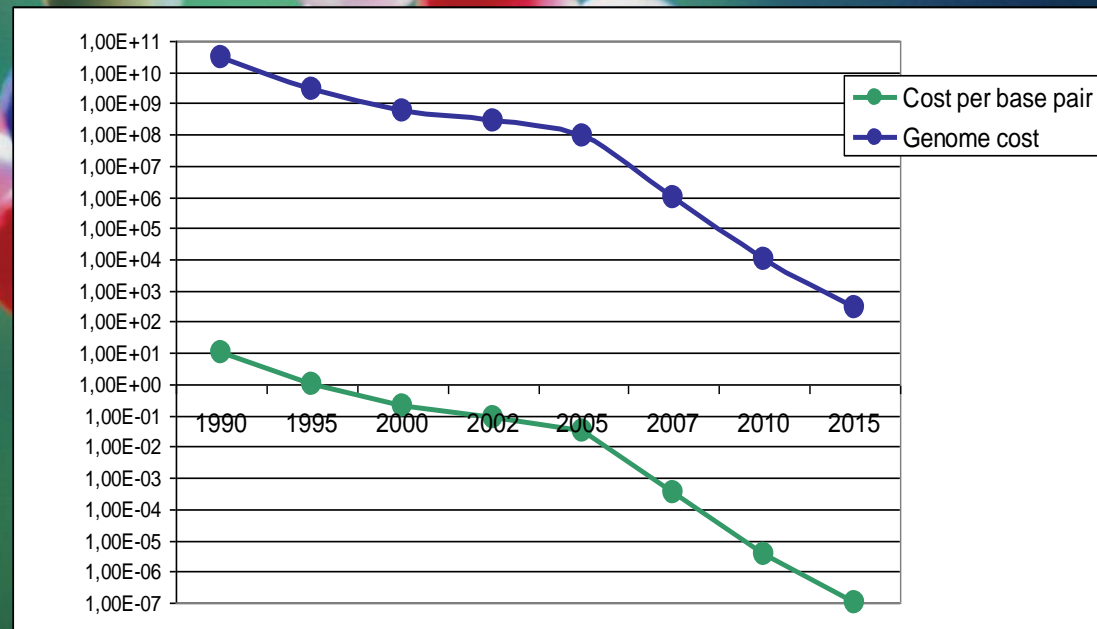
# Genome data

- **Human genome project (2003)**
  - 13 year project
  - \$300 million value with 2002 technology
- **Personal genome (2007)**
  - Genome of James Watson, 2 months
  - \$1 000 000
- **€1000-genome**
  - Expected 2012-2020



GS-FLX Roche  
Applied Science 454

Sequencers





# Tsunami of medical data

sequencing all newborns  
by 2020 (125k births /  
year)

125 PetaByte / year

index of 20  
million  
Biomedical  
PubMed  
records

23 GigaByte

raw NGS data  
of 1 full genome

1 TeraByte

PACS  
UZ Leuven

1,6 PetaByte

Genomics core  
HiSeq 2000 full  
speed exome  
sequencing

1 TeraByte / week

1 small  
animal  
image

1  
GigaByte

1 slice mouse  
brain MSI at  
10  $\mu$ m  
resolution

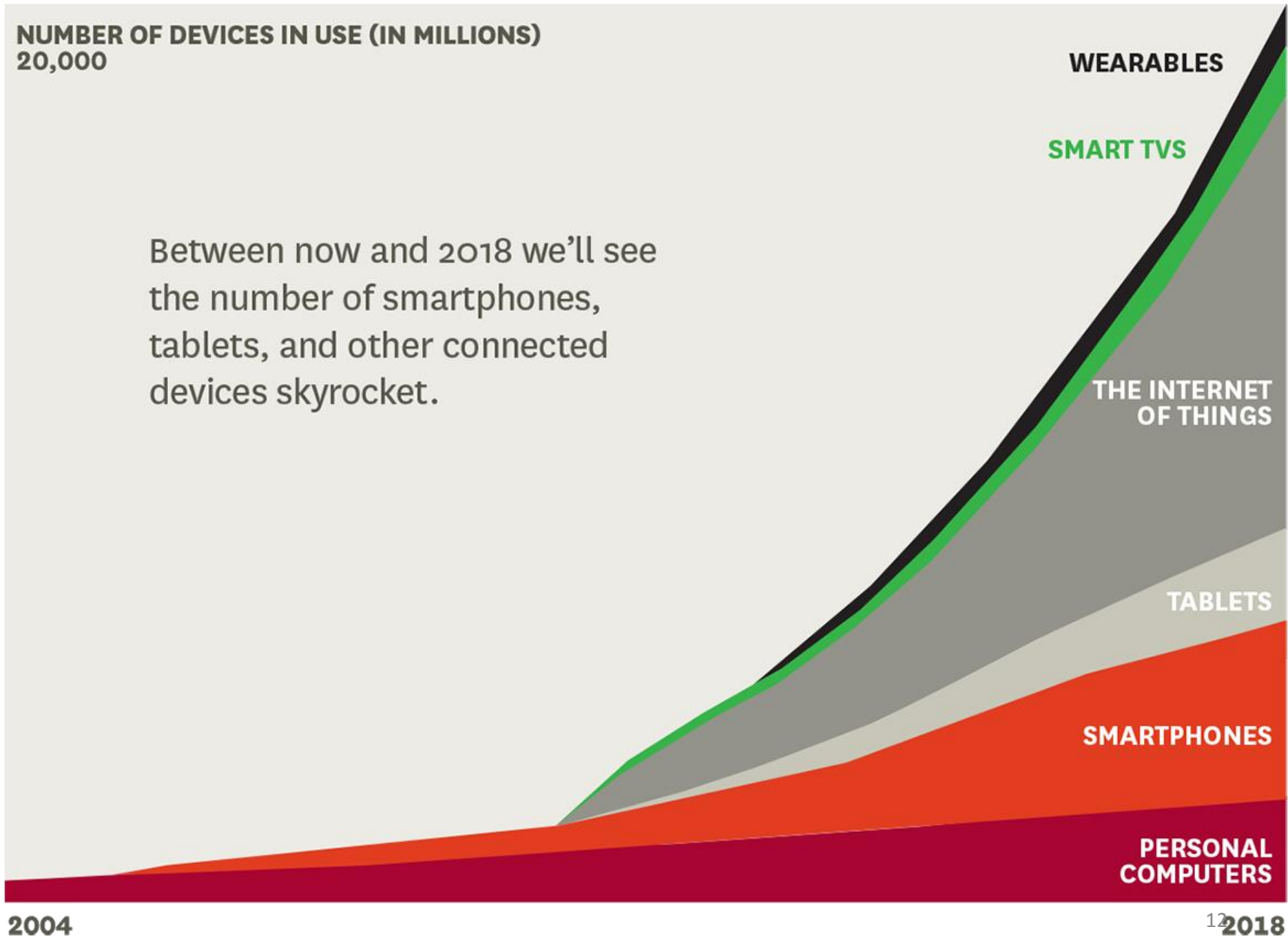
81 GigaByte

1 CD-ROM

750  
MegaByte

**NUMBER OF DEVICES IN USE (IN MILLIONS)**  
20,000

Between now and 2018 we'll see the number of smartphones, tablets, and other connected devices skyrocket.



2004

2018

# Data explosion in finance

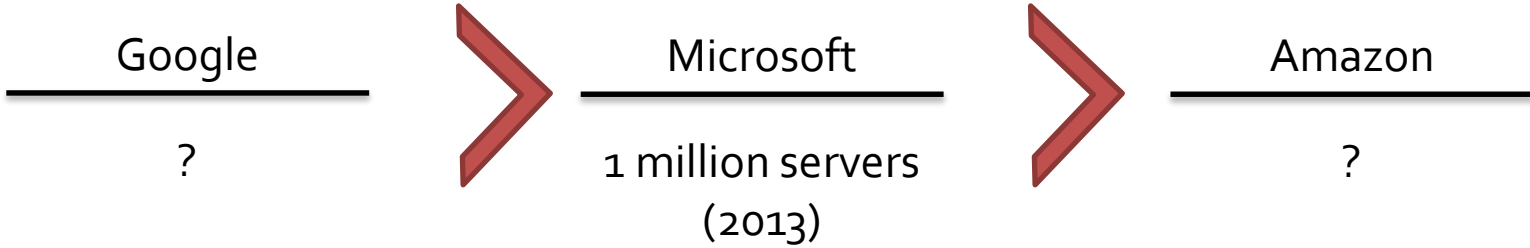


**Growing ~ 30-50% every year,  
half of this is unstructured!**

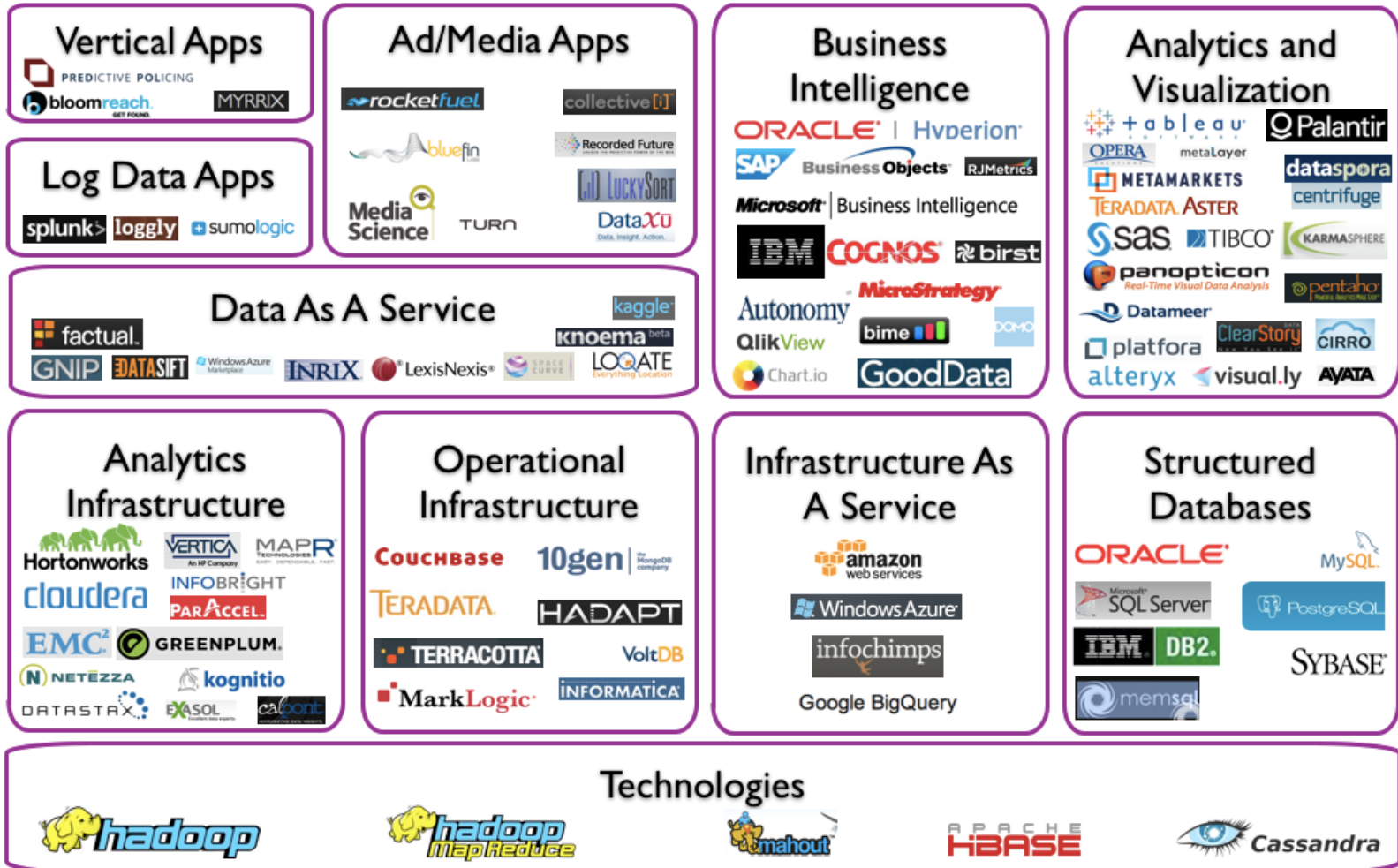
Data storage became 30%  
cheaper, yet budgets for data  
storage are still rising.



