



**iMinds Dept.
MEDICAL IT**

KU Leuven ESAT-STADIUS

Serious Data, Serious Mining

Prof.Dr. Bart De Moor

Bart.DeMoor@iminds.be



Big Data

What

Who

Six dimensions

Data

Compute Infrastructure

Storage Infrastructure

Analytics

Visualization

Security & Privacy

Machine learning as a commodity

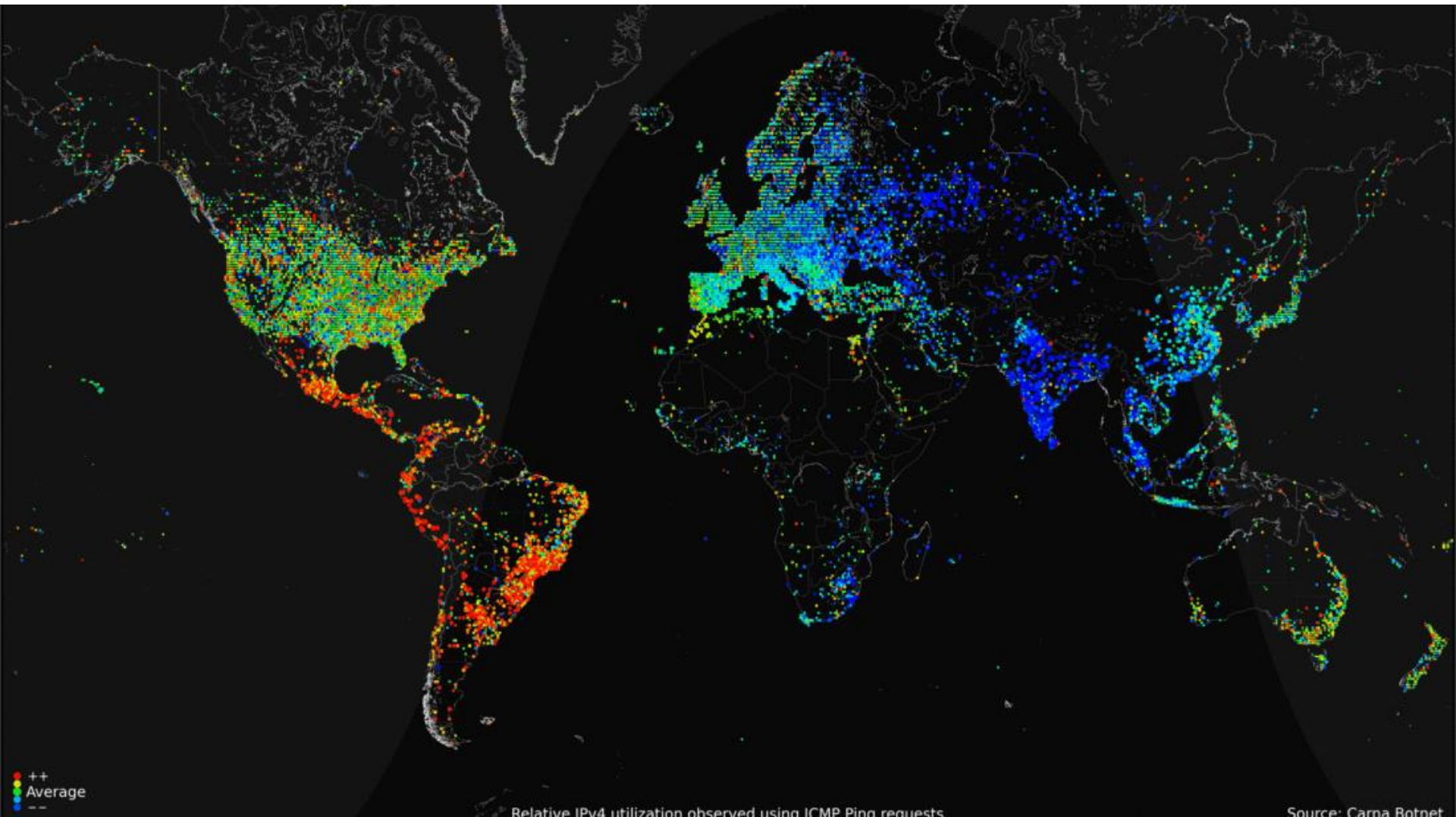
Expertise

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.