# How big is big?



# Big Data Landscape

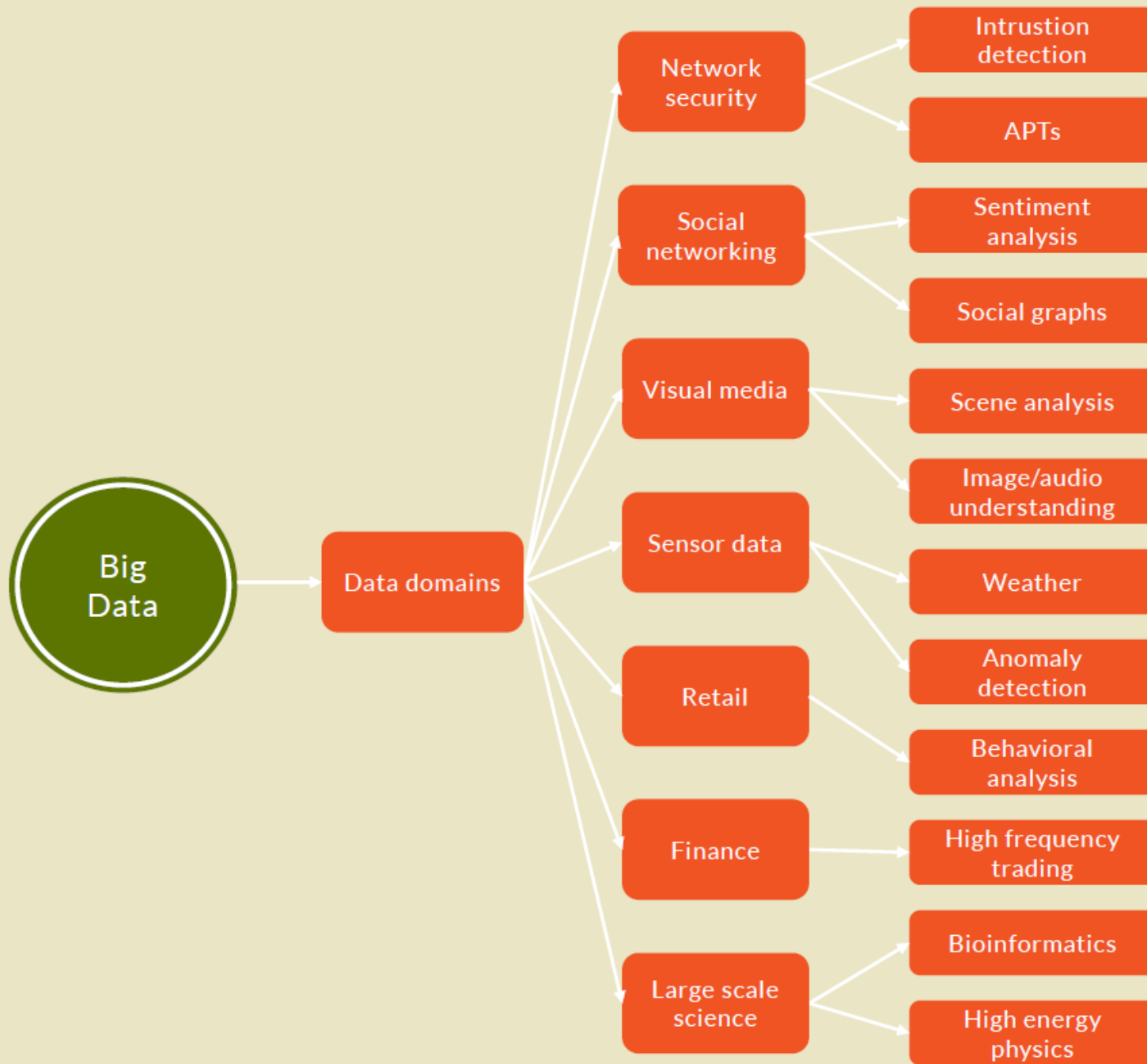


# Big Data Landscape (Version 2.0)



© Matt Turck (@mattturck) and ShivonZilis (@shivonz) Bloomberg Ventures







# Content

## **Big Data**

What

Who

## **Six issues**

Data

Compute Infrastructure

Storage Infrastructure

Analytics

Visualization

Security & Privacy

## **Machine learning as a commodity**

## **Expertise of ESAT-STADIUS, KU Leuven**

Books & Spin-offs

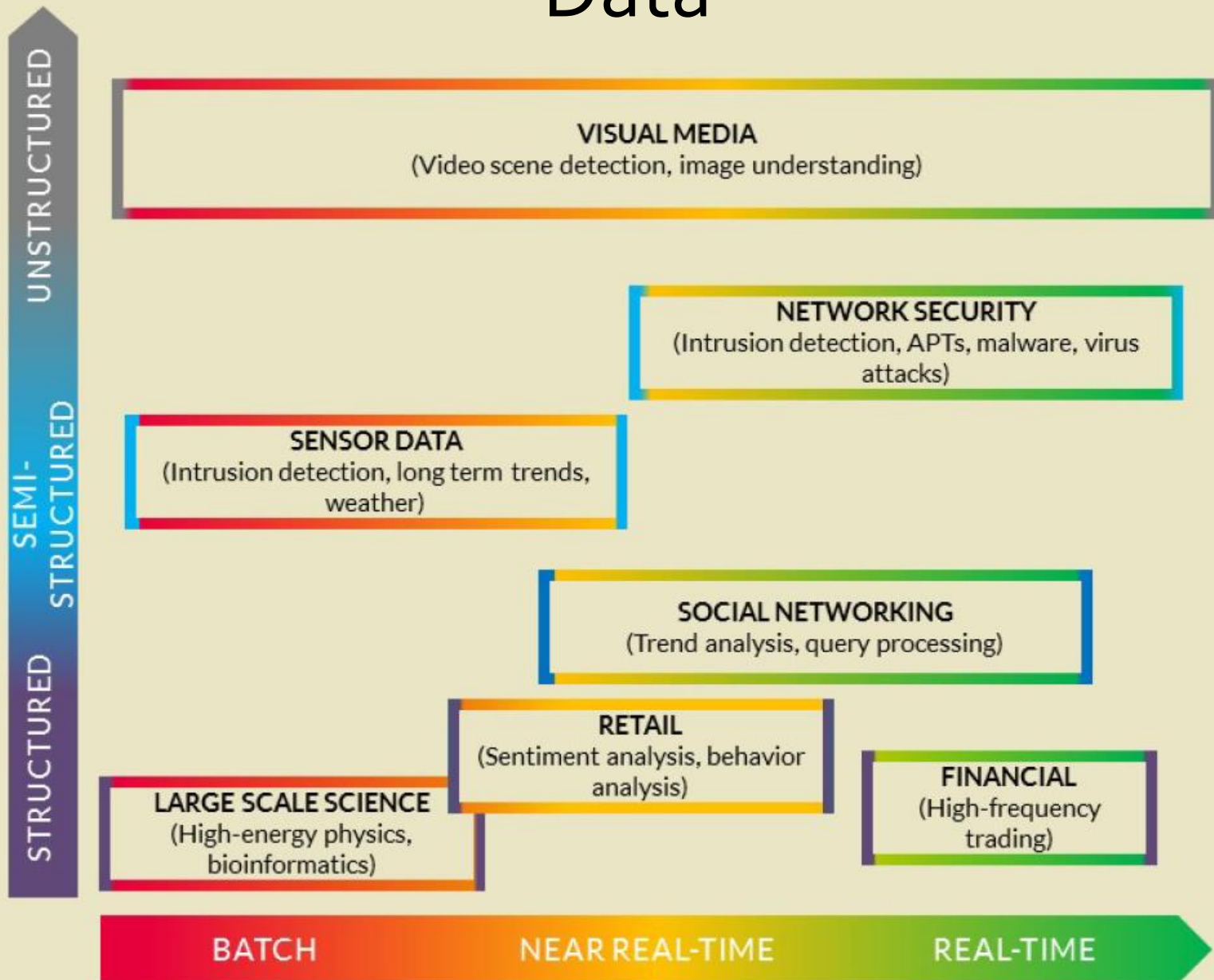
Algorithms

Applications

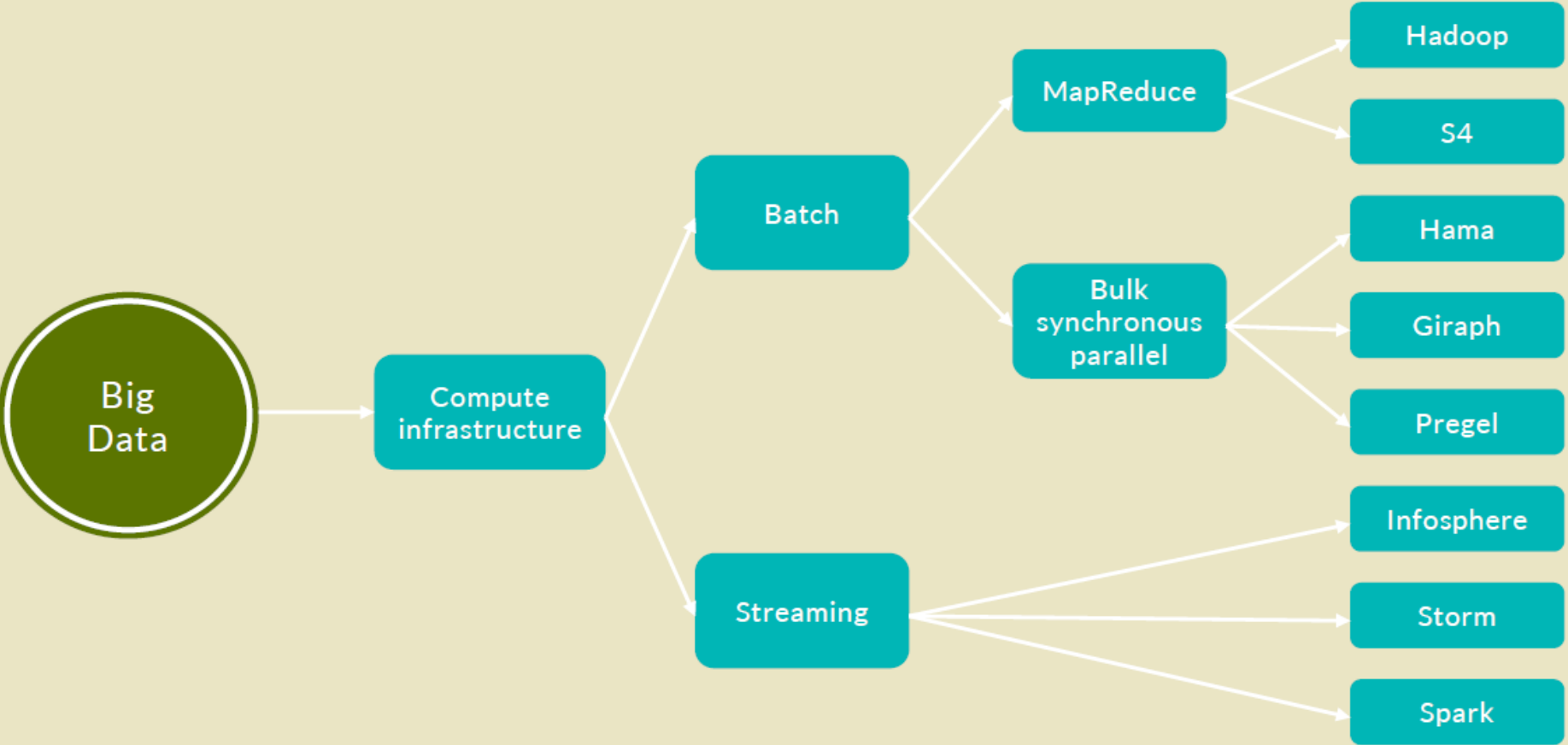




# Data

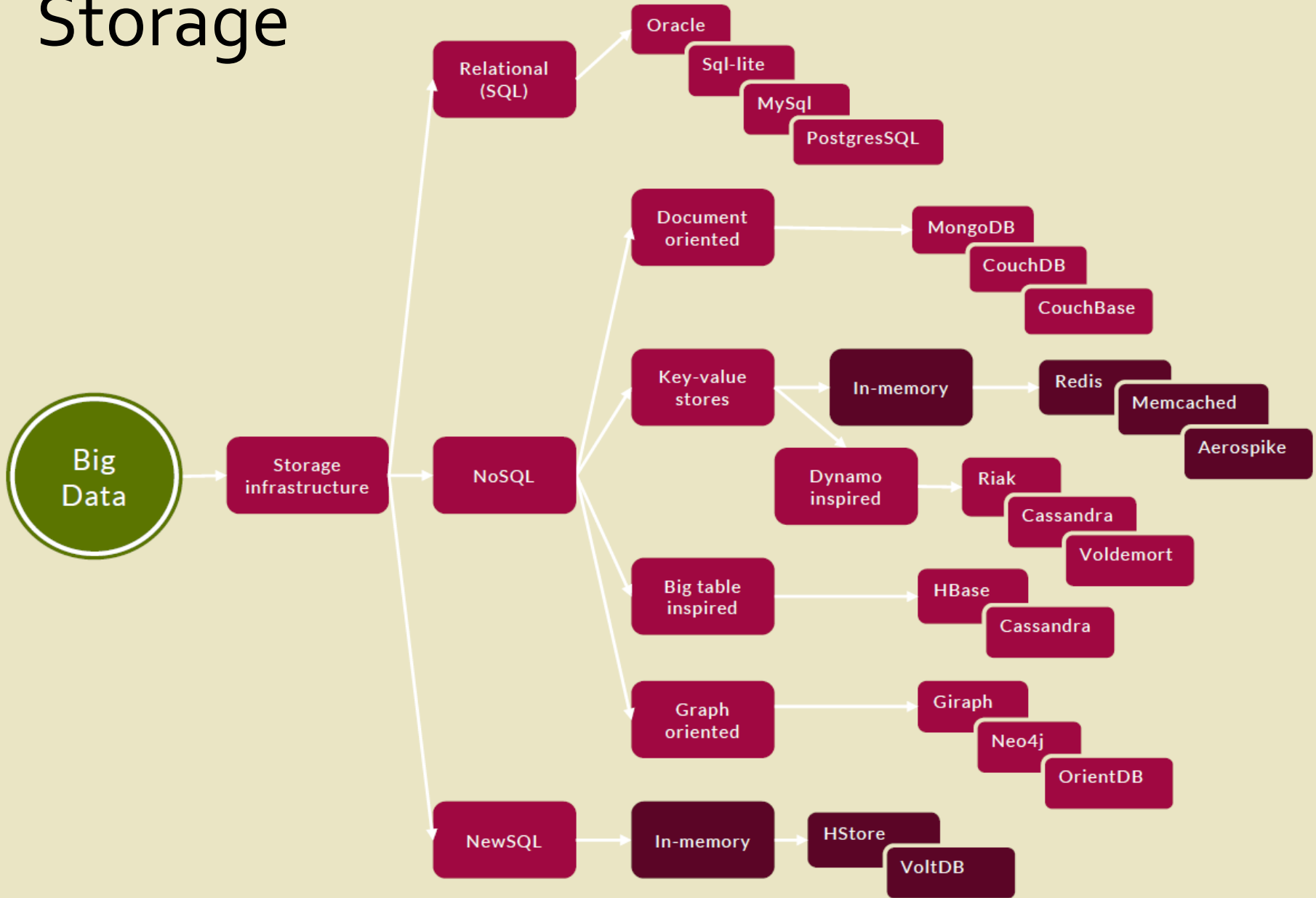


# Compute infrastructure

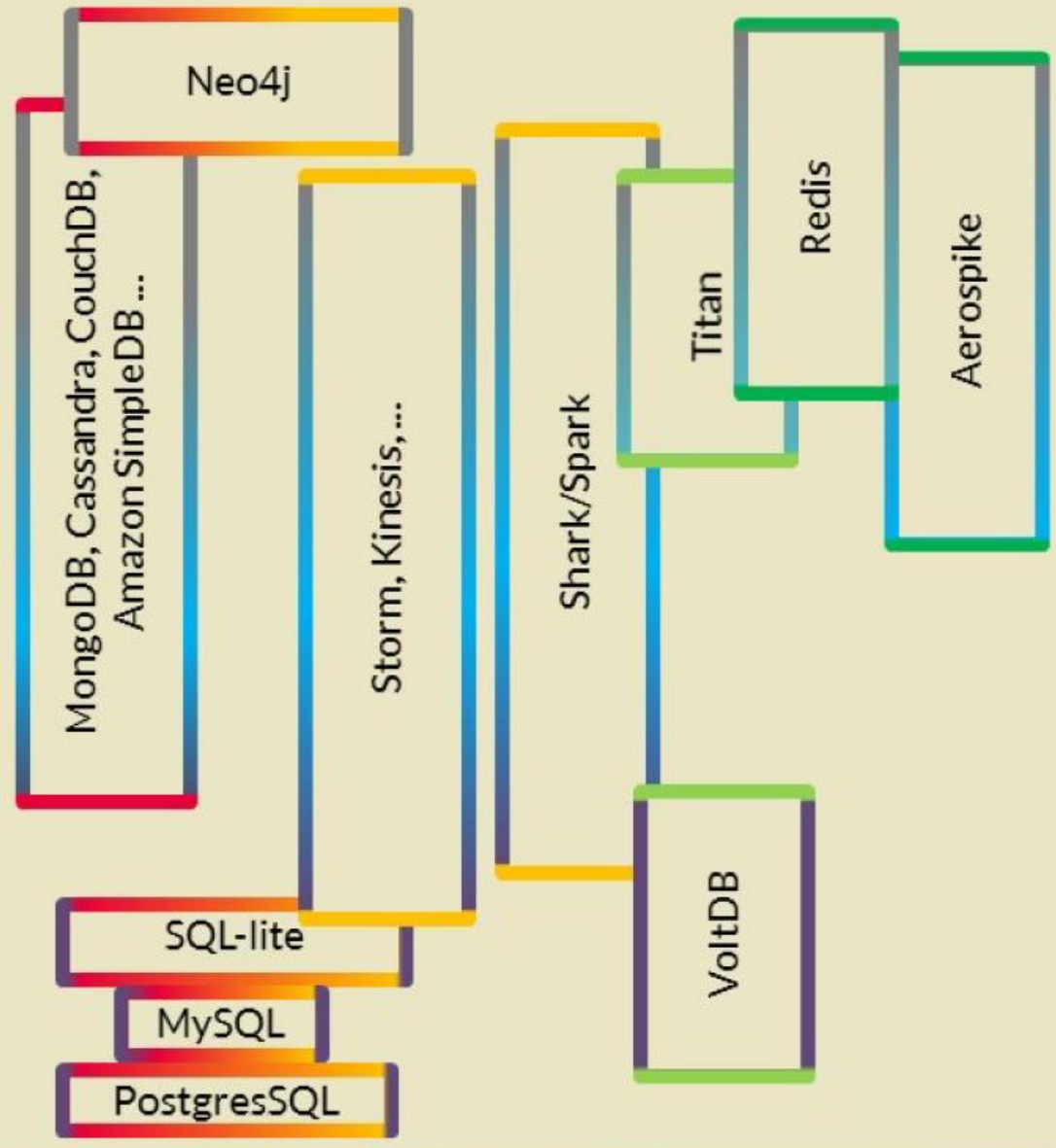




# Storage

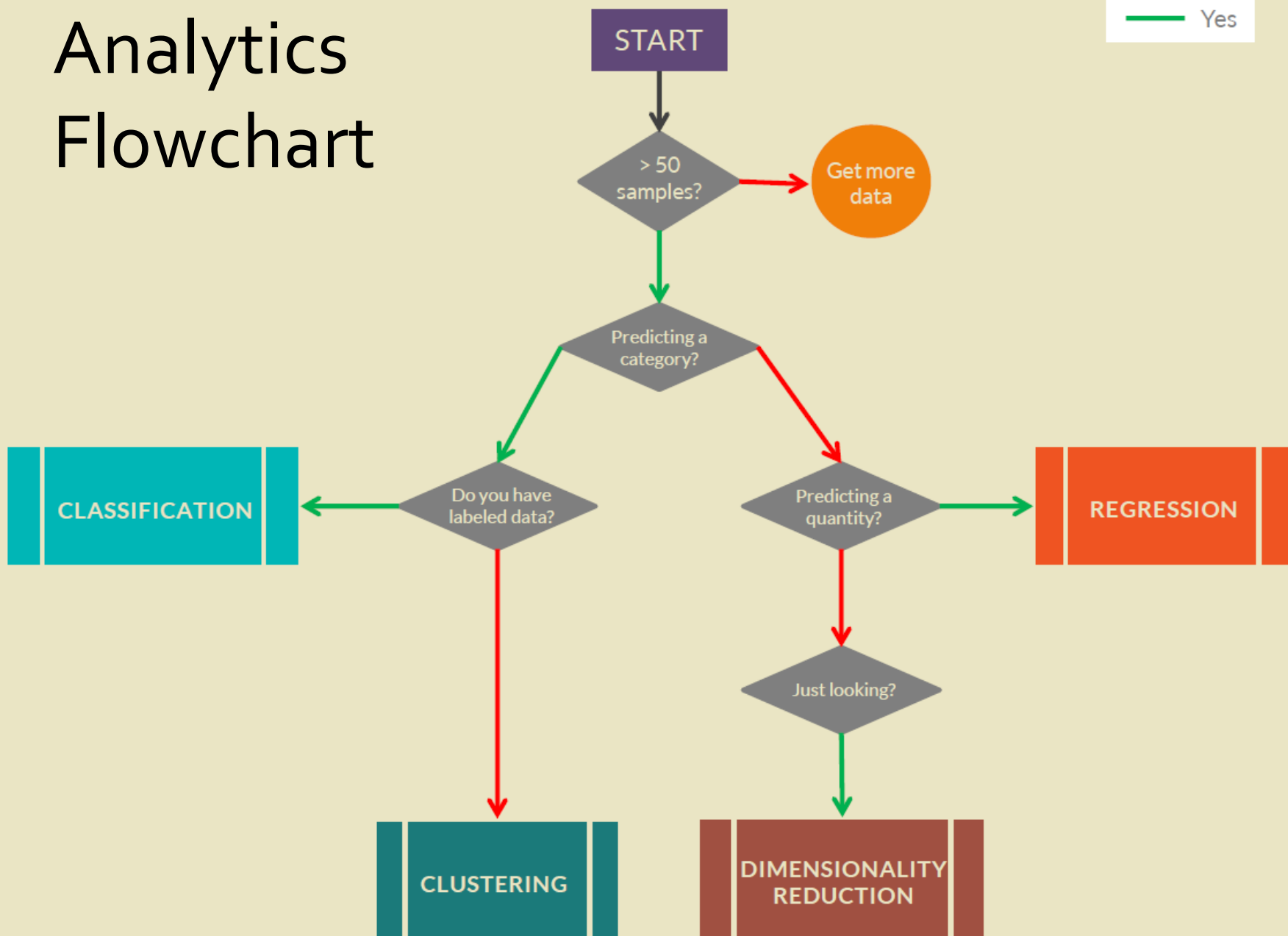


STRUCTURED      SEMI-STRUCTURED      UNSTRUCTURED

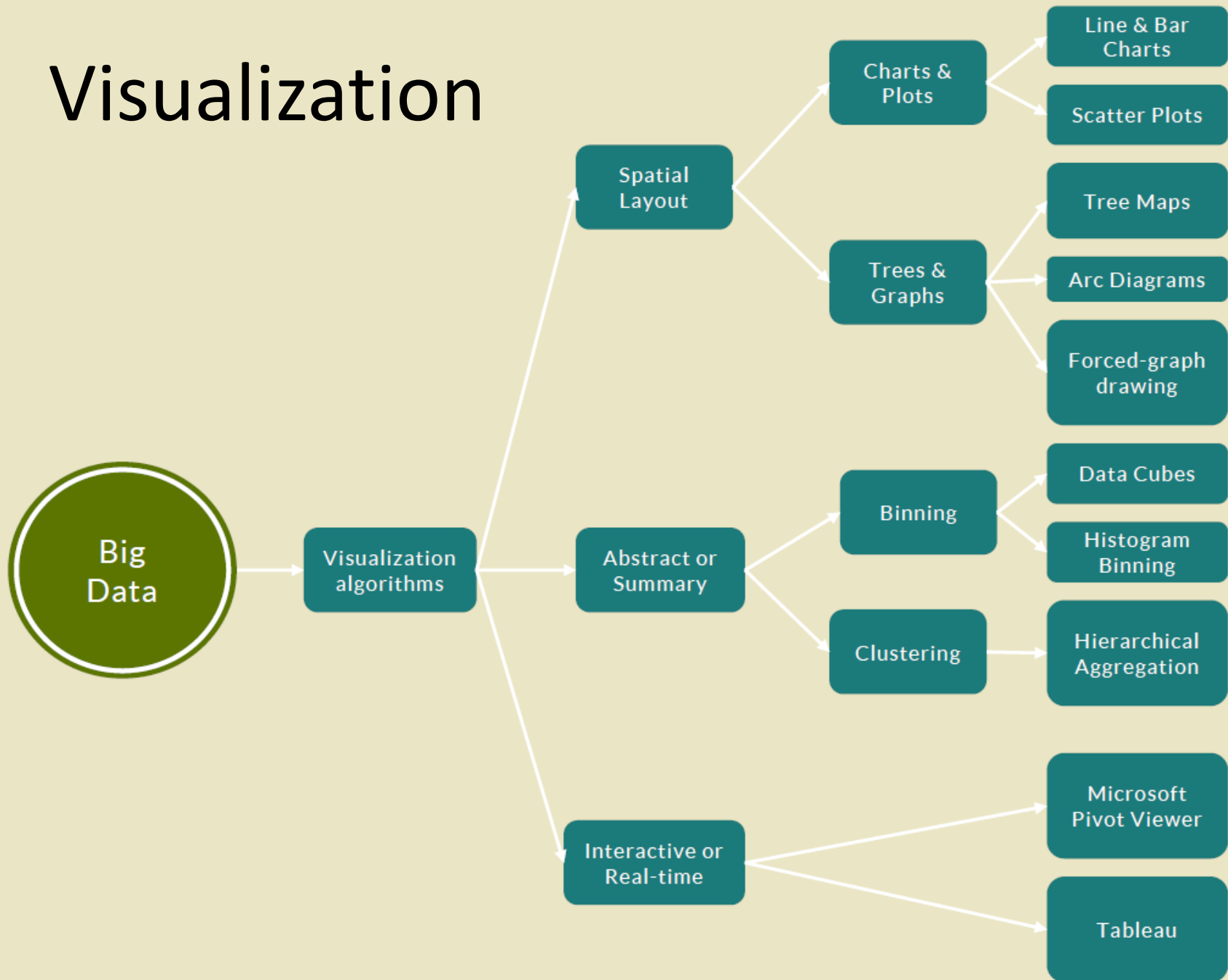


BATCH      NEAR REAL-TIME      REAL-TIME

# Analytics Flowchart

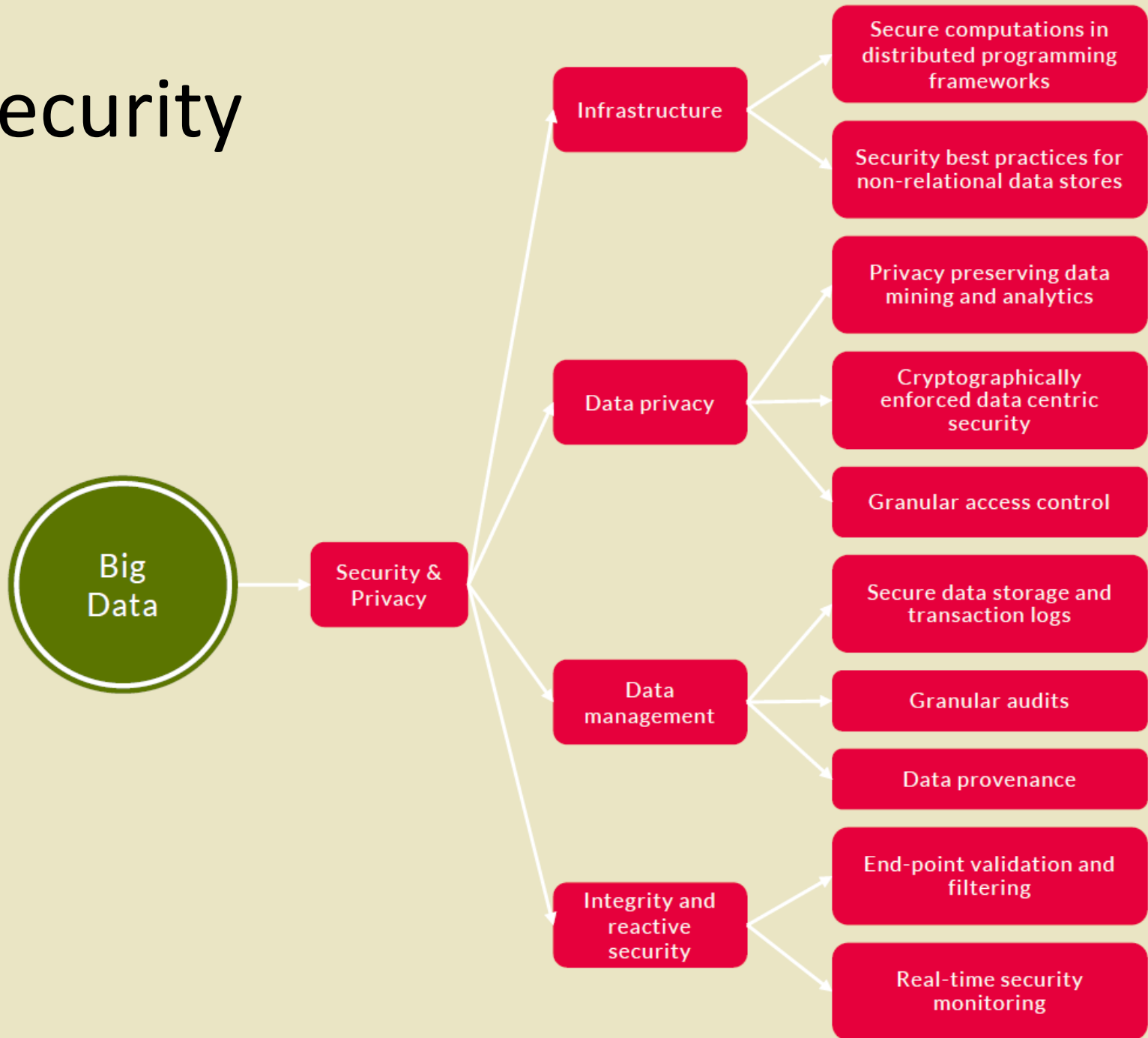


# Visualization





# Security



# Content

## Big Data

What

Who

## Six issues

Data

Compute Infrastructure

Storage Infrastructure

Analytics

Visualization

Security & Privacy

## Machine learning as a commodity

## Expertise of ESAT-STADIUS, KU Leuven

Books & Spin-offs

Algorithms

Applications

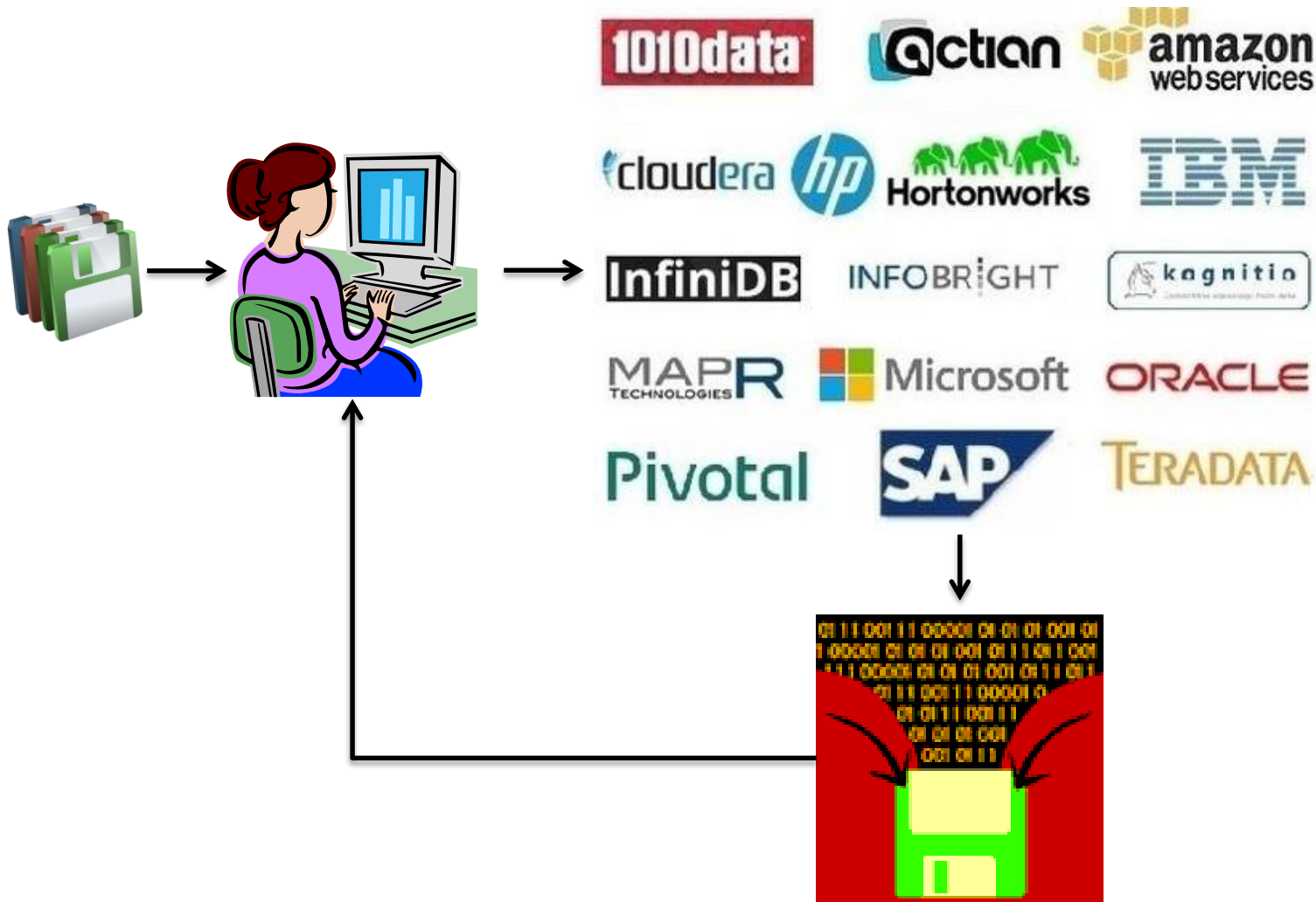


# Big Data Landscape



More and more analytics as a commodity!

# Machine Learning as a commodity





# Big Data Landscape



Many possible applications!

Energy

Industry

Environment

Social networks

Fraud and predictive analysis

Health

...

➔ Focus on Serious Big Data

# Content

## **Big Data**

What

Who

## **Six issues**

Data

Compute Infrastructure

Storage Infrastructure

Analytics

Visualization

Security & Privacy

## **Machine learning as a commodity**

## **Expertise of ESAT-STADIUS, KU Leuven**

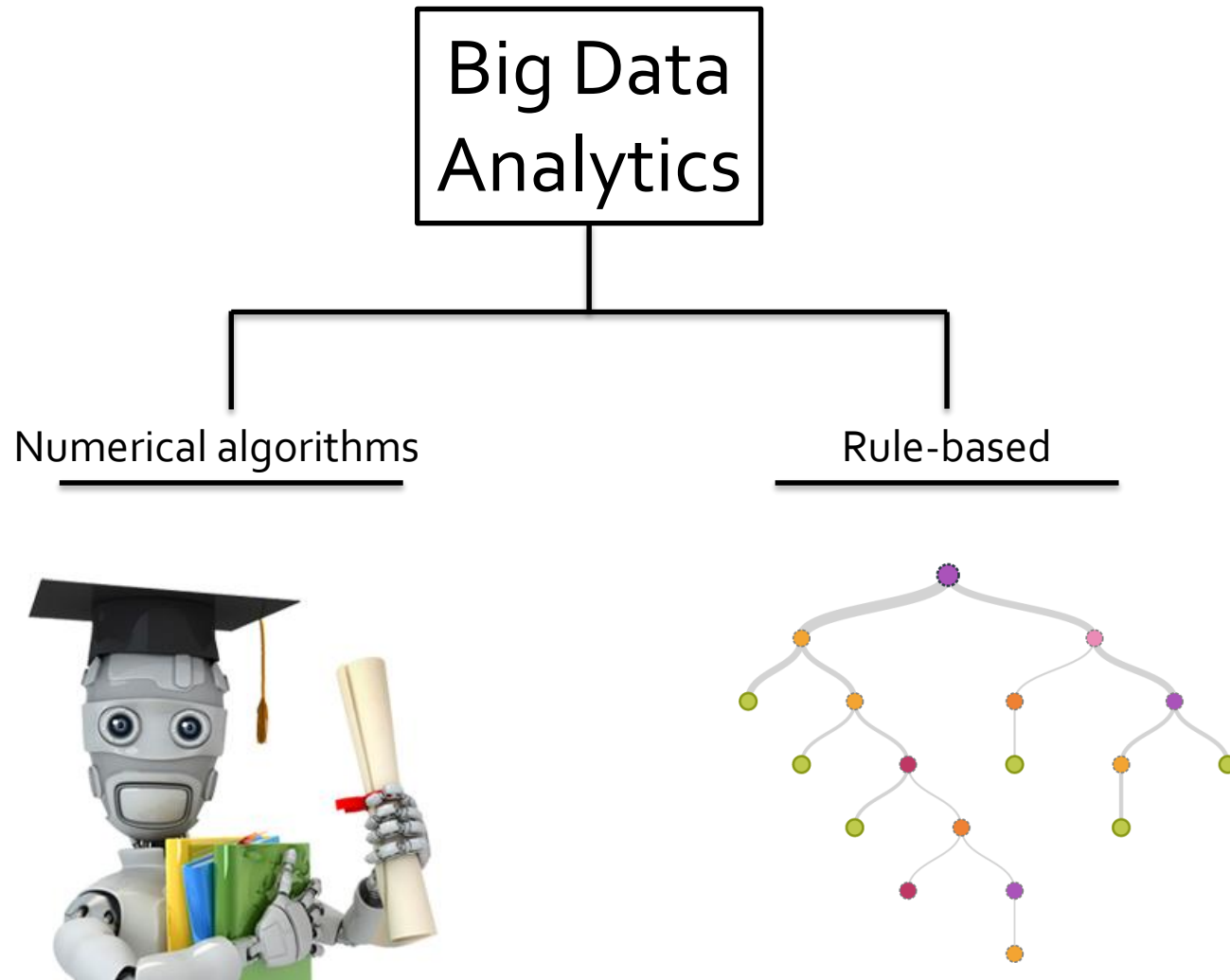
Books & Spin-offs

Algorithms

Applications



# Analytics



# Main tasks

## Prediction

---



Regression

## Segmentation

---



Clustering

Classification

## Anomalies

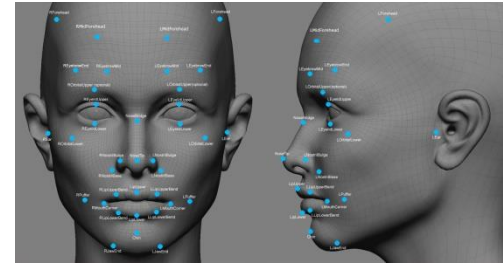
---



Outlier

# What can we do?

Face recognition



Fraud detection



Shopping cart analysis



Just-in-time production



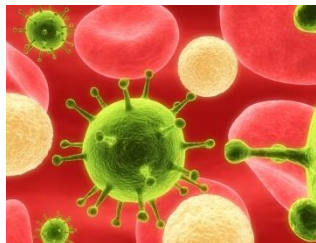
Credit worthiness



Traffic management



Disease spreading

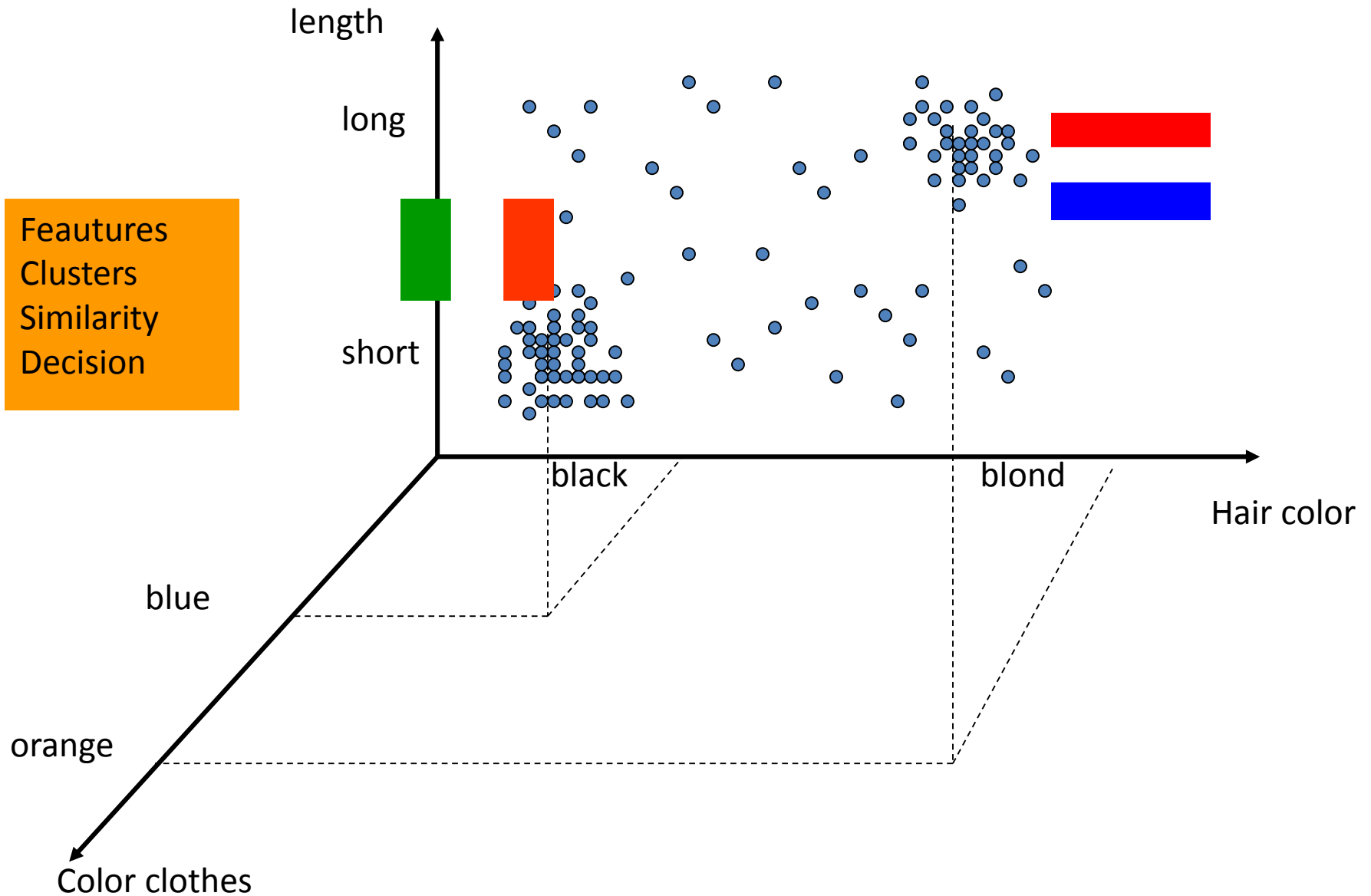


Movie recommendation

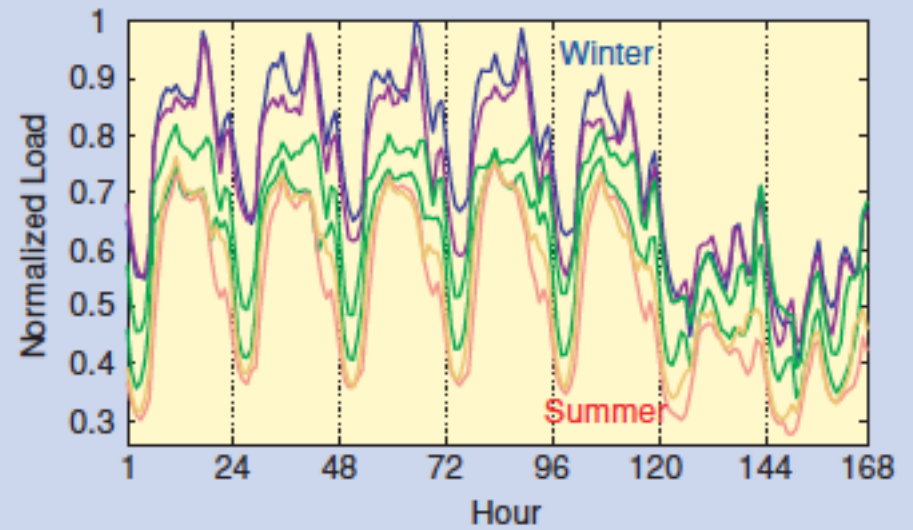
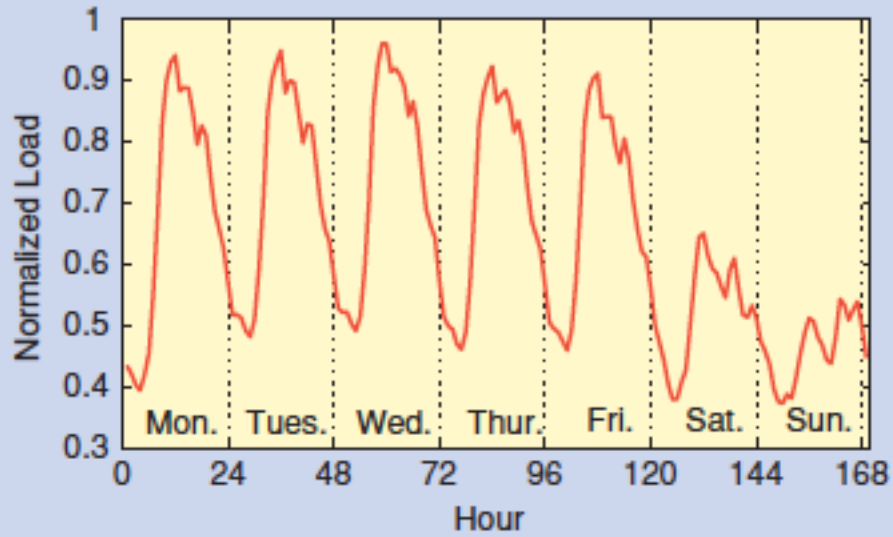




# Clustering/Classification

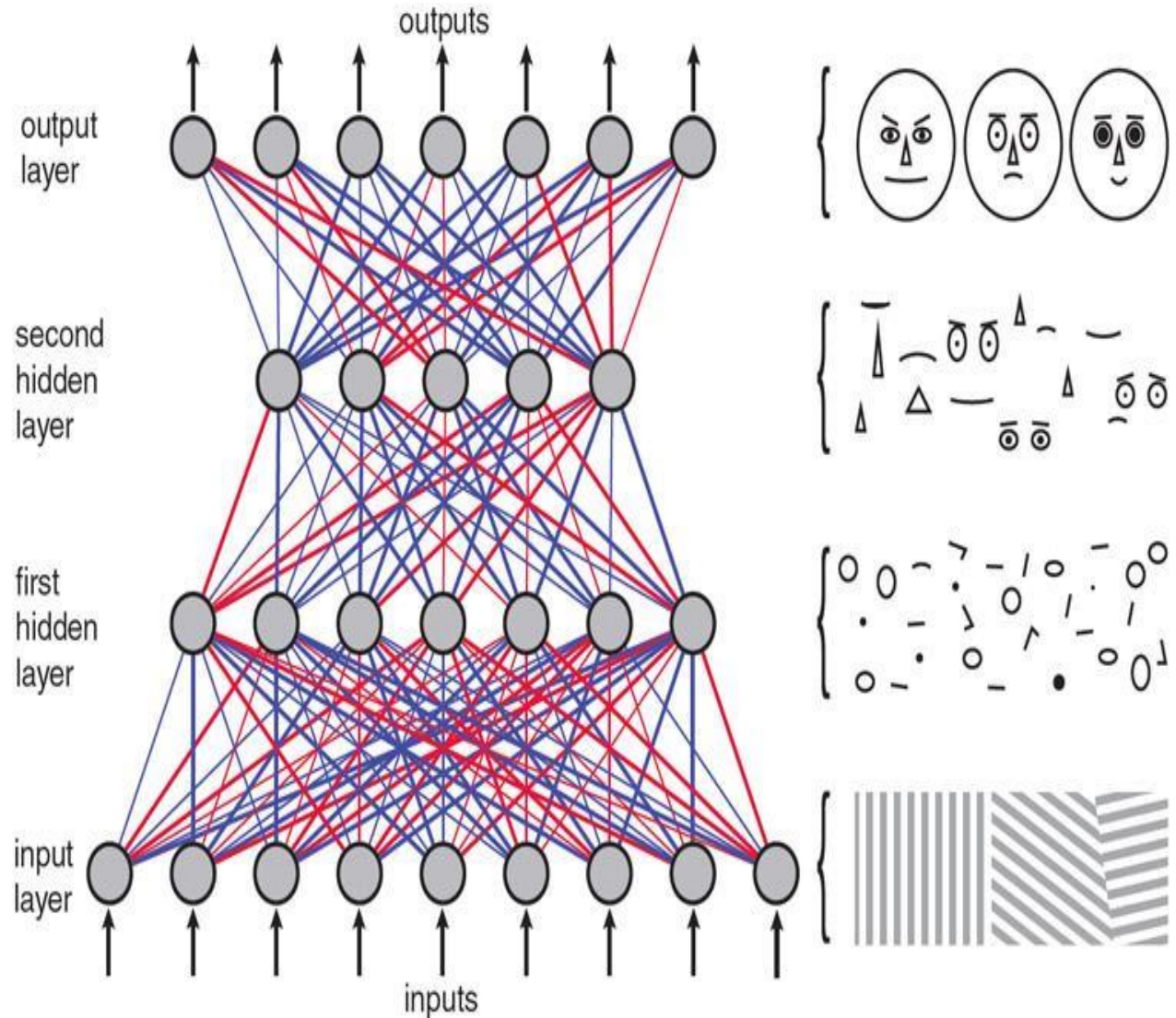


# Forecasting

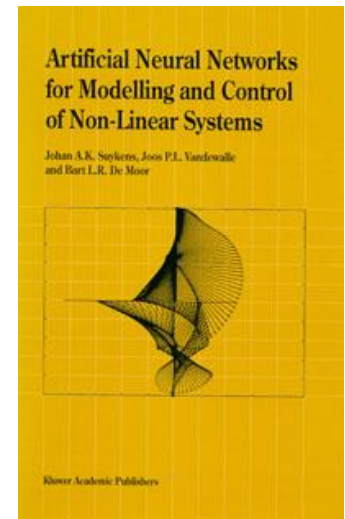
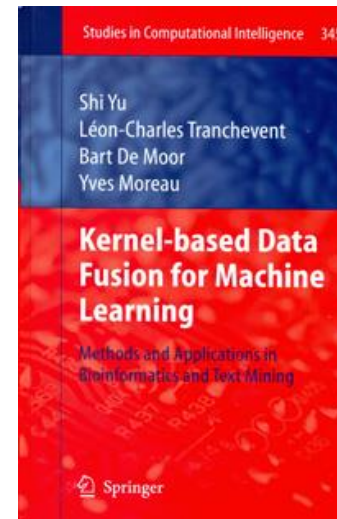
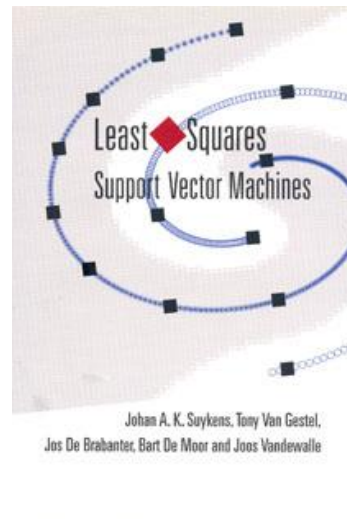
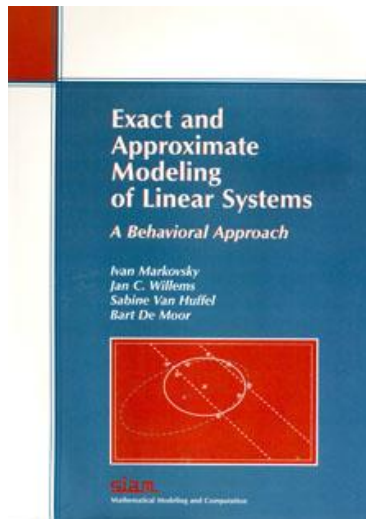
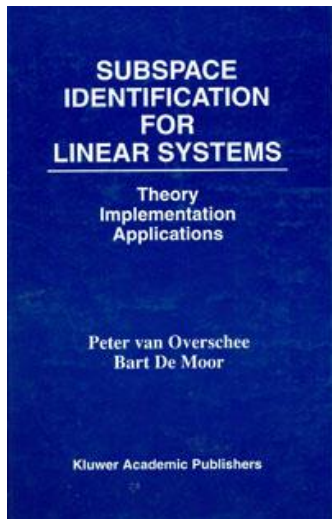


# Deep Learning

- Neural networks.
- New algorithms.
- Multiple layers on top of each other.
- Each layer learns a more complex representation.
- Learn feature hierarchies.



# Stadius - Books

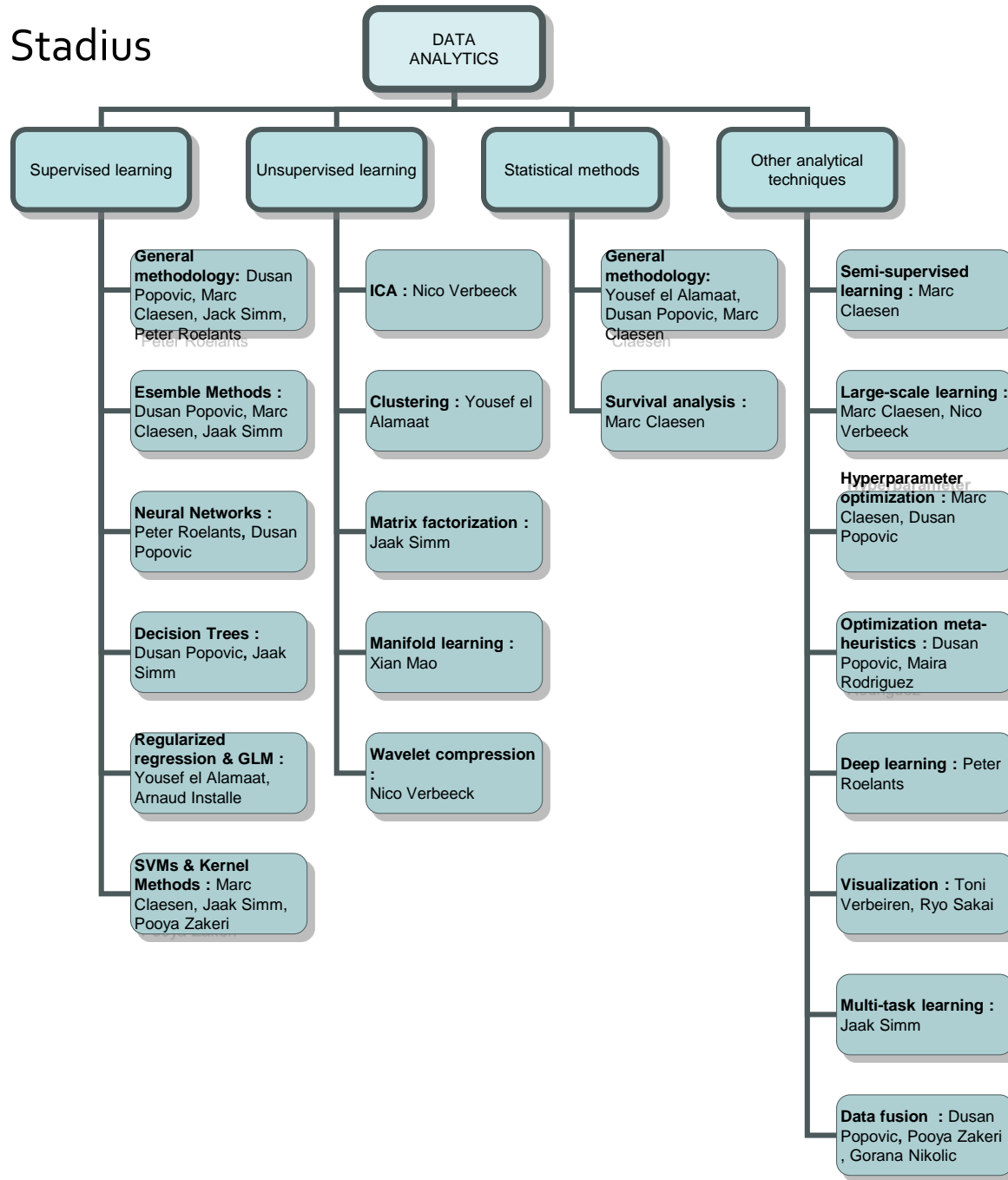


# Stadius - Spin-offs

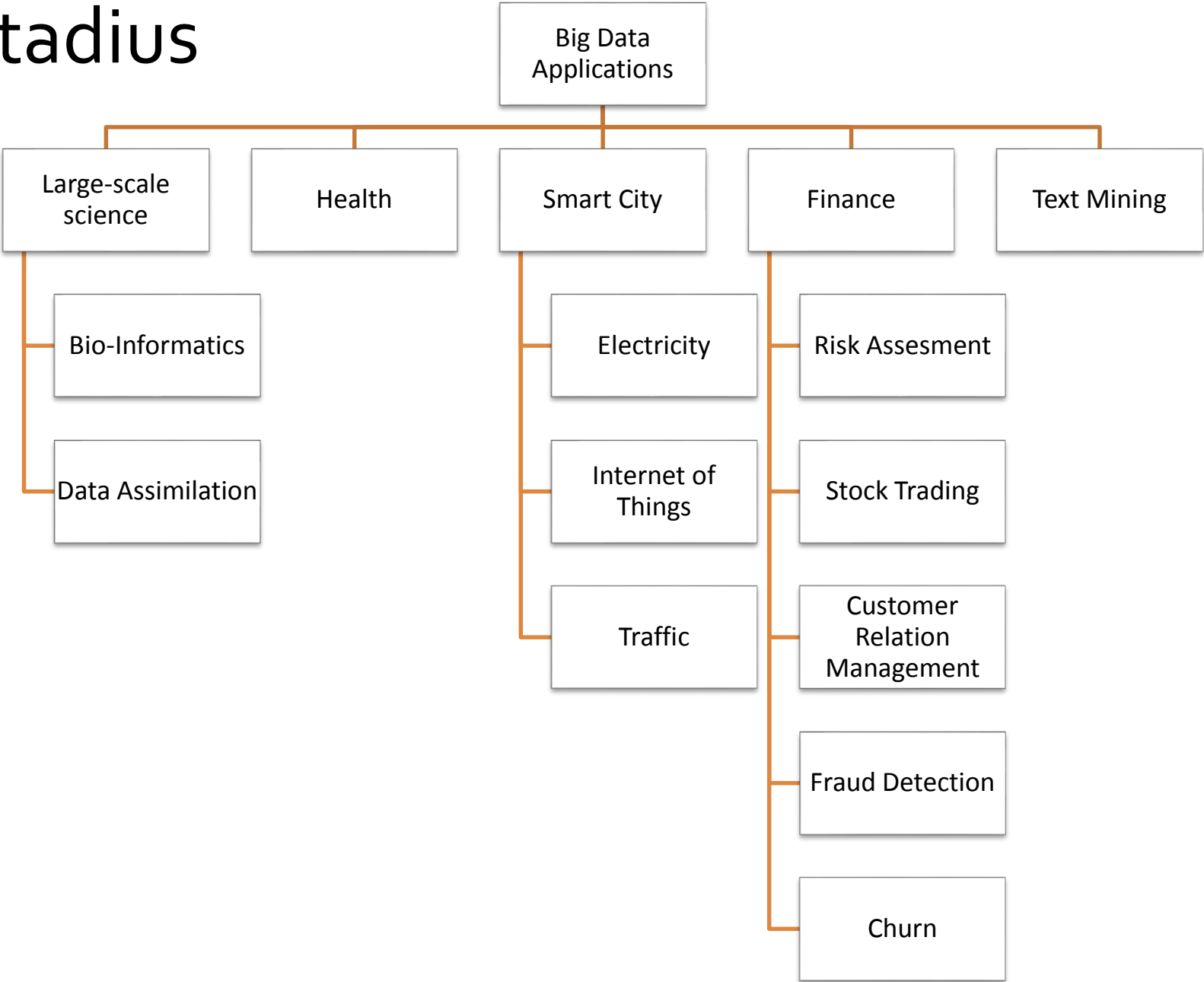


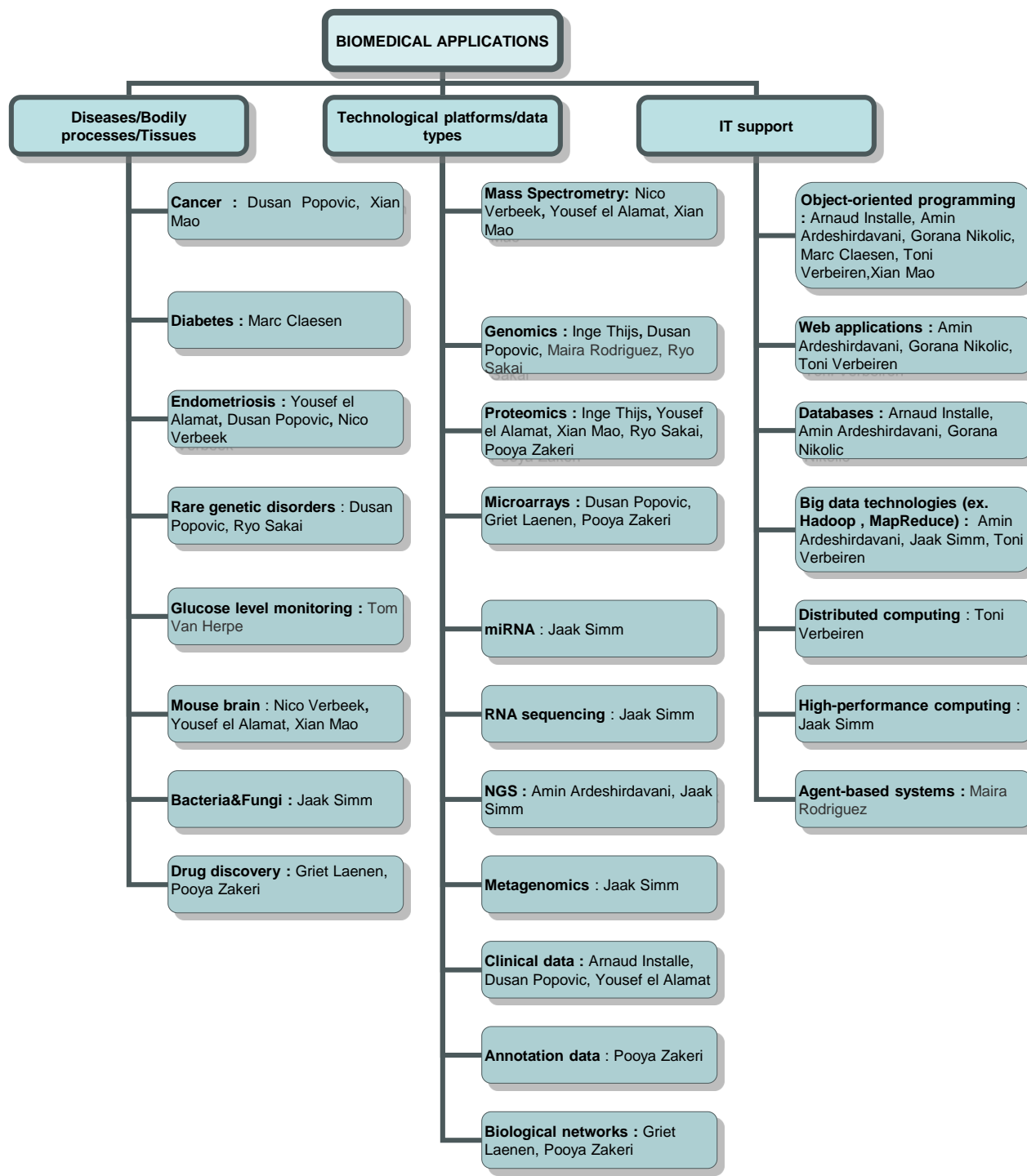


# Algorithms in Stadius



# Stadius





A woman with short brown hair and glasses, wearing a grey top, is pointing at a tablet held by a healthcare worker in blue scrubs. The healthcare worker has her hair in a bun and is wearing a stethoscope. The background is a plain, light-colored wall.

**Energy**

**Industry**

**Environment**

**Social networks**

**Fraud and predictive analysis**

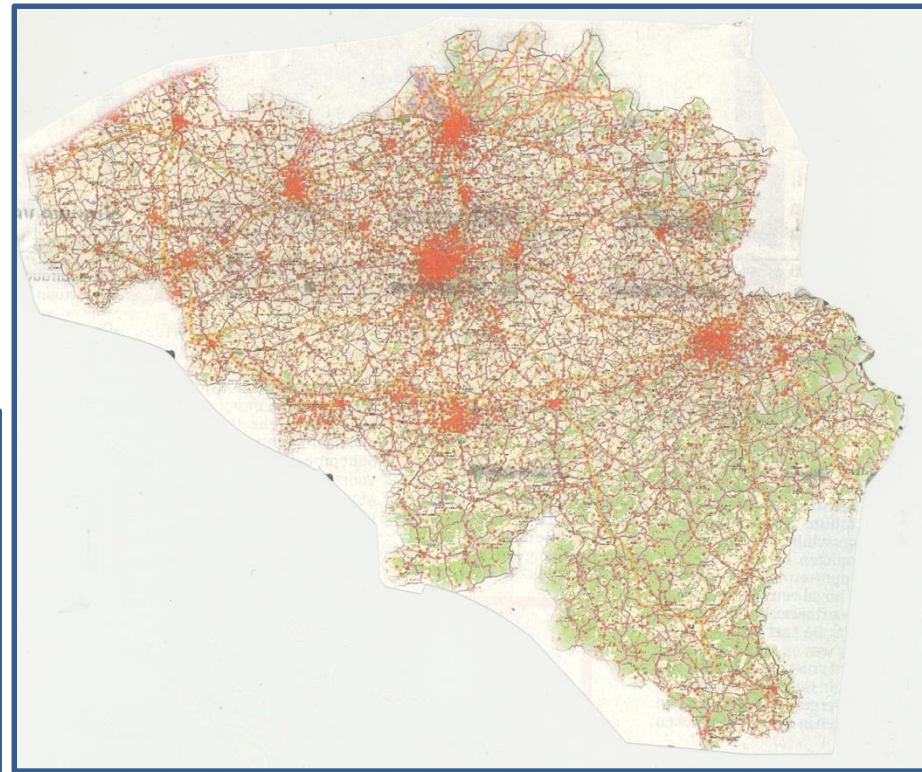
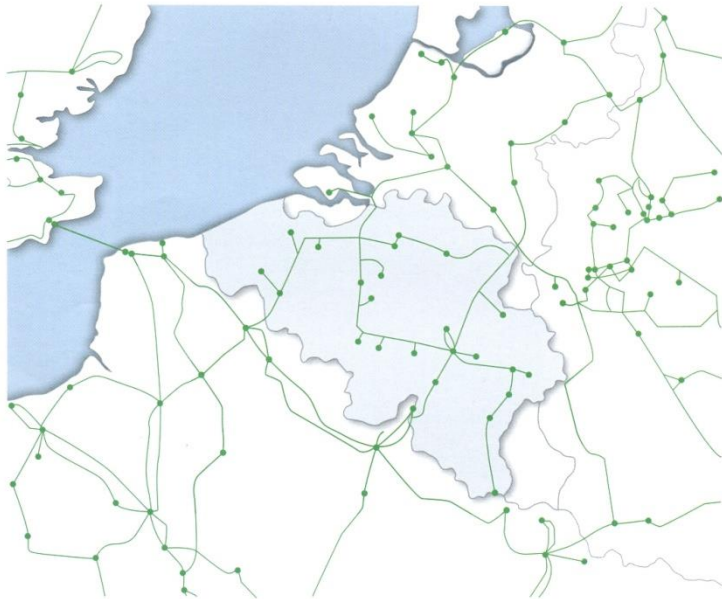
**Health**

# Power grid

## België en Europa

Het Elia-net:  
knooppunt van elektriciteitverkeer in Europa

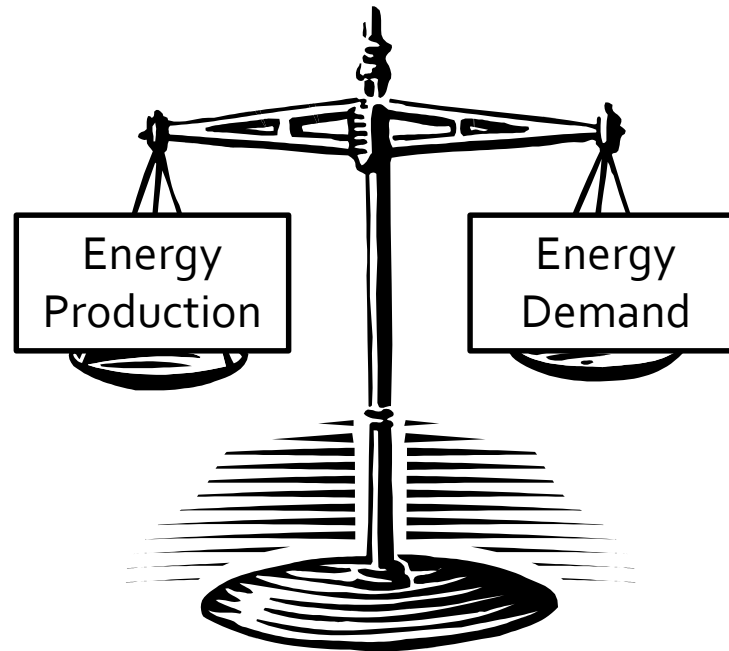
380 kV interconnectienet met hoogspanningsstations





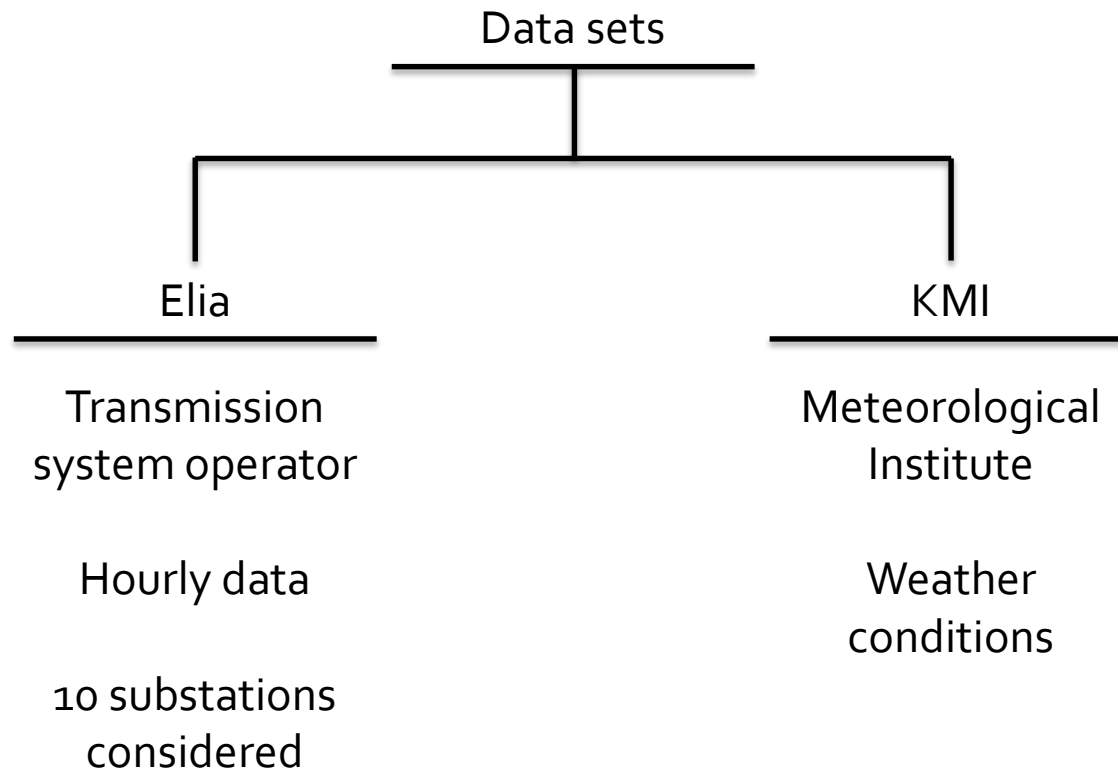
# Electric load forecasting

Problem

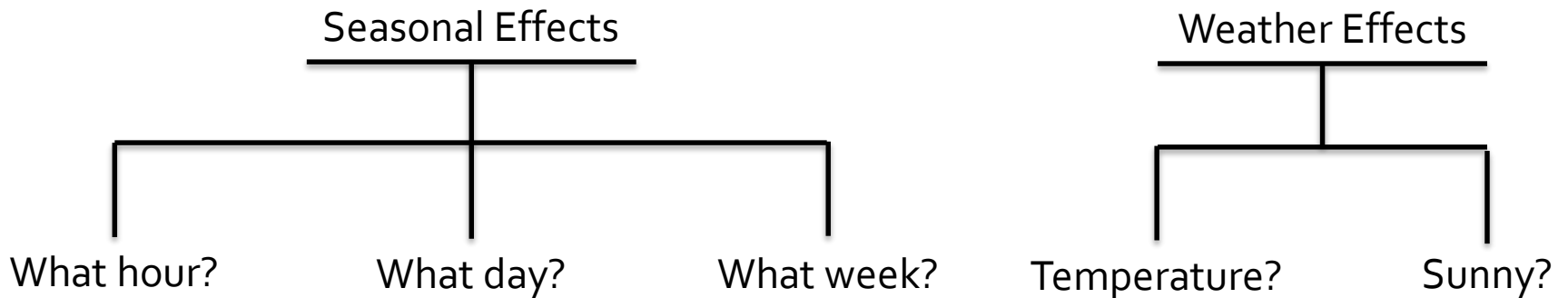


How to forecast  
the demand?

# Electric load forecasting

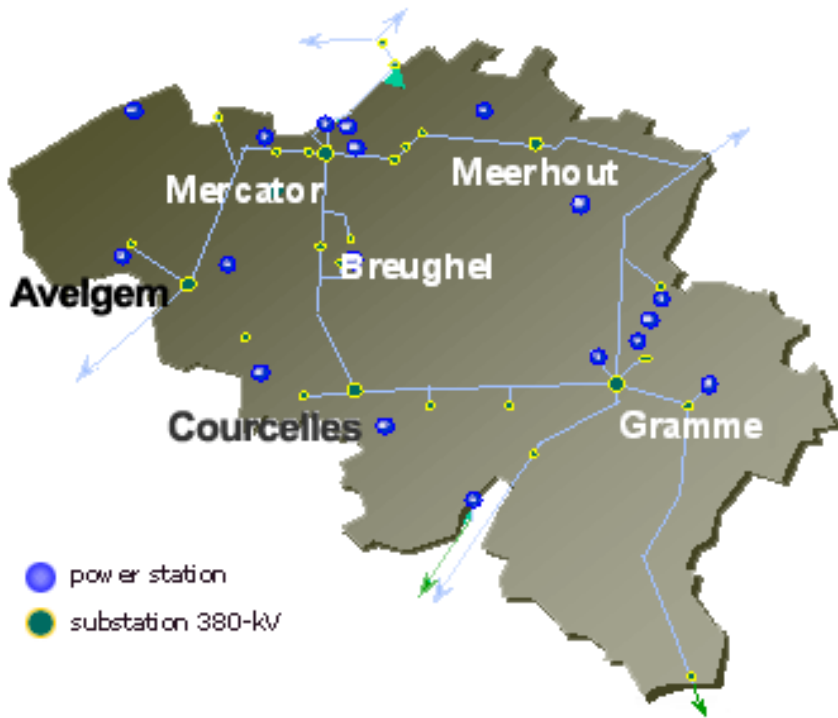


# Electric load forecasting

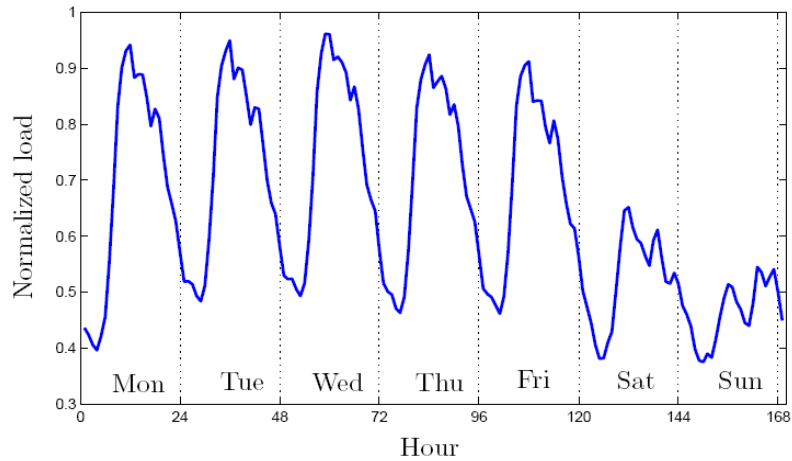


Model update: Every week!

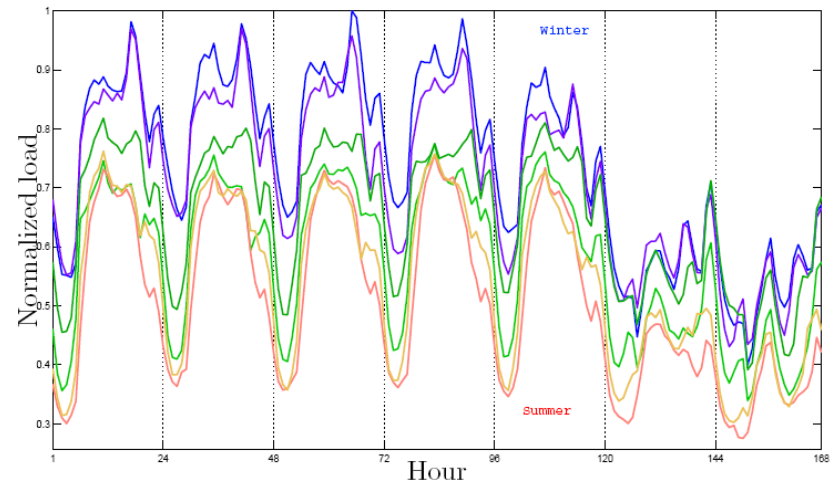
⇒ Accurate forecast



**250 transformer substations**  
**Every 15 min, 5 years**

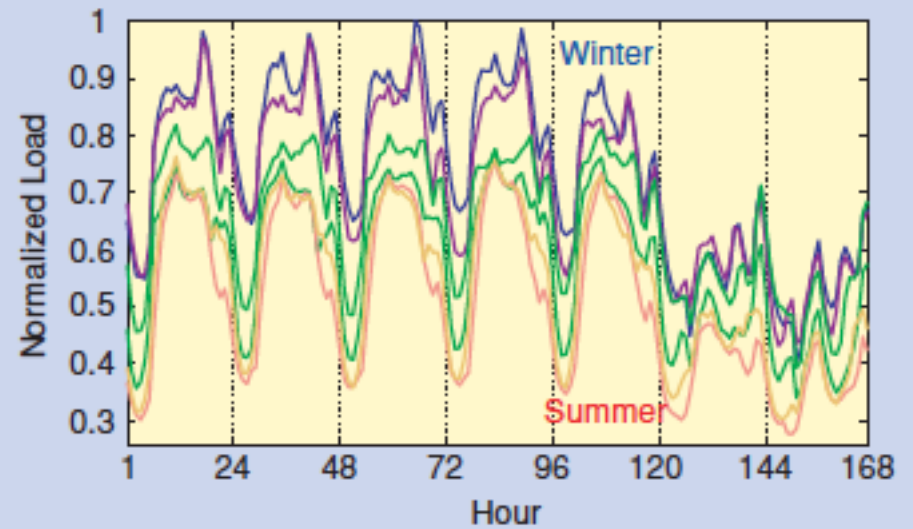
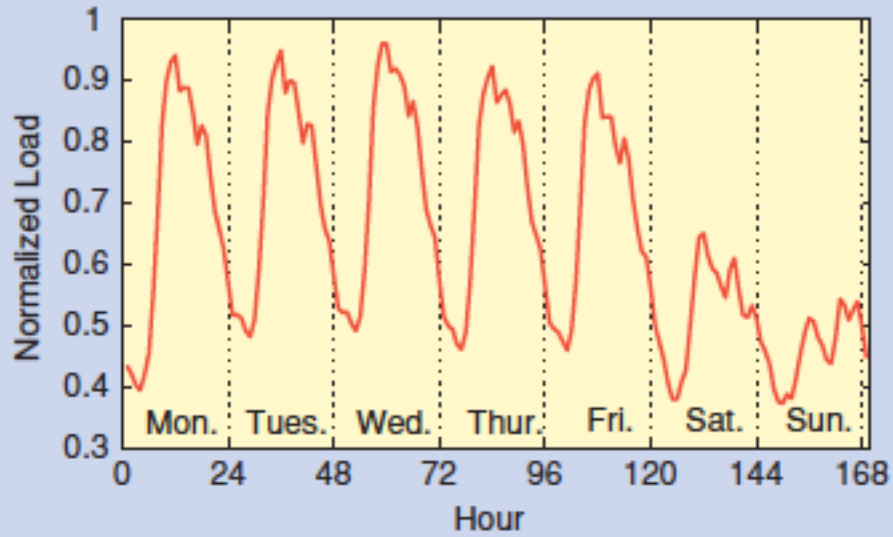


**48**  
**1 post, 1 week**

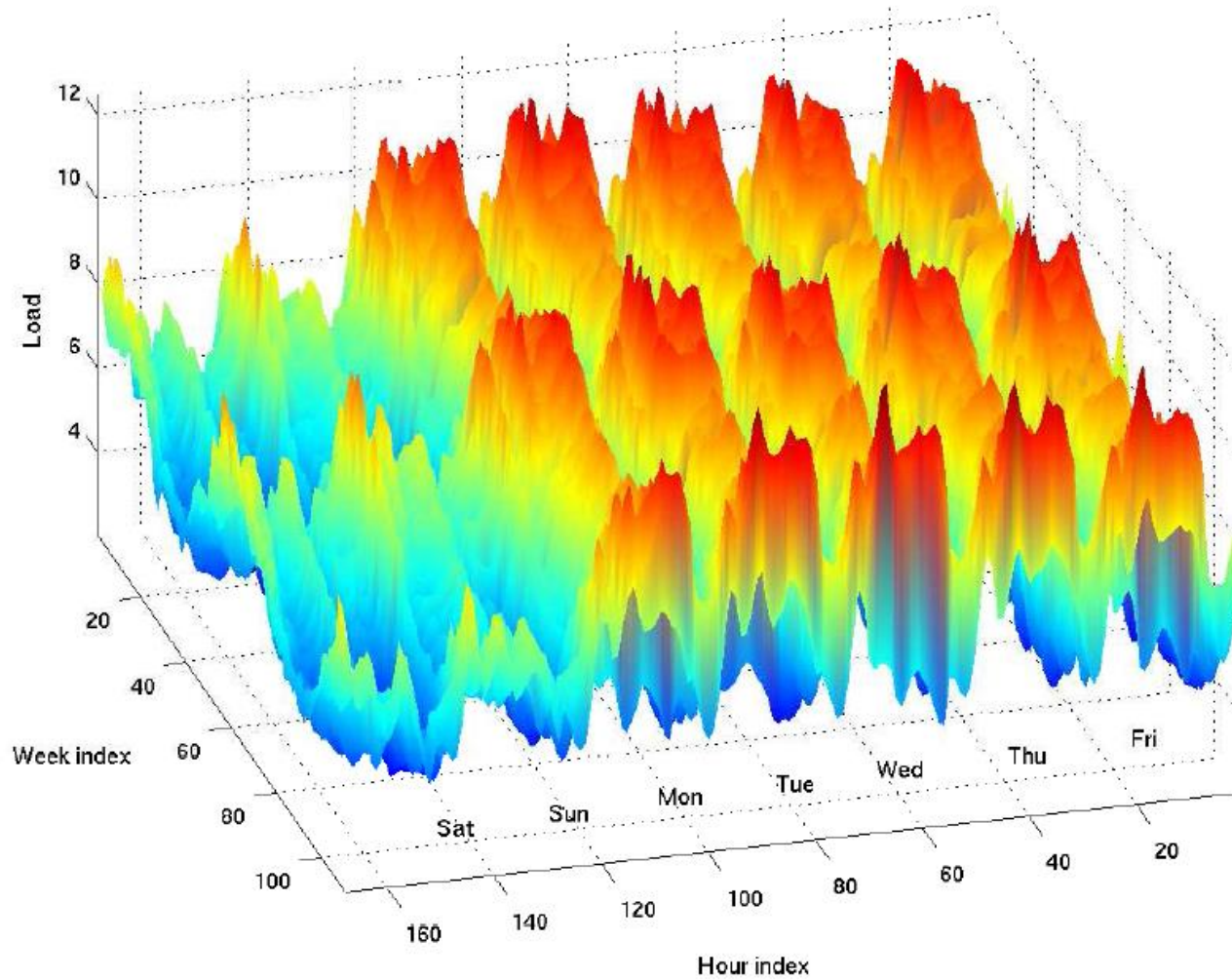


**1 post, four seasons**

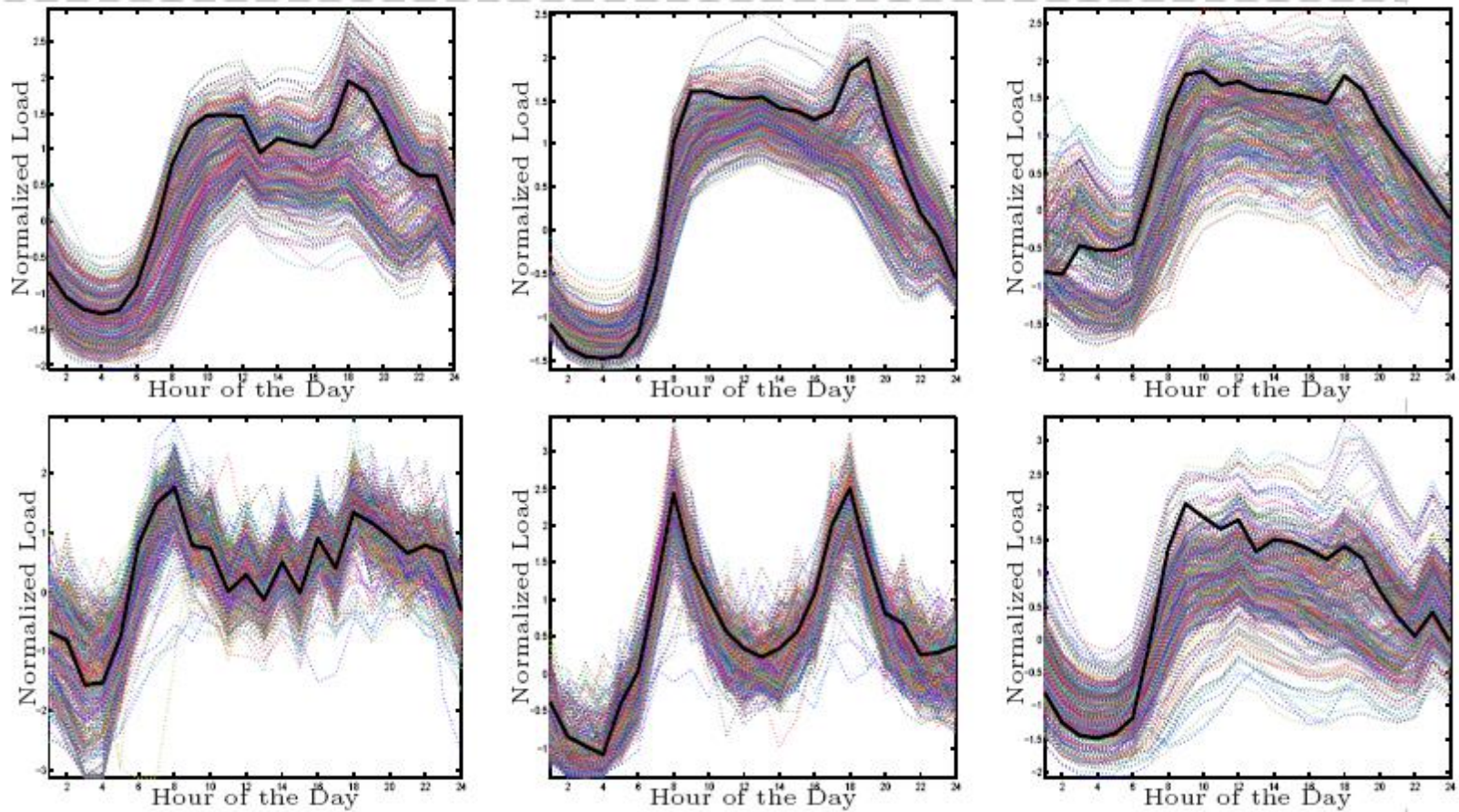
# Electric load forecasting



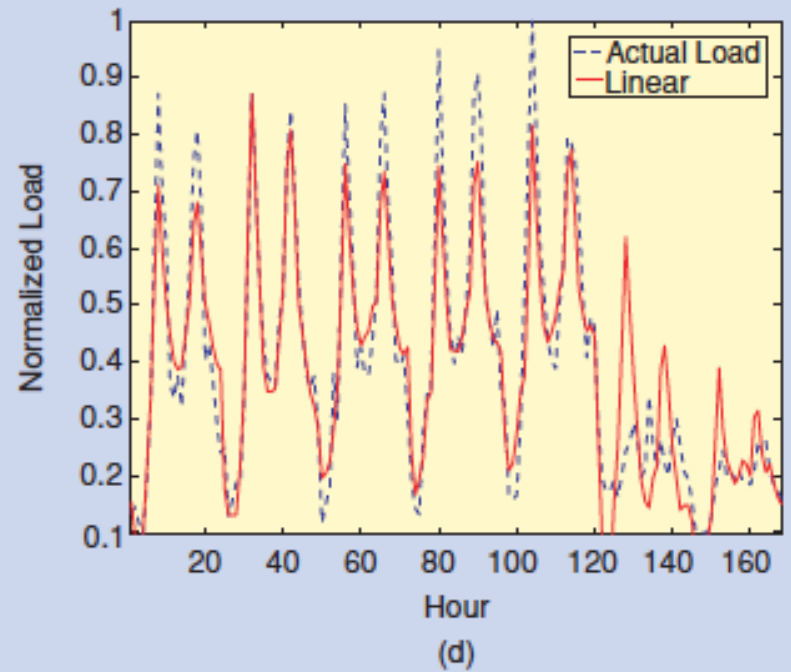
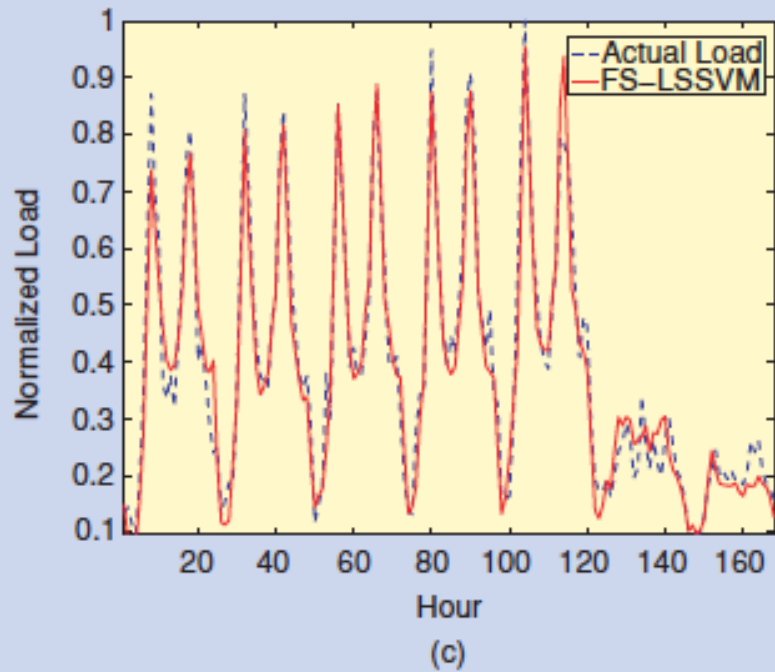
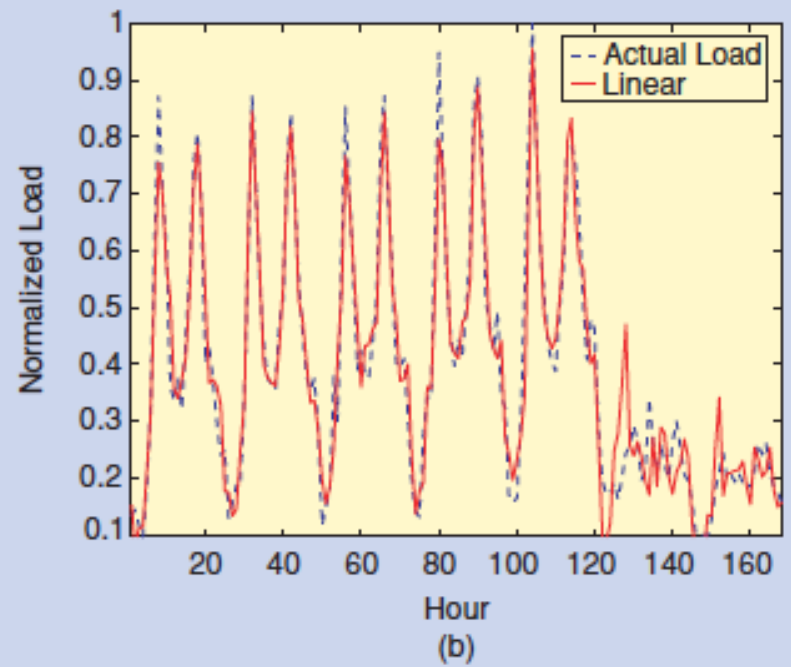
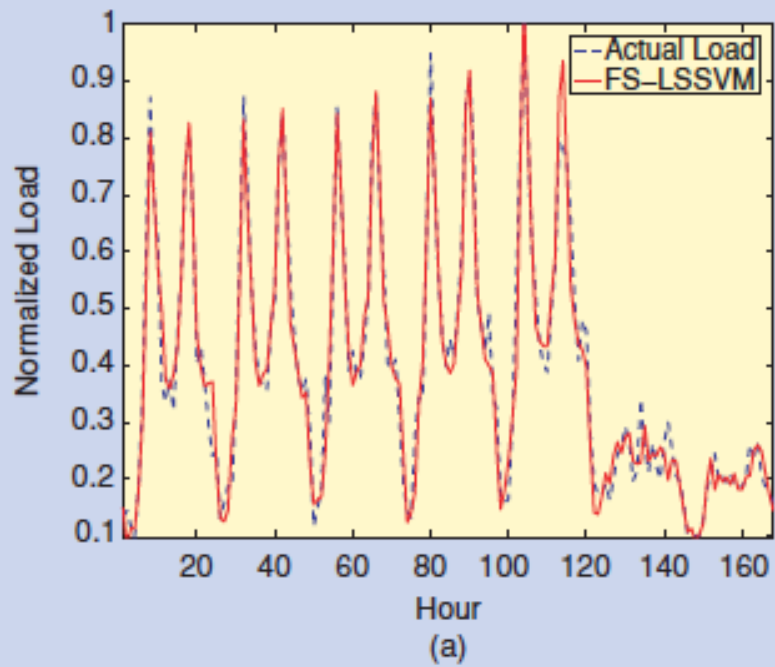




**Seasonalities in the load: day, week, year, holidays**

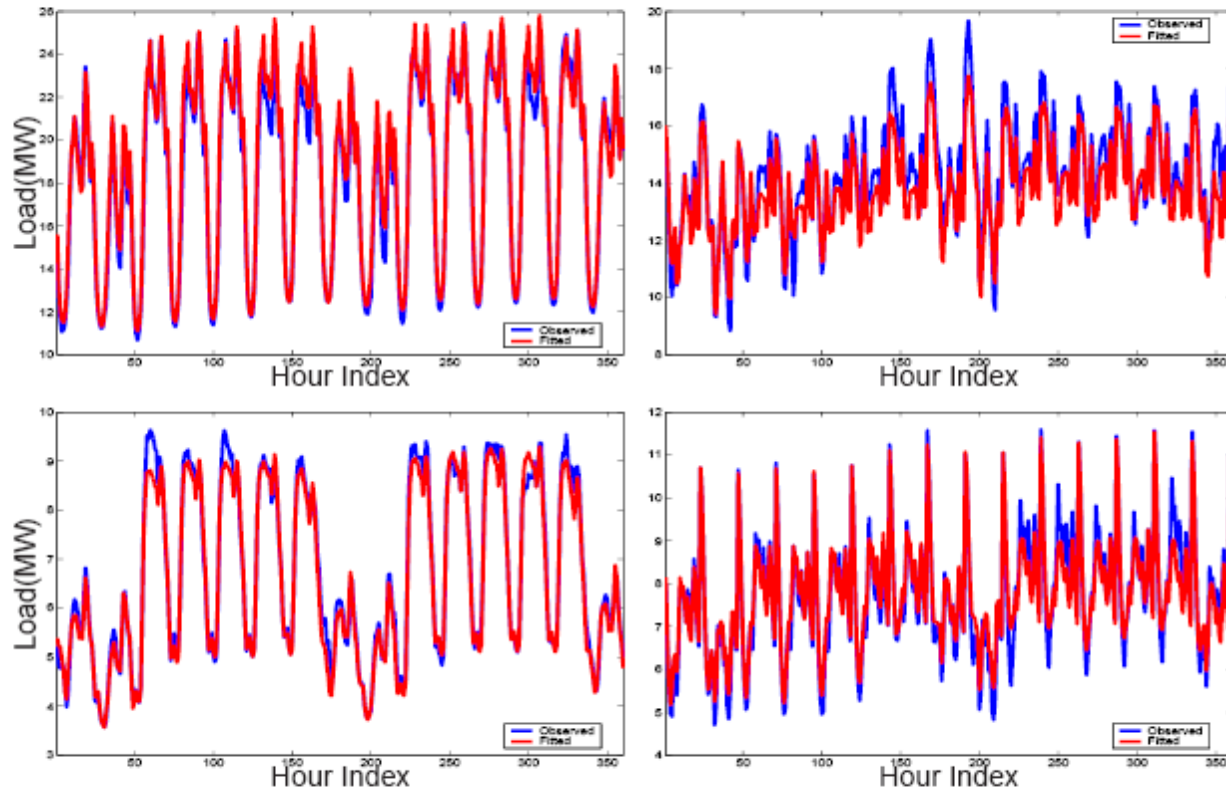


**6 posts, 1 year**  
**Seasonalities, calendar holidays !**



- **Seasonalities = a priori information (regress Monday on Monday !)**
- **Normalization:**
  - **remove effects of temperature, cloudiness,....**
  - **remove effect of holidays – calendar days (use dummies)**
- **Calculate 'eigenprofiles' = daily shape per post**

■ 15-days ahead forecasts for 4 posts:





A woman with short brown hair and glasses, wearing a grey top, is pointing at a tablet held by a healthcare worker in blue scrubs. The healthcare worker has her hair in a bun and is wearing a stethoscope. The background is a plain, light-colored wall.

**Energy**

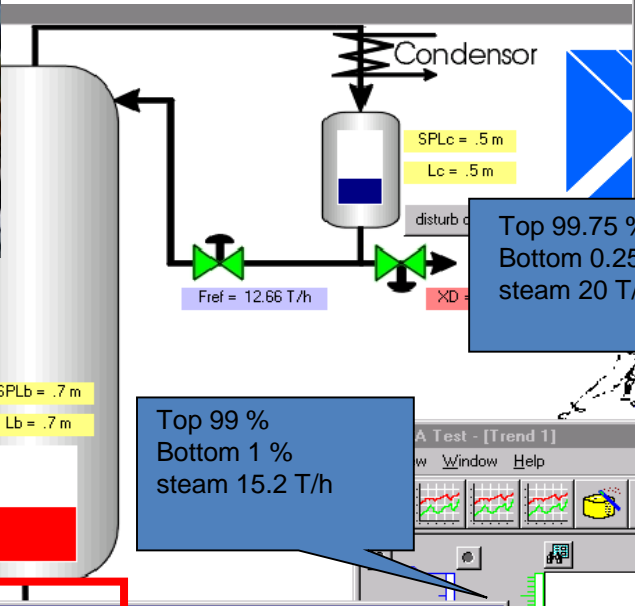
**Industry**

**Environment**

**Social networks**

**Fraud and predictive analysis**

**Health**



Top 99.8 % (99.8%)  
 Bottom 0.31 % (0.2 %)  
 steam 20 T/h

Top 99.75 % (99.8%)  
 Bottom 0.25 % (0.2 %)  
 steam 20 T/h

Top 99 %  
 Bottom 1 %  
 steam 15.2 T/h

4 OK FF = 4 T/h  
 50 OK XF = 50 %

SPLb = .7 m  
 Lb = .7 m

Fst = 20 T/h

Reboiler

INCAView - [Overview]

Log Status: Idle Controller Status, Reason: Turned on by operator request

CV NAME	ENGL0W	OPERLOW	IDEAL	IDEALRANK	OPERUPP	ENGUPP
linX_destil	-5.00	-2.61	-1.61	2	0.68	5.00
linX_bottom	-5.00	-2.00	-1.61	3	0.68	5.00

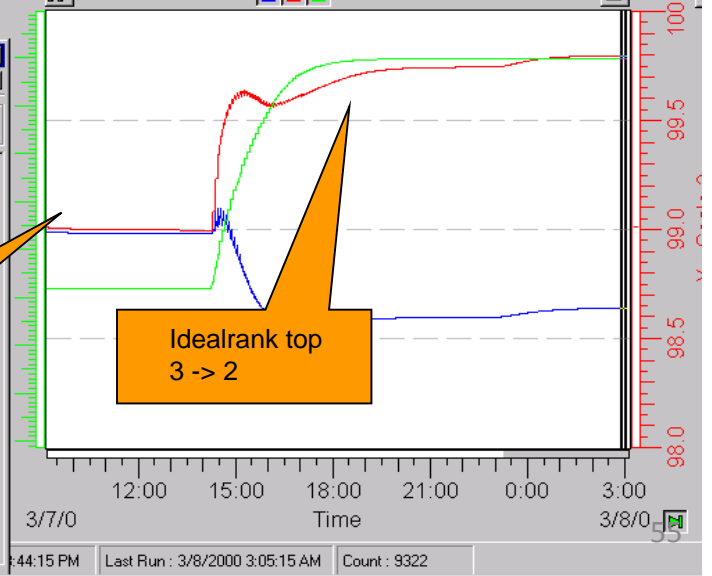
  

MV NAME	ENGL0W	OPERLOW	IDEAL	IDEALRANK
F_reflux	6.50	6.80	9.20	4
F_steam	1.00	2.50	3.00	4

DV NAME	DESCRIPTION	UNIT	PV	USE	CRIT	AUTO	BAD	LBND	UBND
Feedflow	Feed Flow	t/h	4.00	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>

Setpoint changes



Idealrank top  
 3 -> 2

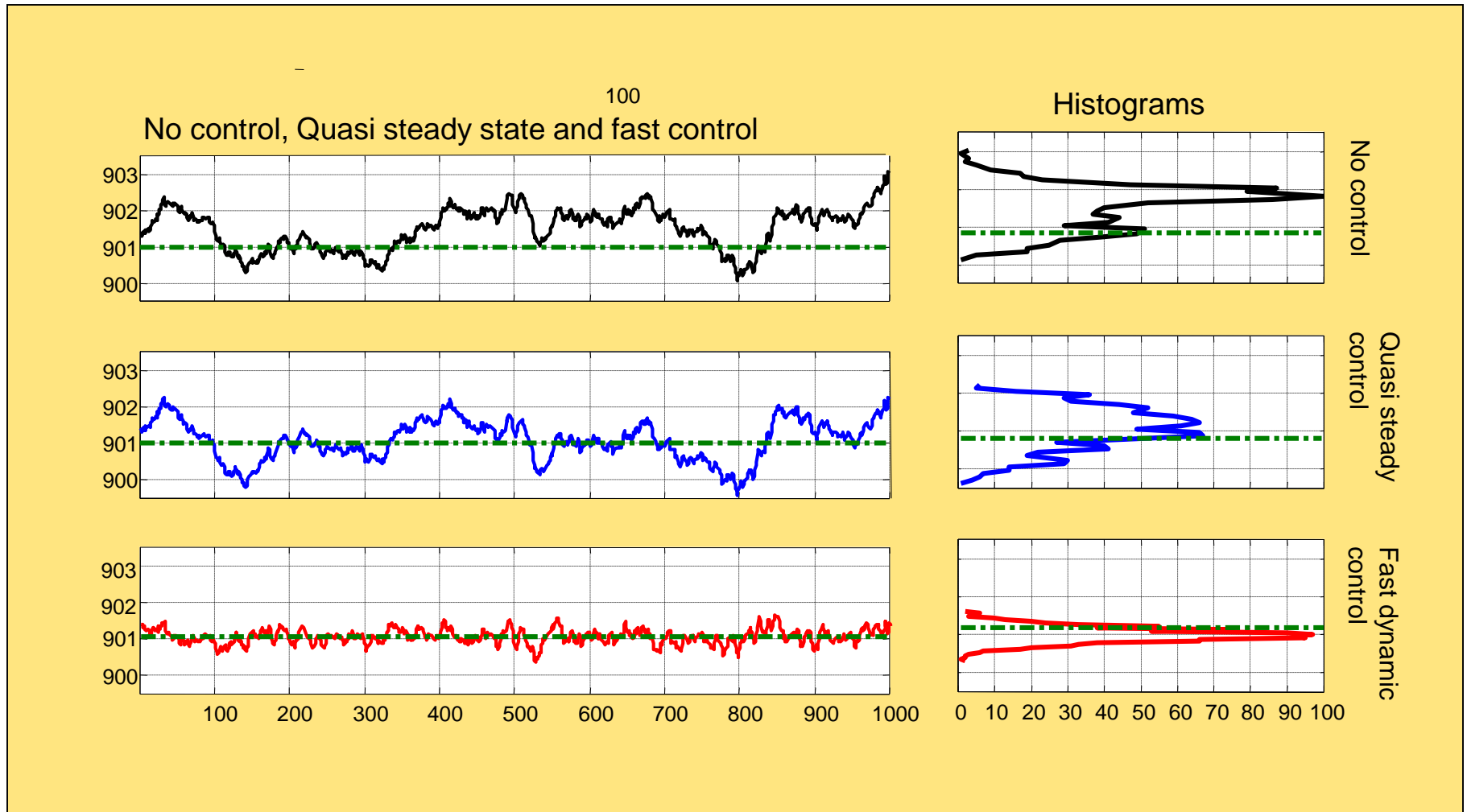
INCAView  
 162 Engine Release Version 1.30, Compiled on Jun 28, 2000, 12:19:42  
 Http://www.ppsol.com  
 © Copyright PPSOL Technology, Inc.

```
to lower opera
: linX_destil
to lower opera
to lower opera
: linX_destil
to lower opera
to lower opera
: linX_destil
et value is clipped to lower opera
et value is clipped to lower opera
```

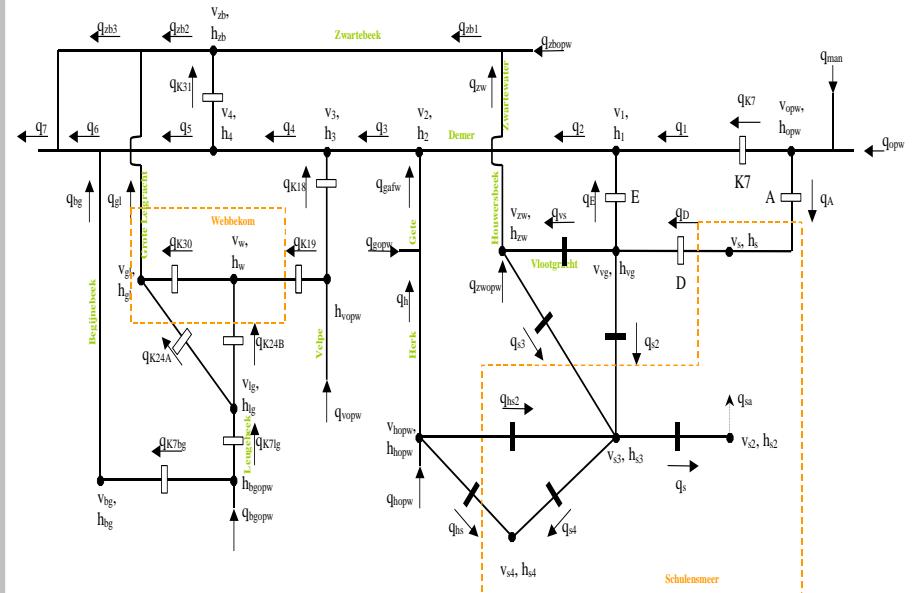
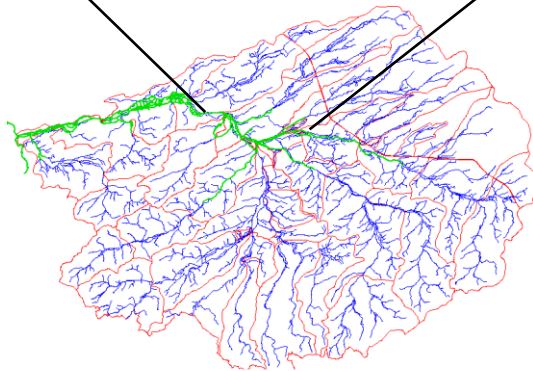
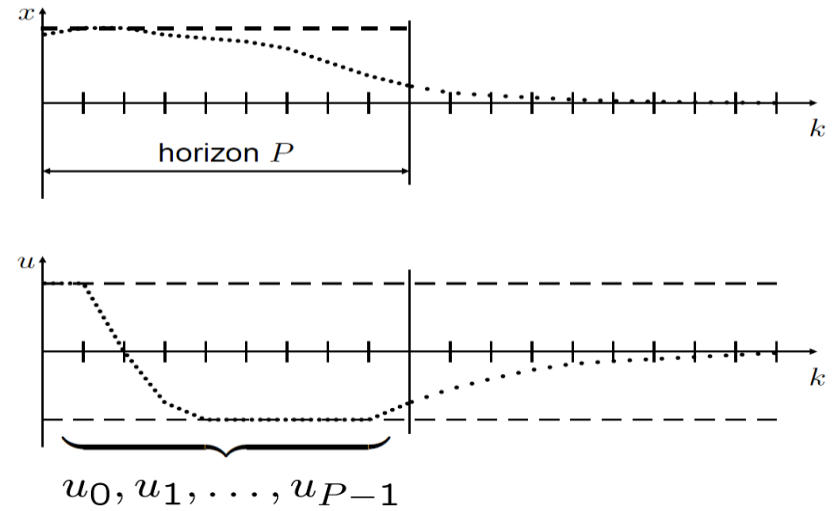
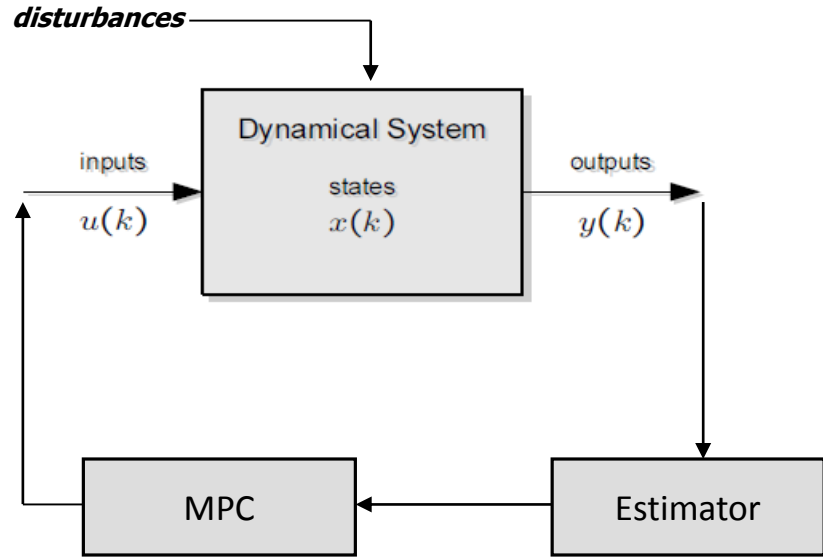
[Time] : 3/8/0 2:58  
 X\_bottom.pv: 0.32  
 X\_destil.pv: 99.800  
 F\_steam.sp: 20.00



# Modelling for control

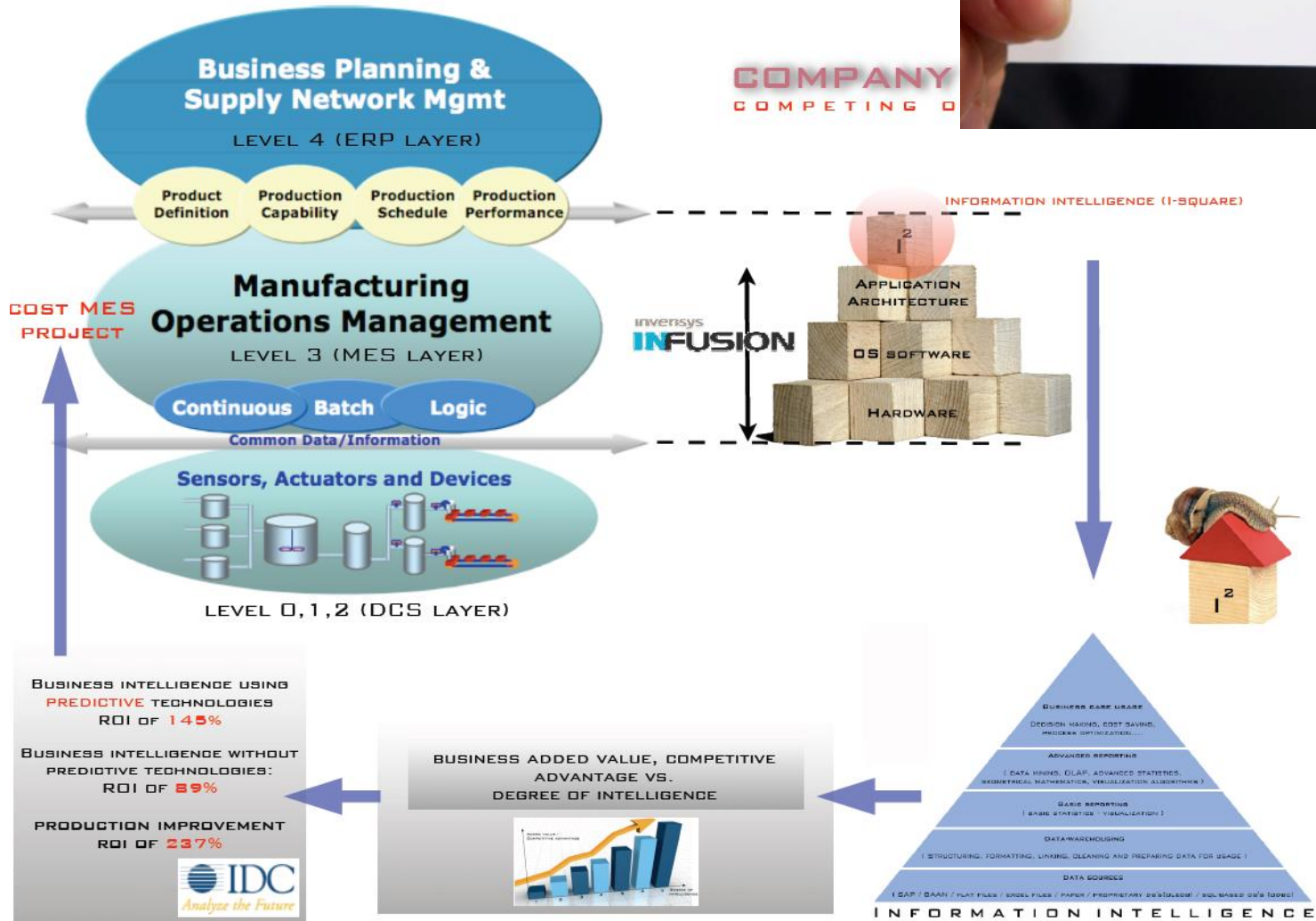


# Model Based Predictive Control for Flood Regulation: Demer





COMPANY  
COMPETING O



# Dsquare website



A photograph showing two women. On the left, an older woman with short brown hair and glasses, wearing a grey turtleneck, is pointing at a tablet. On the right, a younger woman with her hair in a bun, wearing blue scrubs and a stethoscope, is holding the tablet. The background is a plain, light-colored wall.

**Energy**

**Industry**

**Environment**

**Social networks**

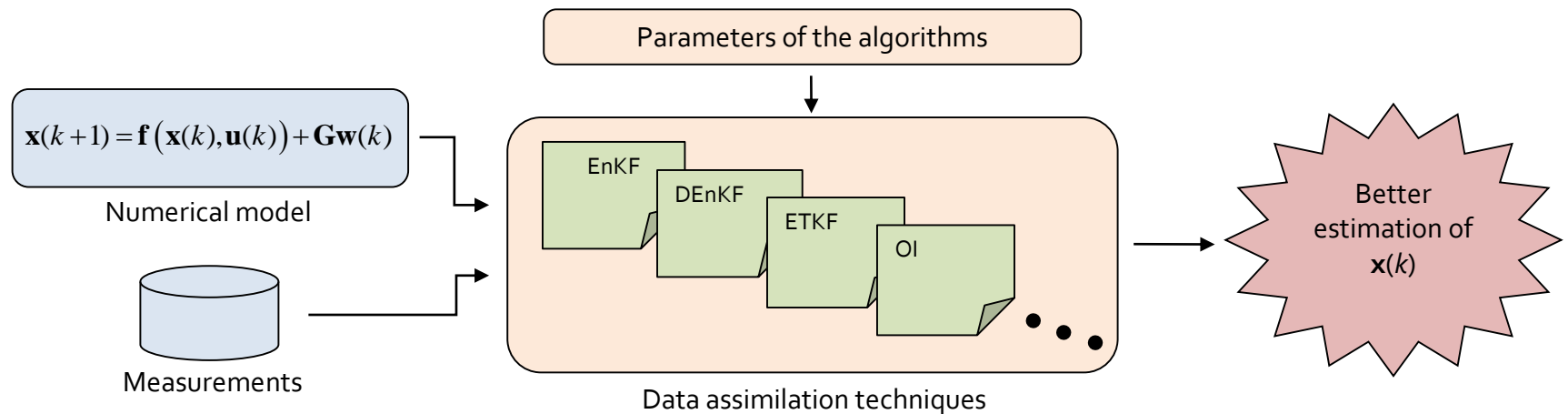
**Fraud and predictive analysis**

**Health**



# Data Assimilation

Data assimilation is the common name given to several numerical techniques that combine **the outputs of a numerical model** with **observational data** in order to improve the quality of the model predictions.

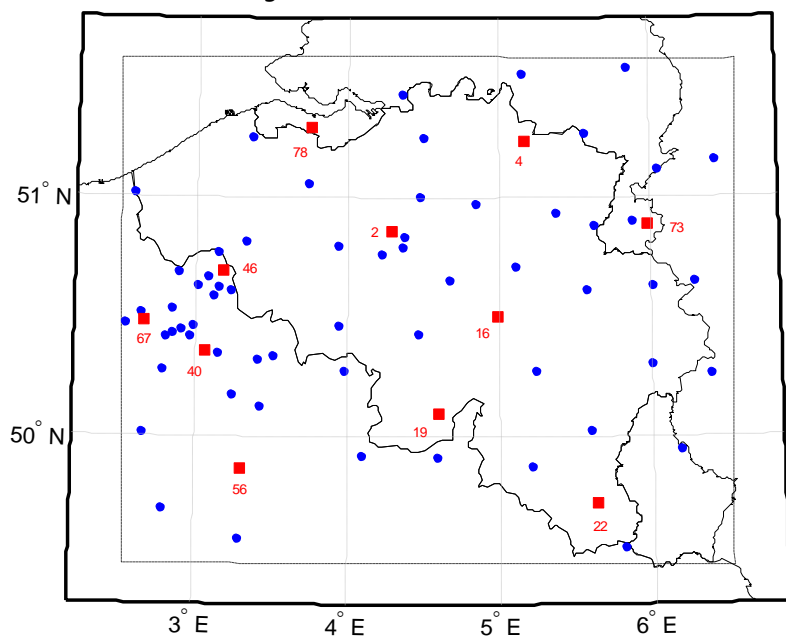


Some data assimilation techniques: 3DVAR, 4DVAR, Ensemble Kalman Filter (EnKF) and its variants, Optimal Interpolation (OI), particle filters, etc.

# Data Assimilation

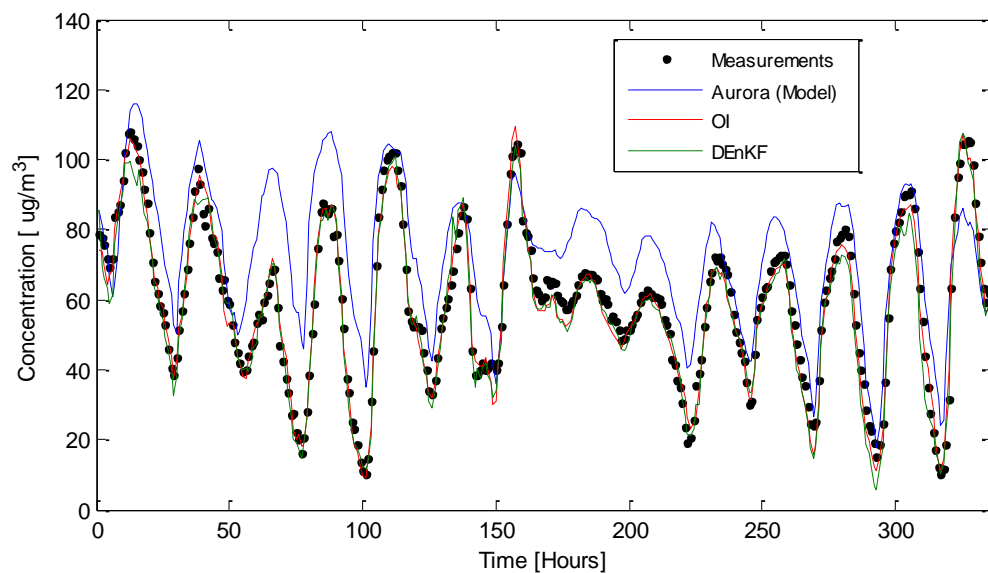
The Deterministic Ensemble Kalman Filter (DEnKF) and the OI technique have been used to improve the O<sub>3</sub> estimates of the Air-quality model AURORA.

O<sub>3</sub> air-quality stations



- - Assimilation stations
- - Validation stations

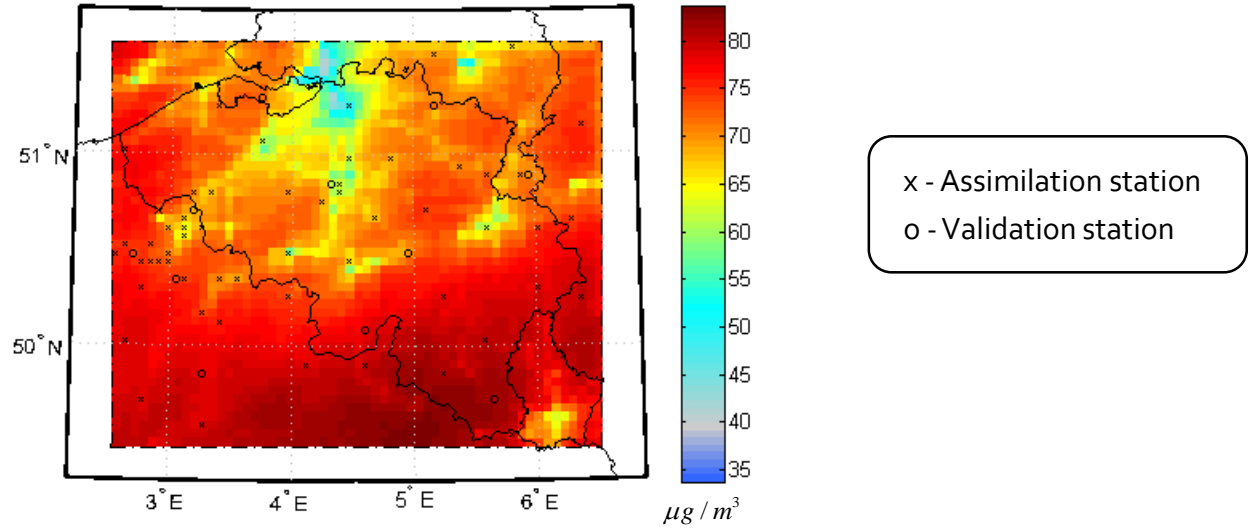
Average of the O<sub>3</sub> concentration over the validation stations



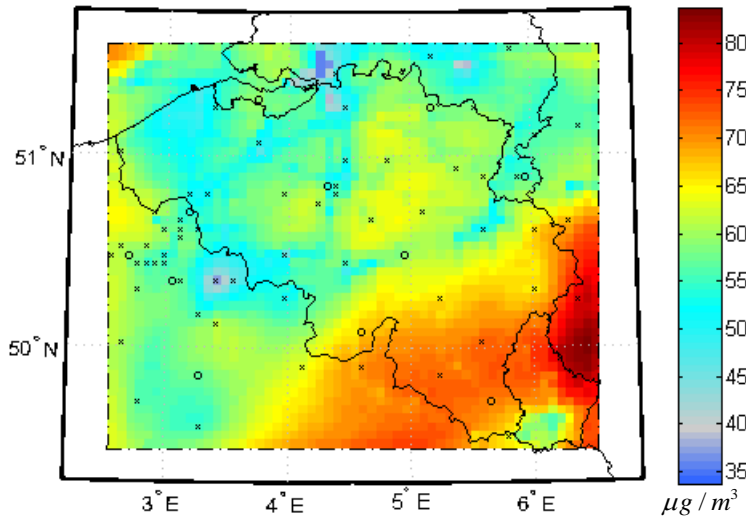
Starting date: May 28<sup>th</sup>, 2005 at midnight

# Average of the O<sub>3</sub> concentration field over the 14 day period

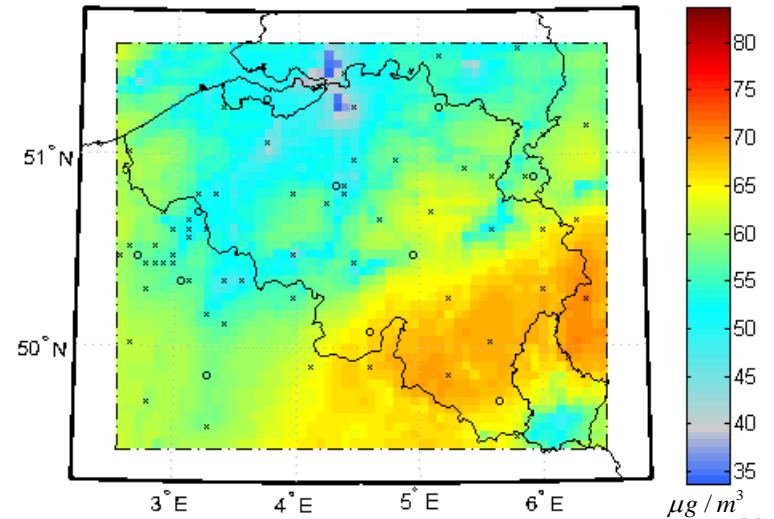
## Free-run of Aurora



## Optimal Interpolation



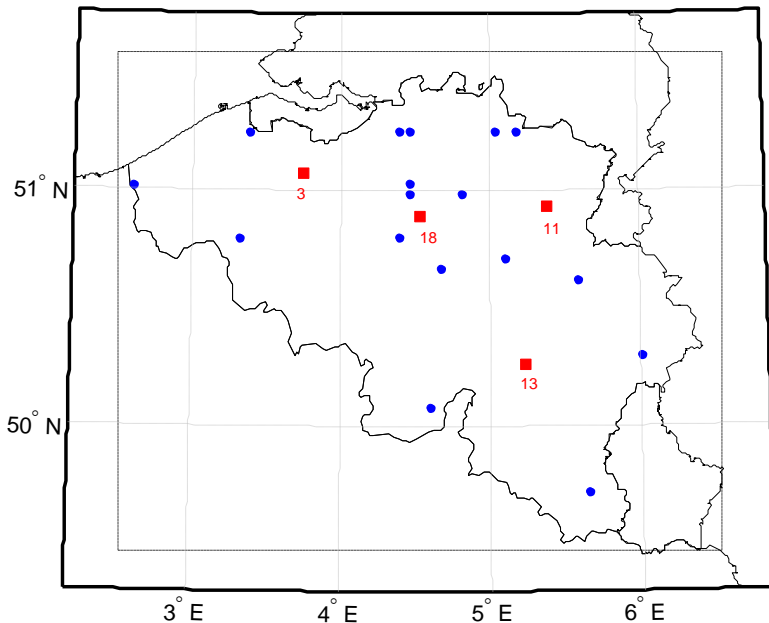
## DEnKF



# Data Assimilation

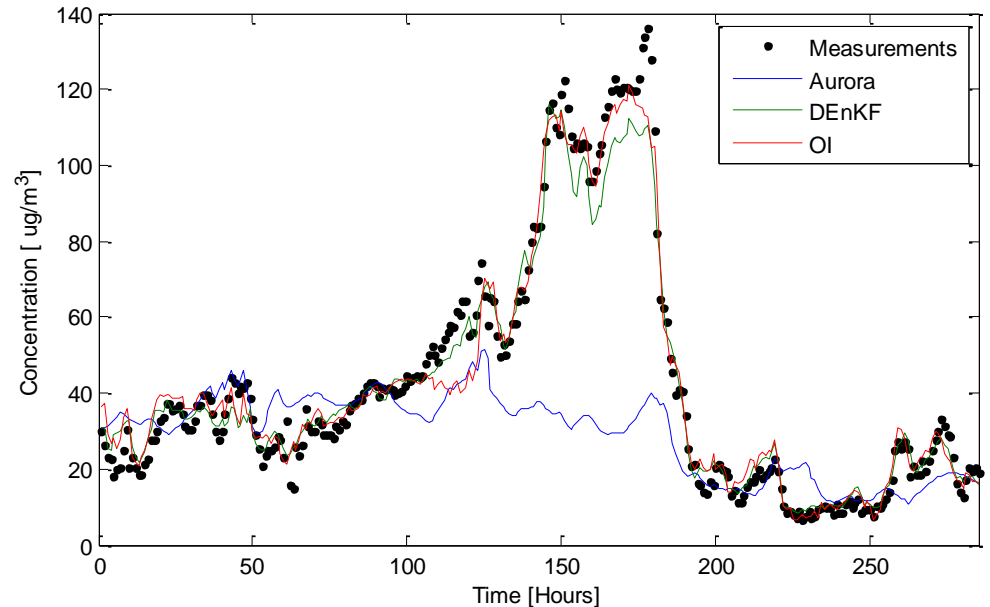
The Deterministic Ensemble Kalman Filter (DEnKF) and the OI technique have been used to improve the PM<sub>10</sub> estimates of the Air-quality model AURORA.

PM<sub>10</sub> air-quality stations



- - Assimilation stations
- - Validation stations

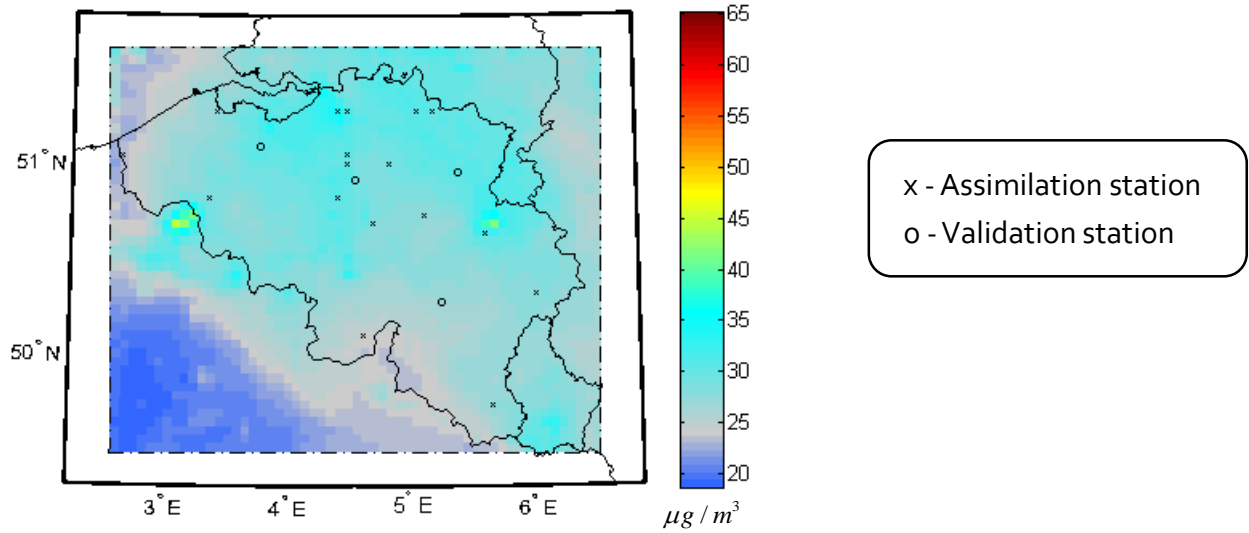
Average of the PM<sub>10</sub> concentration over the validation stations



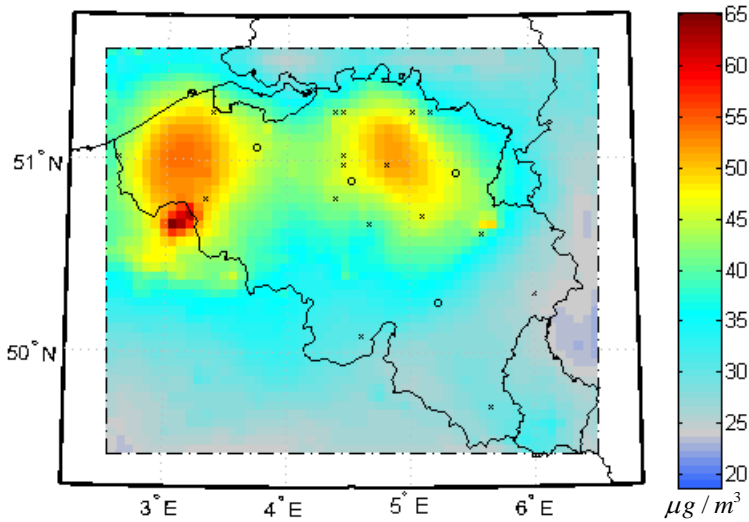
Starting date: January 20<sup>th</sup>, 2010 at midnight

# Average of the PM<sub>10</sub> concentration field

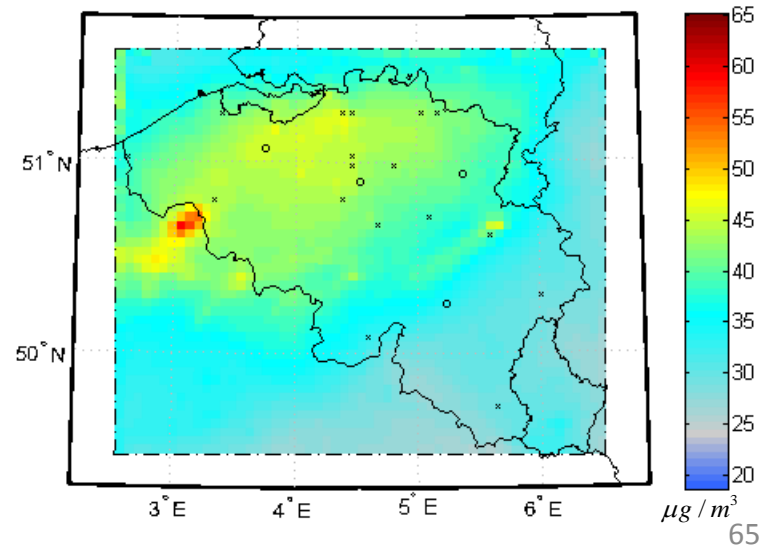
## Free-run of Aurora



## Optimal Interpolation

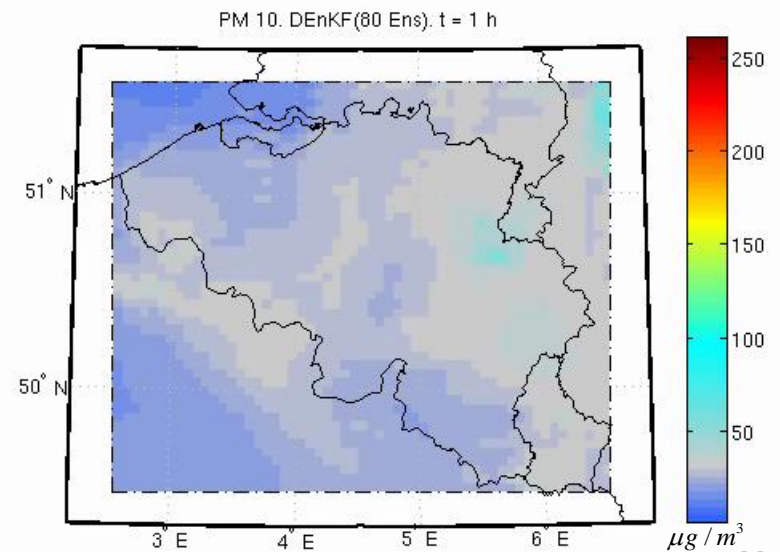
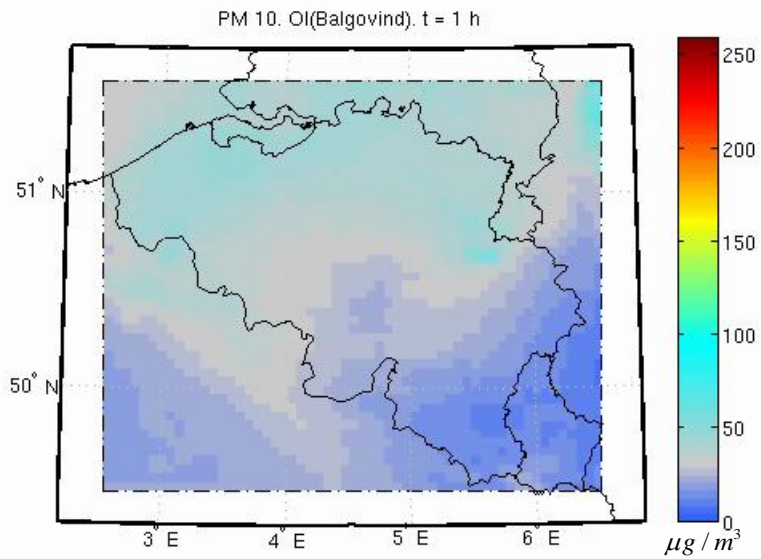
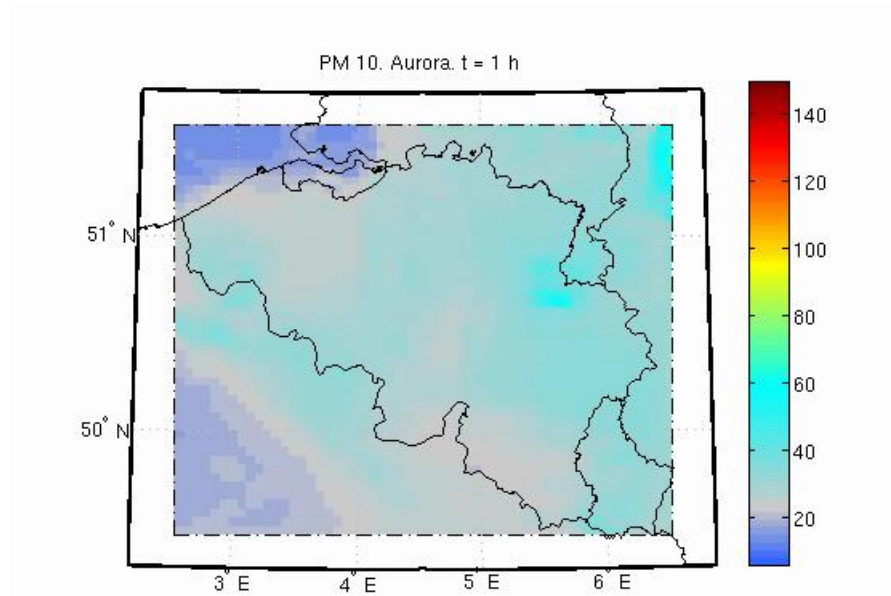


## DEnKF





# Average of the PM<sub>10</sub> concentration field



A photograph showing two women. On the left, an older woman with short brown hair and glasses, wearing a grey top, is pointing at a tablet. On the right, a younger woman with her hair in a bun, wearing blue scrubs and a stethoscope, is holding the tablet. The image is split vertically down the middle.

**Energy**

**Industry**

**Environment**

**Social networks**

**Fraud and predictive analysis**

**Health**

# Journal Clustering



Find all about specific topic?





# Journal Clustering

## Weighted Hybrid Clustering by Combining Text Mining and Bibliometrics on a Large-Scale Journal Database

We propose a new hybrid clustering framework to incorporate text mining with bibliometric journal set analysis. The framework integrates two different approaches: clustering (ensemble and kernel-based clustering). To improve the flexibility and the efficiency of processing large-scale data, we propose an information-based weighting scheme to leverage the effect of multiple data sources in hybrid clustering. Three different algorithms are extended by the proposed weighting scheme and they are employed on a large journal set retrieved from the Web of Science (WoS) database. The clustering performance of the proposed algorithms is systematically evaluated using multiple evaluation methods, and they were cross-compared with alternative methods. Experimental results demonstrate that the proposed weighted hybrid clustering strategy is superior to other methods in clustering performance and efficiency. The proposed approach also provides a more refined structural mapping of journal sets, which is useful for monitoring and detecting new trends in different scientific fields.

### Introduction

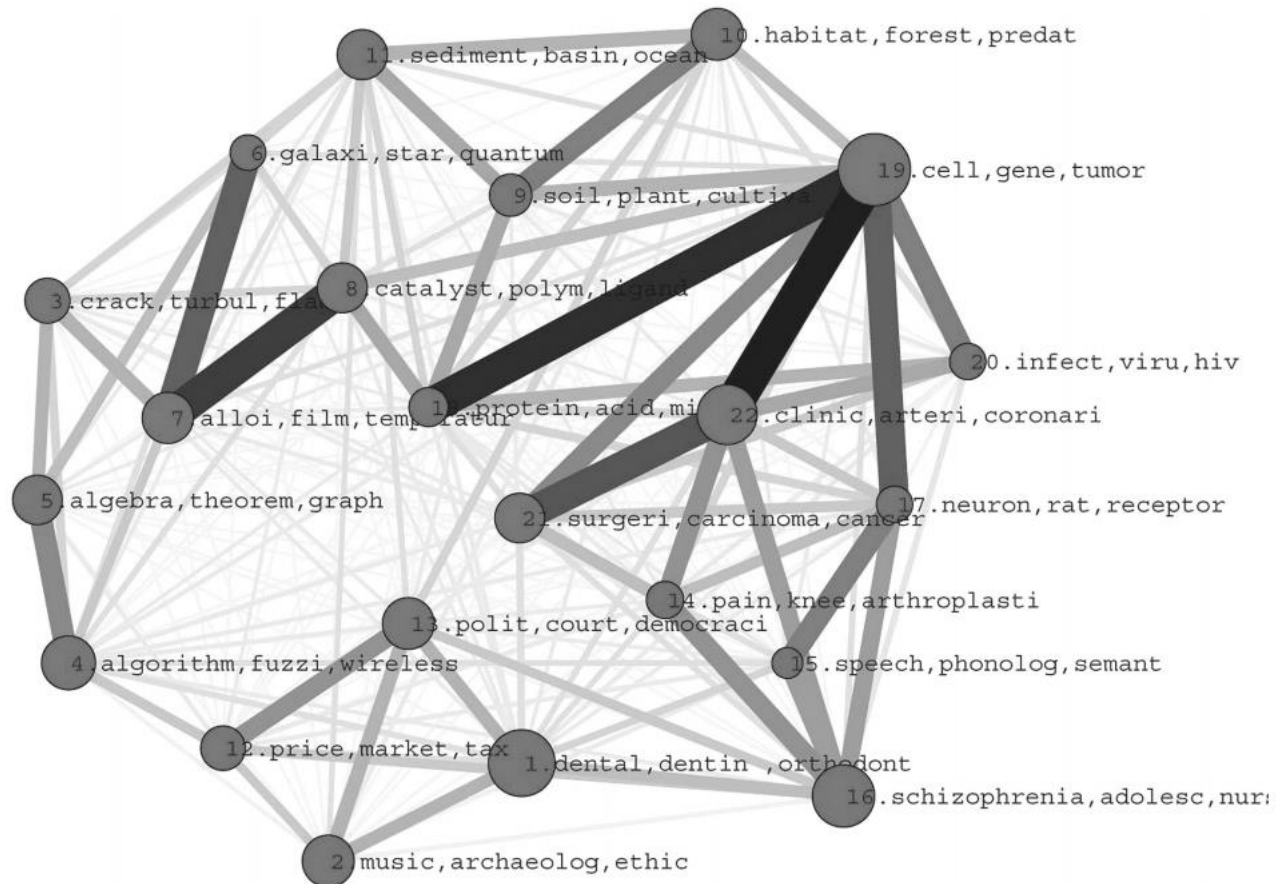
In scientometrics, information from journals can be categorized lexically or with citations. An important area of scientometric research is the clustering or mapping of scientific publications. The widely used method of cocitation clustering was introduced independently by Small (1973, 1978) and Marshakova (1973). Cross-citation-based cluster analysis for science mapping is different; while the former is usually based on links connecting individual documents, the latter requires aggregation of documents to units like journals or subject fields among which cross-citation links are established. Some advantages of this method (for instance, the possibility to analyze directed information flows) are undermined by possible biases. For example, bias could be caused by the use of predefined units (journals, subject categories, etc.), implying already certain structural classification. Journal cross-citation clustering has been used by Leydesdorff (2006), Leydesdorff and Rafols (2009), and Boyack, Börner, and Klavans (2009), while Moya-Aneón et al. (2007) applied subject cocitation analysis to visualize the structure of science and its dynamics.

The integration of lexical similarities and citation links has also attracted interest in other fields such as search engine

Received July 7, 2009; revised October 31, 2009; accepted December 30, 2009

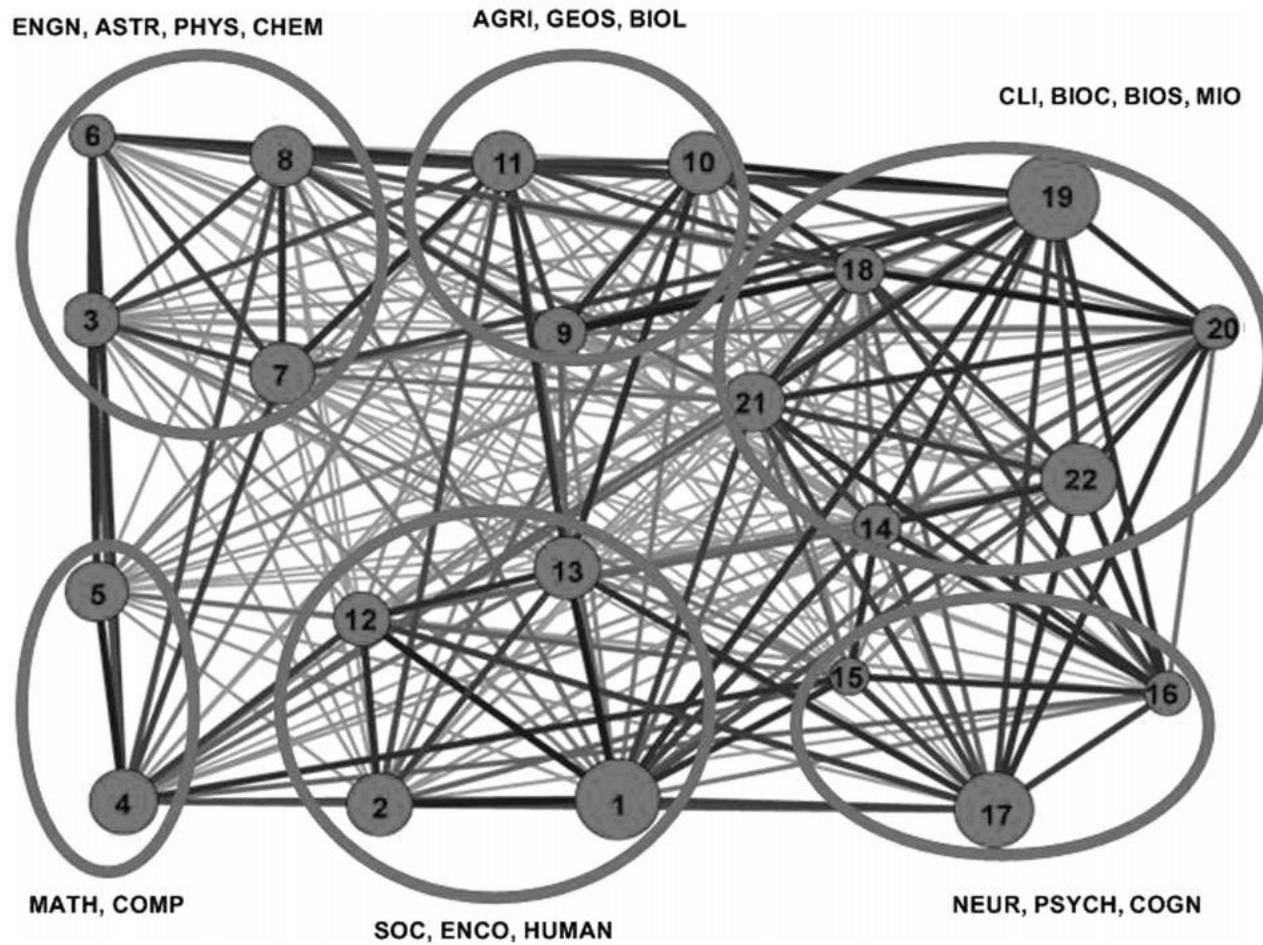
© 2010 ASIST & Published online 11 March 2010 in Wiley InterScience

# Journal Clustering





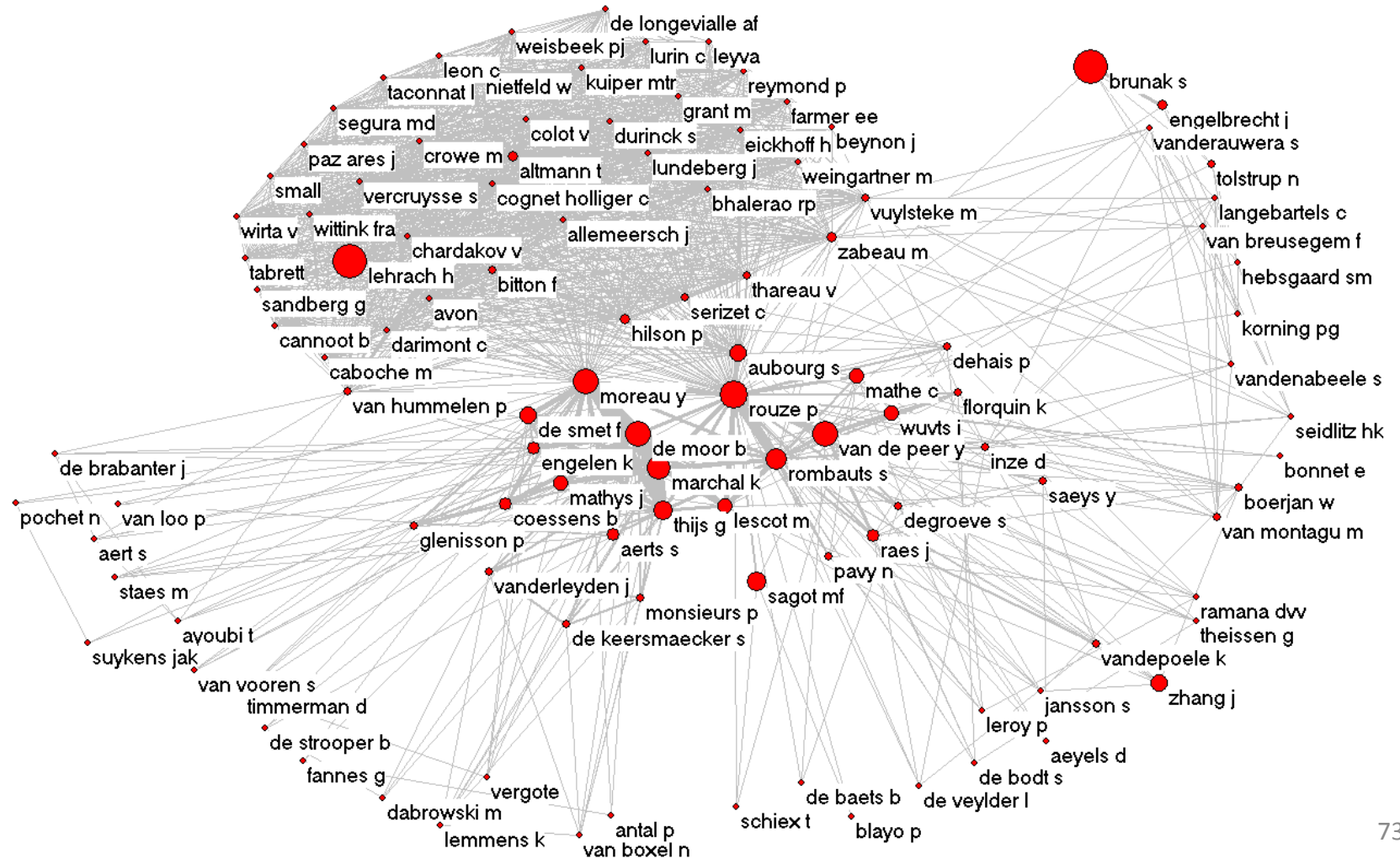
# Journal Clustering



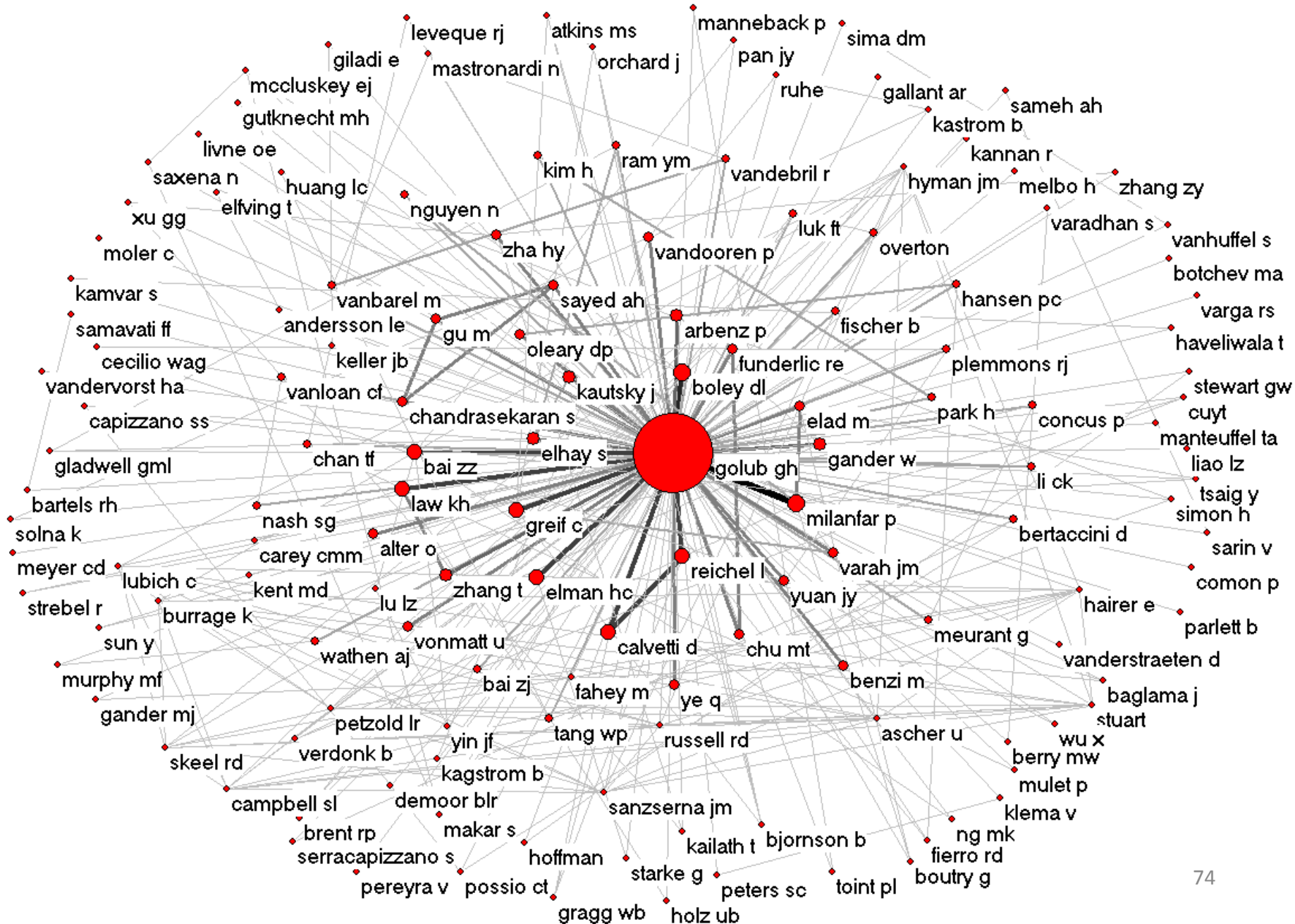
# Author Collaboration Clustering



# Author Collaboration Clustering

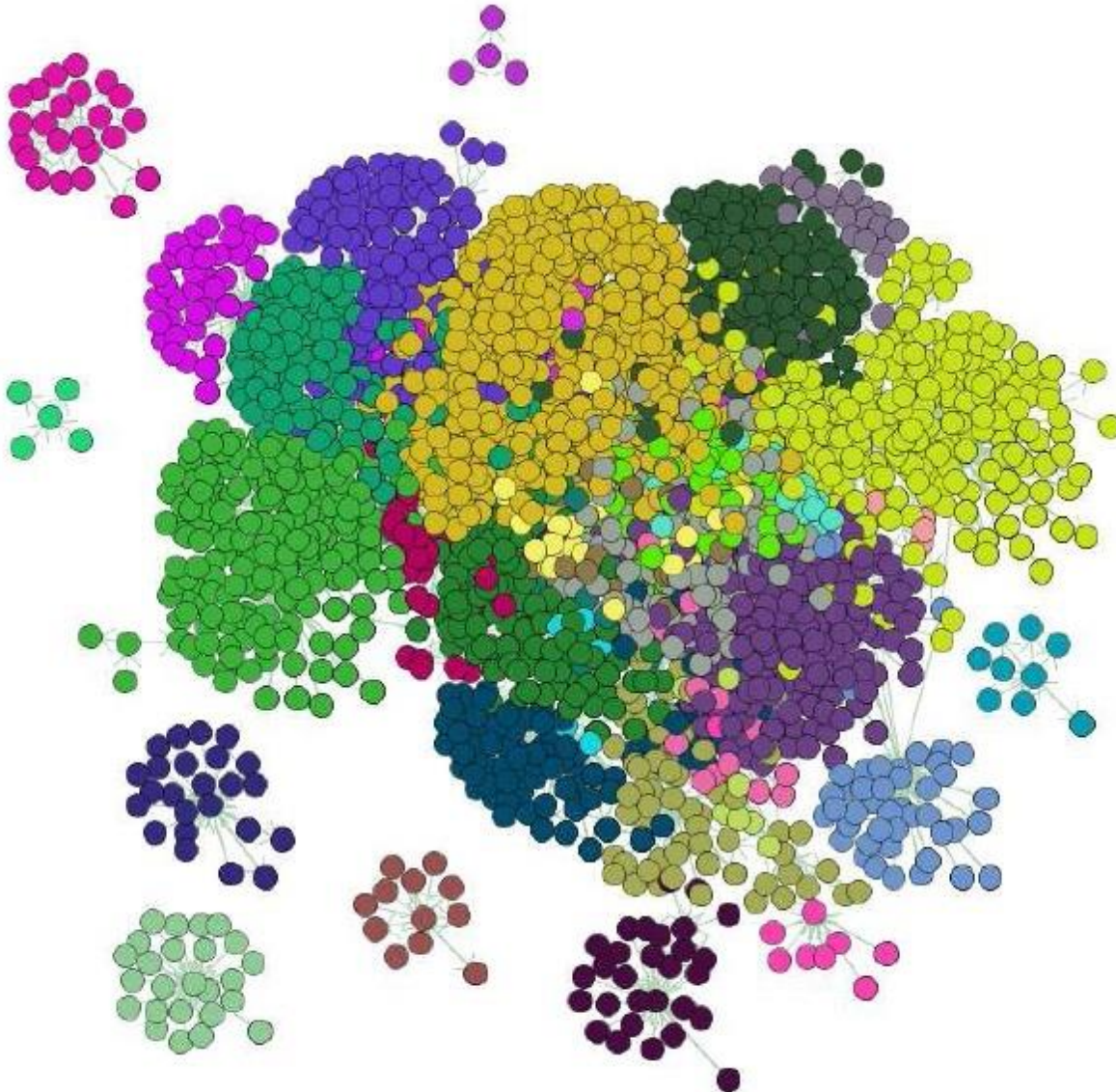


# Golub around the world commemoration February 29 2008





# Web of Science based literature network for Lennart Ljung



138 seed papers  
+ all cited and citing publications

Result: 4943 nodes, 6216 edges

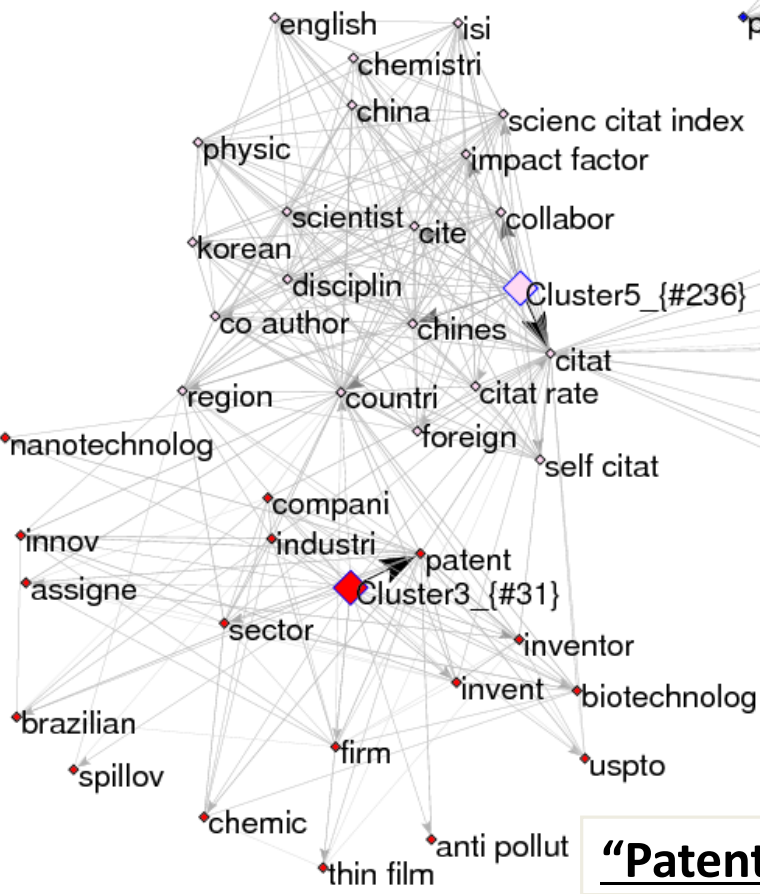
Link based clustering identifies  
topically homogeneous clusters.

13 papers are written  
by another L. Ljung.

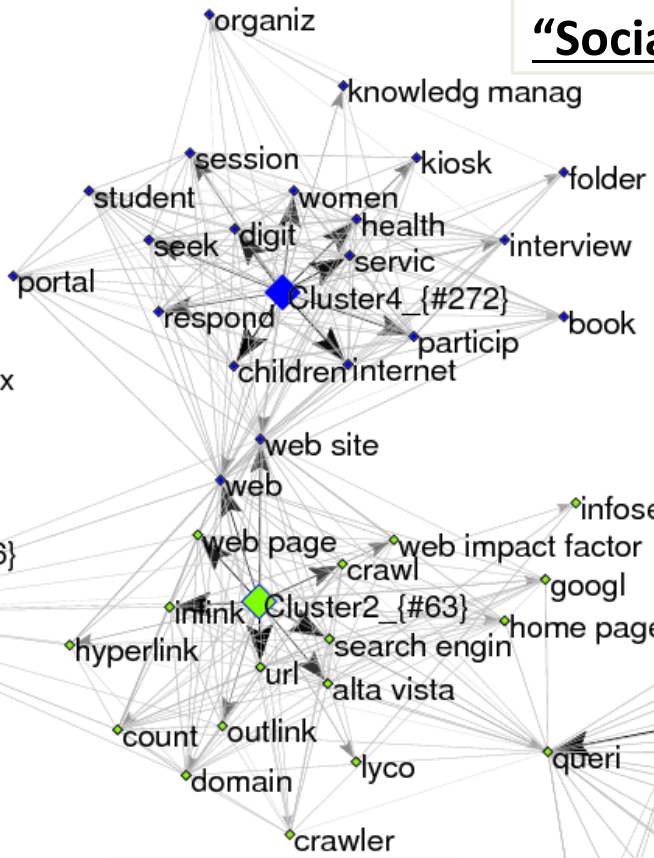


# Community detection

**“Bibliometrics”**



**“Social aspects”**



**“Webometrics”**

**“Information retrieval”**

**“Patent analysis”**

A woman with short brown hair and glasses, wearing a grey top, is pointing at a tablet held by a healthcare worker in blue scrubs. The healthcare worker has her hair in a bun and is wearing a stethoscope. The background is a plain, light-colored wall.

**Energy**

**Industry**

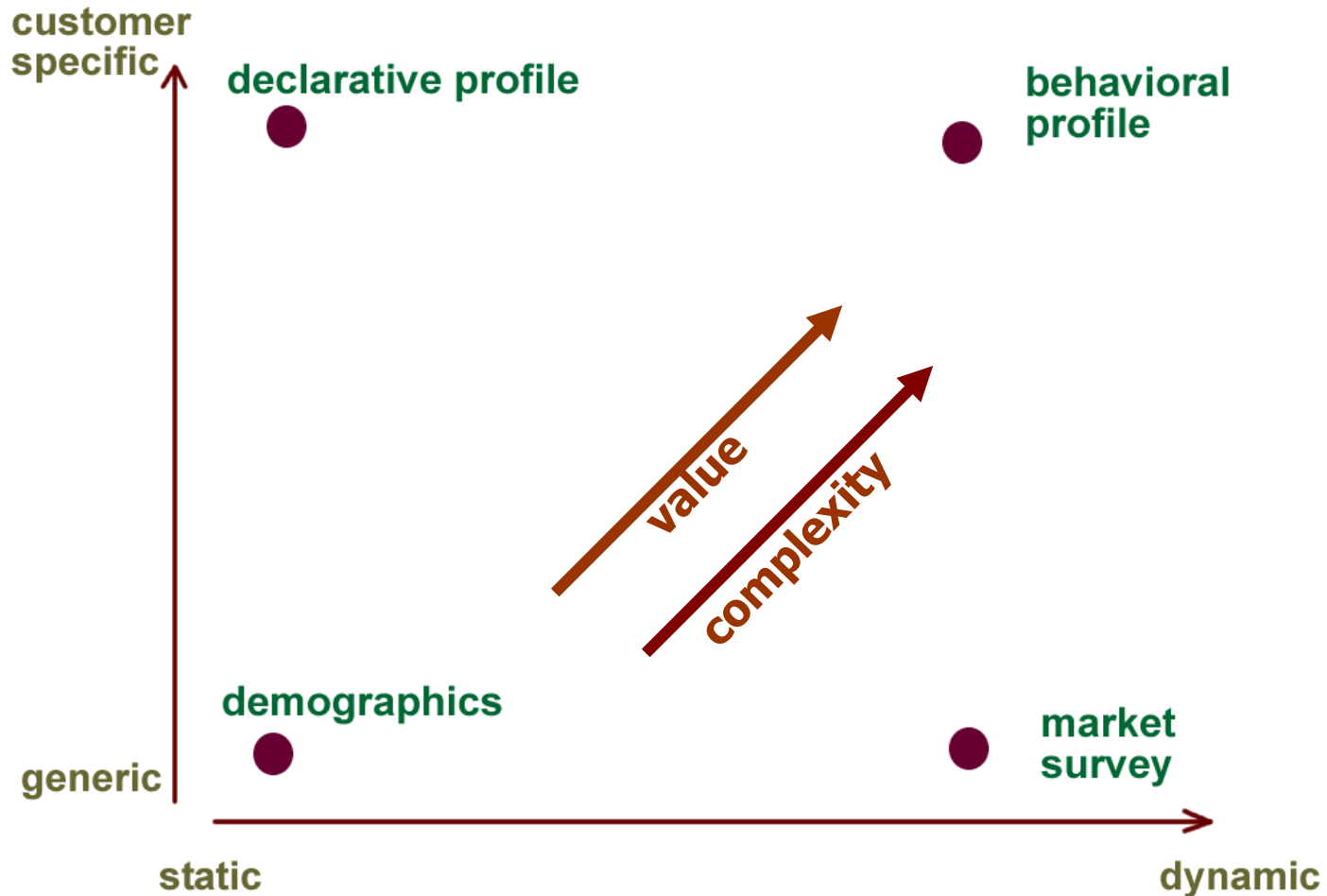
**Environment**

**Social networks**

**Fraud and predictive analysis**

**Health**

# Customer Intelligence

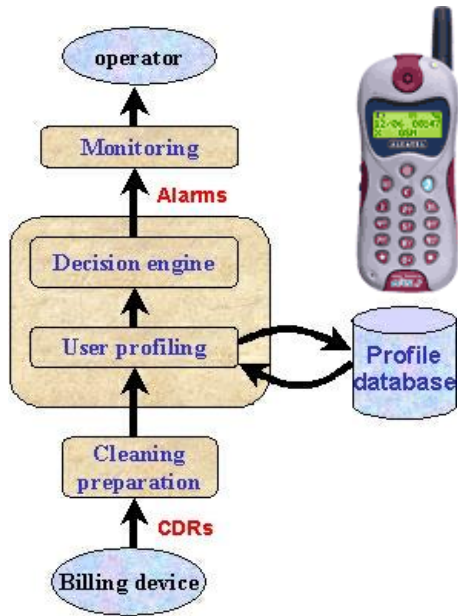




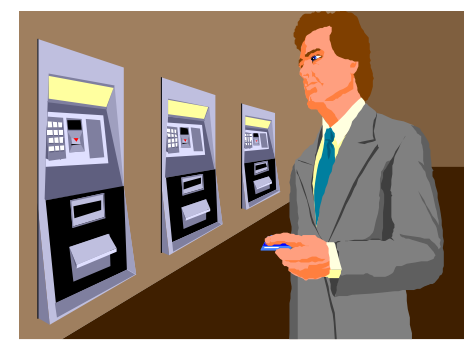
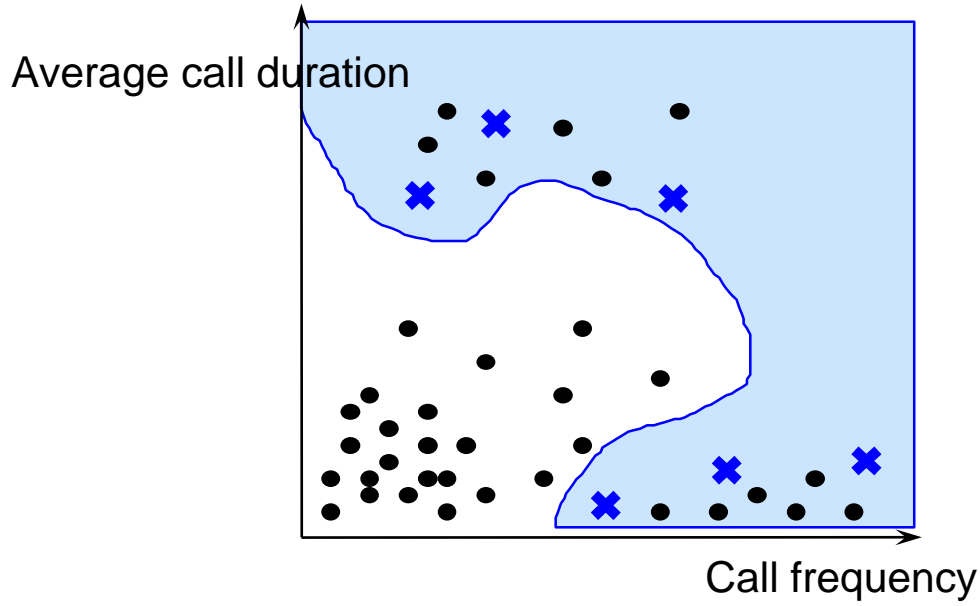
# Customer Intelligence



# Fraud detection on mobile phone network



	Short Duration	Long Duration	High Frequency	International	Same Destination	Off Peak	Call Forwarding	Behaviour Change
Direct calling		X	X	X			X	
ABX fraud	X		X		X	X		X
Freephone fraud	X		X		X			X
Premium rate fraud		X	X		X			X
Subscription fraud			X					
Handset theft		X	X	X	X			X





A photograph showing two women. On the left, an older woman with short brown hair and glasses, wearing a grey top, is pointing at a tablet. On the right, a younger woman with her hair in a bun, wearing blue scrubs and a stethoscope, is holding the tablet. The background is a plain, light-colored wall.

**Energy**

**Industry**

**Environment**

**Social networks**

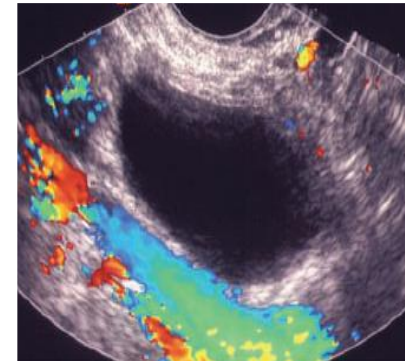
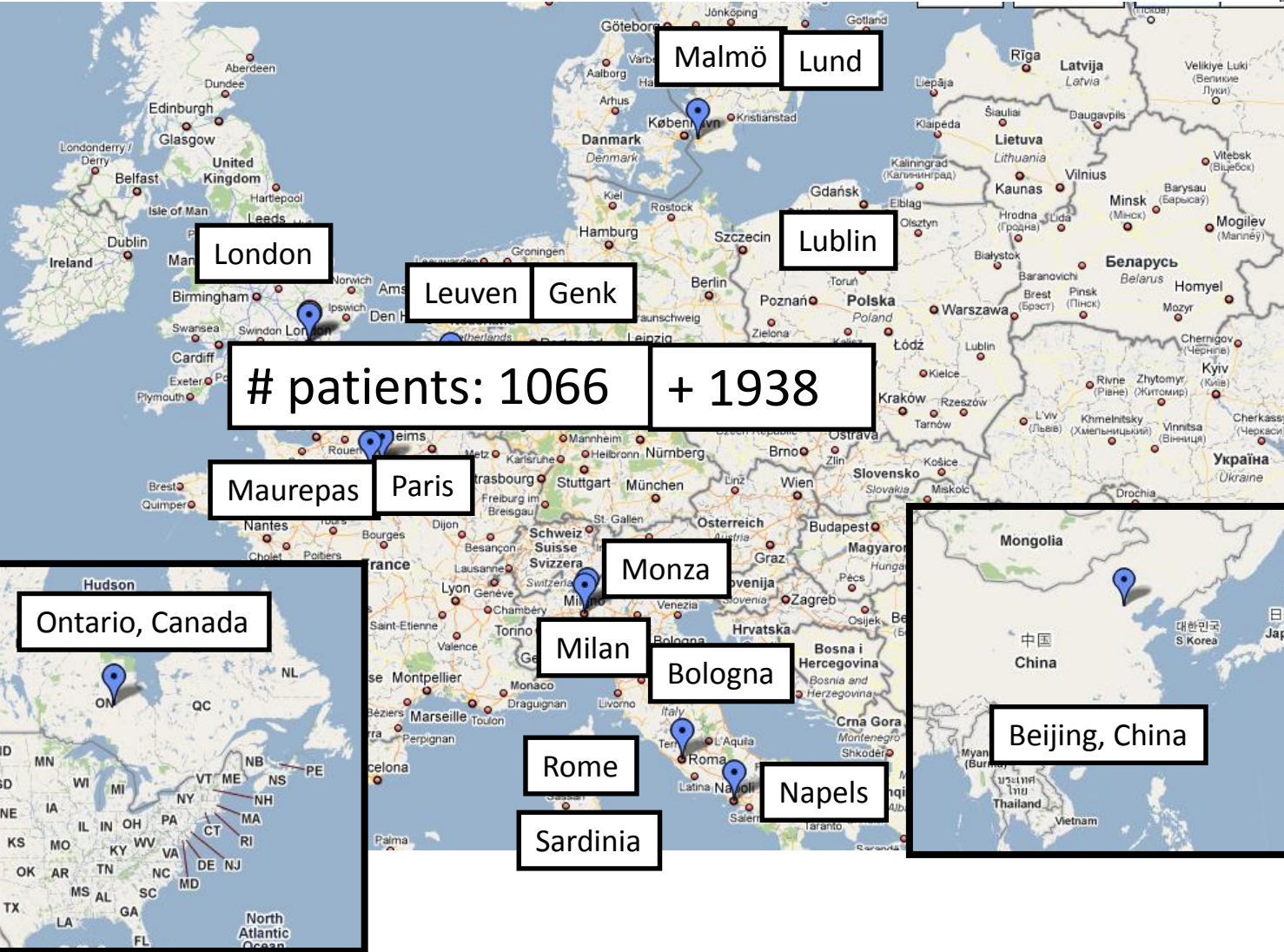
**Fraud and predictive analysis**

**Health**

# Participatory

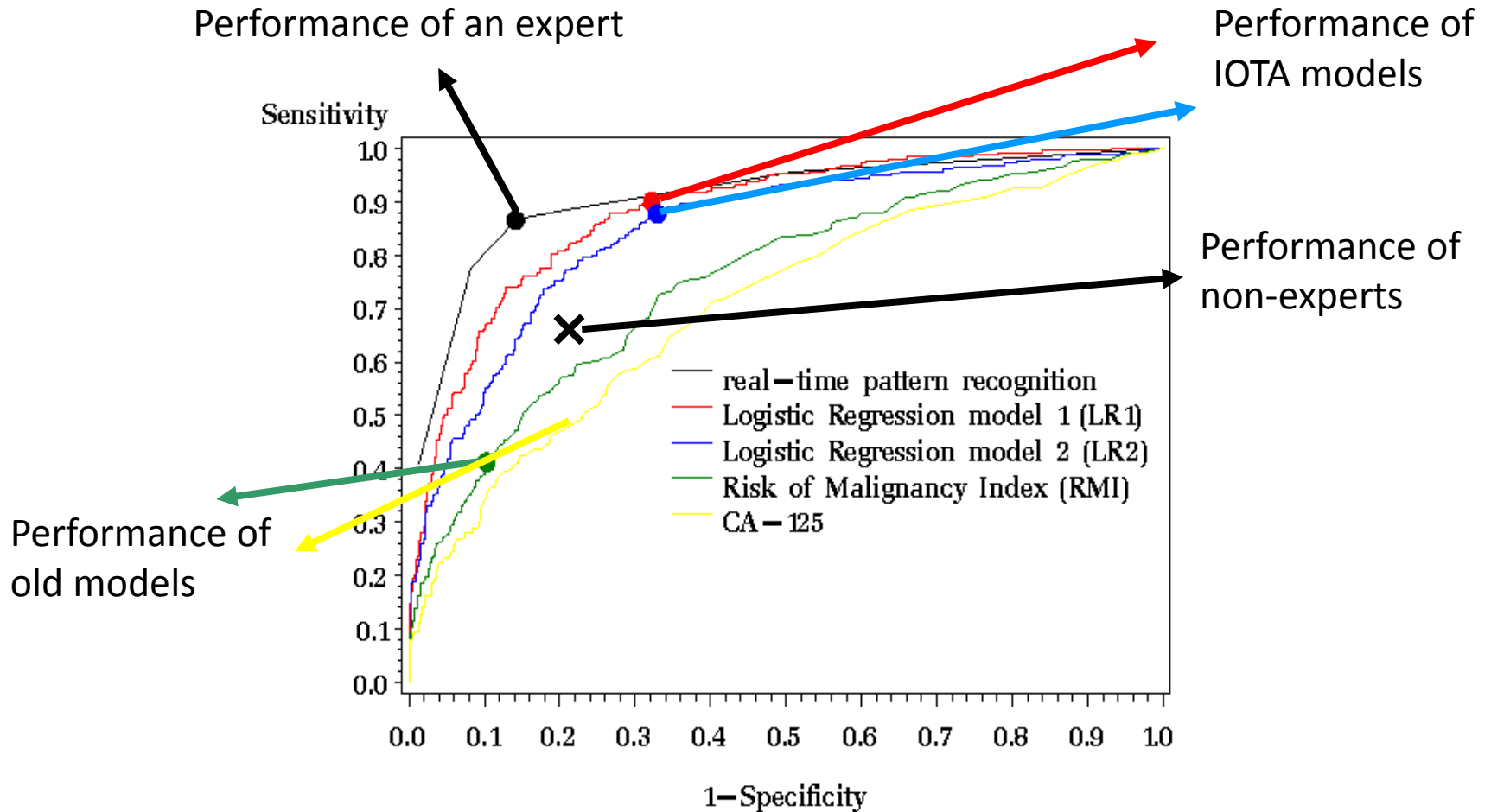


IOTA app:  
population  
based  
assessment  
of ovarian  
tumour malignancy:



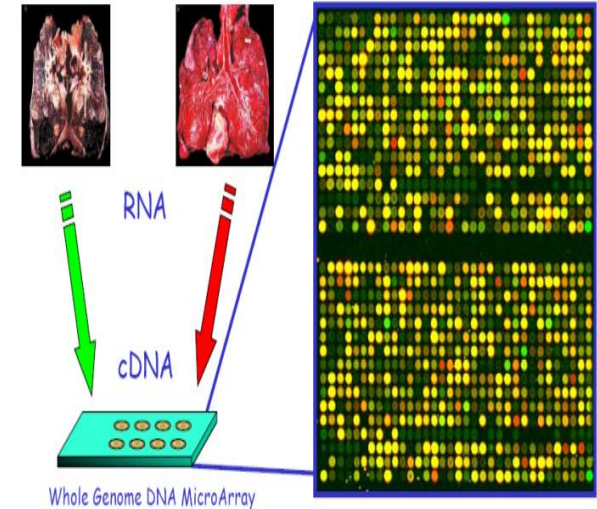
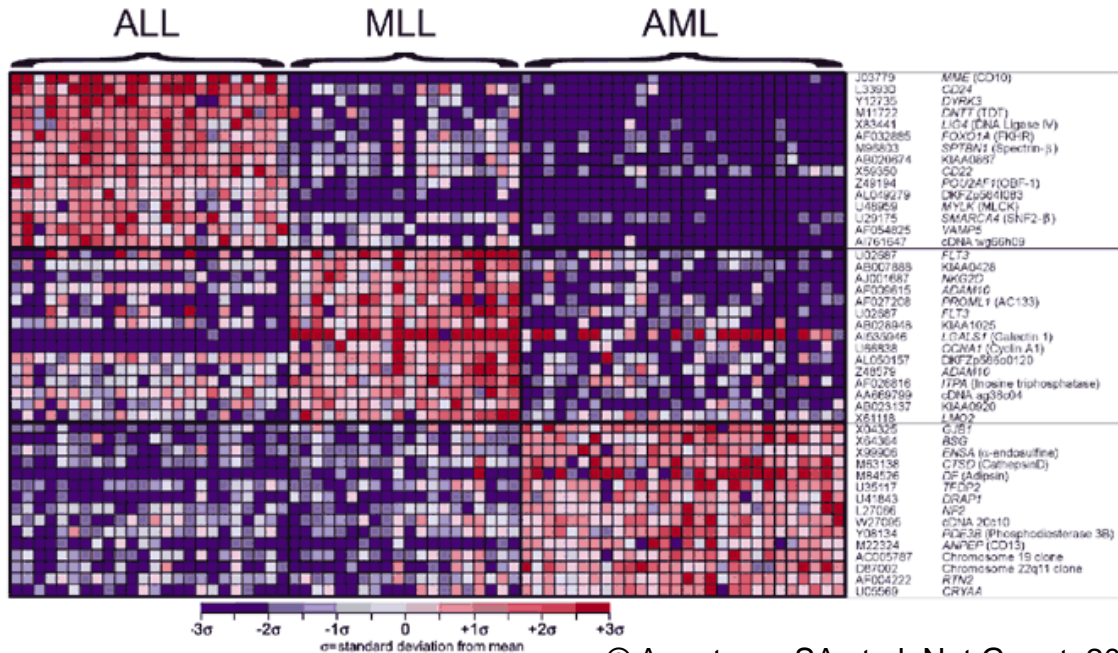


# Performance



**You share, we care !**

# Genomic markers for Leukemia



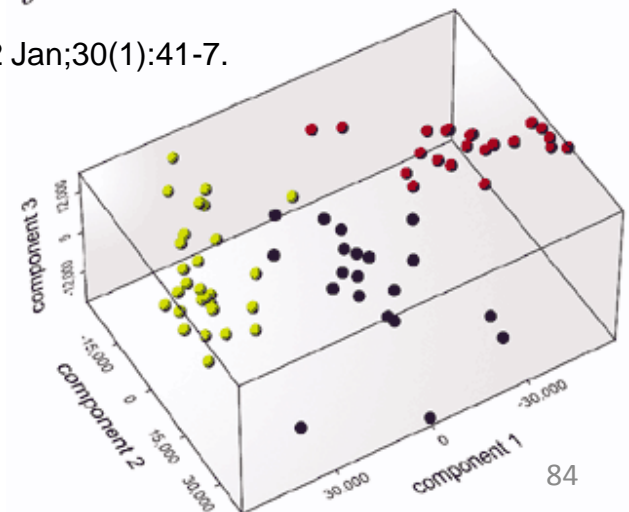
b

12 600 genes

72 patients

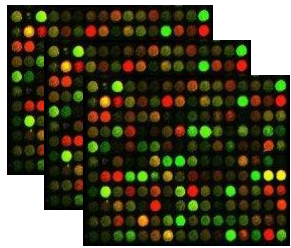
- 28 Acute Lymphoblastic Leukemia (ALL)
- 24 Acute Myeloid Leukemia (AML)
- 20 Mixed Linkage Leukemia (MLL)

© Armstrong SA et al. Nat Genet. 2002 Jan;30(1):41-7.

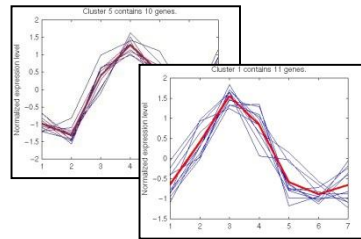


# Genomic Data Fusion

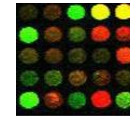
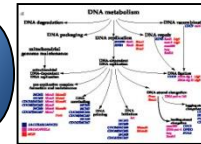
High-throughput genomics



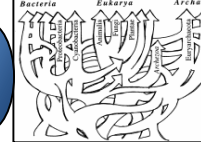
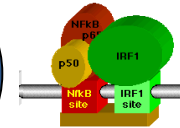
Data analysis



Information sources



After champagne, at 11, she saw Michael Jackson performing on television and told Angelil that she wanted to be that big. Fine, said Angelil, who advised her to take 18 months off, during which she underwent a massive makeover that included plastic surgery, shorter hair and caps for the long incisors that had propelled a Quebec beauty



Candidate genes

Name	Ensembl
TTR	ENSG00000118271
PAH	ENSG00000171759
G6PC	ENSG00000131482
IGF1	ENSG00000017427
ALB	ENSG00000163631
CRP	ENSG00000132693
HABP2	ENSG00000148702
IF	ENSG00000138799
FST	ENSG00000134363
ARAF1	ENSG00000078061
HMGA2	ENSG00000149948
C9	ENSG00000113600
PCBP2	ENSG00000111406
HOXB6	ENSG00000108511
RERE	ENSG00000142599
HOXA11	ENSG00000005073
CLIC1	ENSG00000096238
ERCC3	ENSG00000163161
ERCC3	ENSG00000163161
TLL2	ENSG00000095587
SYT4	ENSG00000132872
SYT4	ENSG00000132872
PIK4CB	ENSG00000143393
PKD2	ENSG00000118762
	ENSG00000081026
ANKRD3	ENSG00000183421
F13A1	ENSG00000124491
BPAG1	ENSG00000151914
KCNN3	ENSG00000143603
GRIN2A GRIN2B	ENSG00000150086
SIM1	ENSG00000112246
	ENSG00000174891
	ENSG00000089195
C14orf10	ENSG00000092020
STX8	ENSG00000170310
	ENSG00000107671
MSH5	ENSG00000096474
CRH	ENSG00000147571
MID1	ENSG00000101871
	ENSG00000184508
	ENSG00000113460
TGFB3	ENSG00000119699
C1QR1	ENSG00000125810
NR4A2	ENSG00000153234
PDGFC	ENSG00000145431
PDGFC	ENSG00000145431
NR3C2	ENSG00000151623
NFYA	ENSG00000001167
	ENSG00000101898
C8orf4	ENSG00000176907
TM4SF13	ENSG00000106537
MMP3 MMP1	ENSG00000149968
	ENSG00000135112

Candidate prioritization

Rank	En	Ex	Ip	Ke	GO	TeAvg	Pval
1	TTR	G6PC	PAH	G6PC	IGF1	TTR	TTR
2	IGF1	TTR	IGF1	PAH	PAH	IGF1	PAH
3	CRP	ALB	TTR	RERE	G6PC	CRP	G6PC
4	HOXB6	HABP2	ALB	ERCC3	TTR	HOXB6	IGF1
5	ALB	PAH	HDC	ERCC3		ALB	ALB
6	NR4A2	IF	TLL2	ANKRD3	HMGA2		CRP
7	PAH		C10R1	ARAF1	HDC	NR4A2	HABP2
8	HOXA11	IGF1	G6PC	PKD2	F13A1	PAH	IF
9	NFYA	CRP	HABP2	MTMR1	KCNN3	HOXA11	C13orf7
10	C9	ARAF1	IF	HDC	CLIC1	NFYA	TTR
							ARAF1

Validation



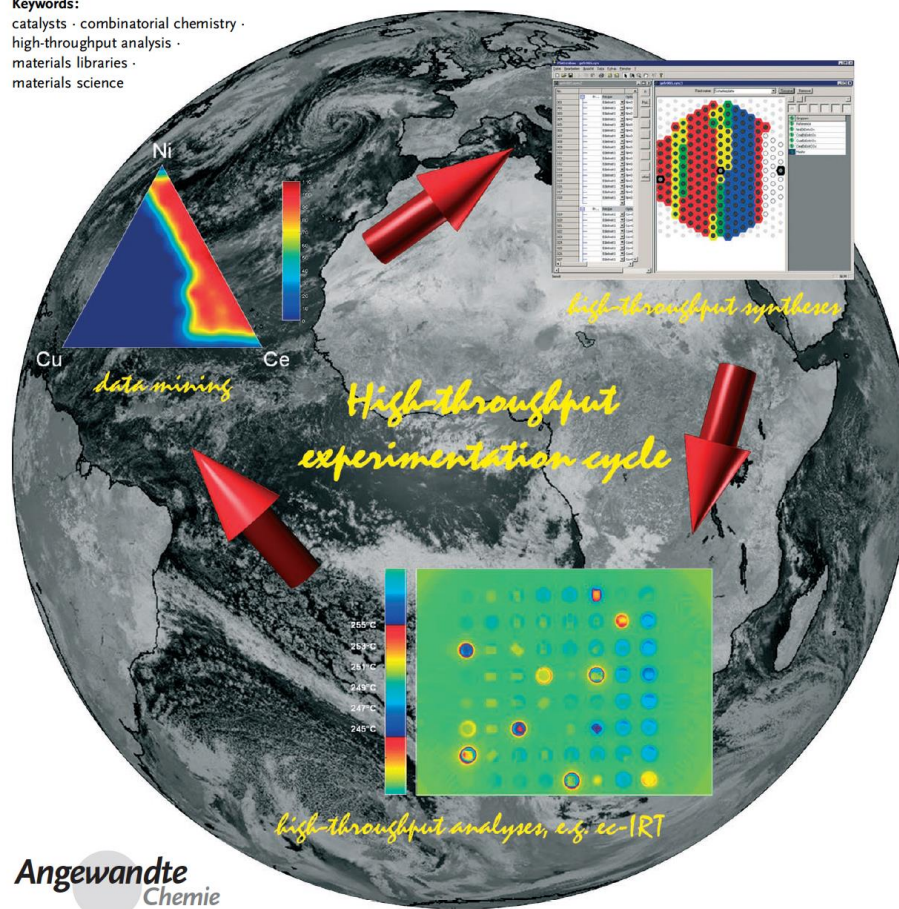


# Combinatorial and High-Throughput Materials Science

Wilhelm F. Maier,\* Klaus Stöwe, and Simone Sieg

**Keywords:**

catalysts · combinatorial chemistry ·  
high-throughput analysis ·  
materials libraries ·  
materials science



**Angewandte**  
Chemie

6016 www.angewandte.org

© 2007 Wiley-VCH Verlag GmbH & Co. KGaA, Weinheim

Angew. Chem. Int. Ed. 2007, 46, 6016–6067



Contents lists available at SciVerse ScienceDirect

# Electrochimica Acta

journal homepage: [www.elsevier.com/locate/electacta](http://www.elsevier.com/locate/electacta)



## Review

# A review of high throughput and combinatorial electrochemistry

T.H. Muster<sup>a,\*</sup>,<sup>1</sup>, A. Trinchi<sup>a</sup>, T.A. Markley<sup>a</sup>, D. Lau<sup>a</sup>, P. Martin<sup>b</sup>, A. Bradbury<sup>a</sup>, A. Bendavid<sup>b</sup>, S. Dligatch<sup>b</sup>

<sup>a</sup> CSIRO Division of Materials Science and Engineering, Private Bag 33, Clayton South, Victoria 3169, Australia

<sup>b</sup> CSIRO Division of Materials Science and Engineering, PO Box 218, Lindfield, NSW, 2070, Australia

## ARTICLE INFO

### Article history:

Received 20 June 2011

Accepted 2 September 2011

Available online 14 September 2011

### Keywords:

Electrochemistry

Combinatorial

High-throughput

Multielectrode

Review

## ABSTRACT

Many 21st century technological solutions are reliant on the development of new materials with improved properties, and increasingly on materials that can be optimised to perform more than one function. High-throughput and combinatorial methodologies are being used more frequently to discover and design improved materials in a time efficient manner for a variety of applications. A number of technological challenges involve the field of electrochemistry, such as battery development, electrocatalysis, photocatalysis, corrosion protection, sensor development, photovoltaics and light-emitting materials. This review focuses on the utilisation of high-throughput and combinatorial methods that have incorporated, or are associated with, electrochemical methods. In many cases electrochemical determinations are well-suited for high-throughput methodologies, enabling direct quantitative analysis of properties. However, in other circumstances electrochemical measurements are complicated by additional factors. Hence the limitations of high-throughput and combinatorial electrochemistry are also discussed within.

Top Curr Chem (2014) 345: 139–180  
DOI: 10.1007/128\_2013\_486  
© Springer-Verlag Berlin Heidelberg 2013  
Published online: 28 November 2013

## Data Mining Approaches to High-Throughput Crystal Structure and Compound Prediction

Geoffroy Hautier

**Abstract** Predicting unknown inorganic compounds and their crystal structure is a critical step of high-throughput computational materials design and discovery. One way to achieve efficient compound prediction is to use data mining or machine learning methods. In this chapter we present a few algorithms for data mining compound prediction and their applications to different materials discovery problems. In particular, the patterns or correlations governing phase stability for experimental or computational inorganic compound databases are statistically learned and used to build probabilistic or regression models to identify novel compounds and their crystal structures. The stability of those compound candidates is then assessed using ab initio techniques. Finally, we report a few cases where data mining driven computational predictions were experimentally confirmed through inorganic synthesis.

**Keywords** Ab initio computations · Crystal structure prediction · Data mining · High-throughput computing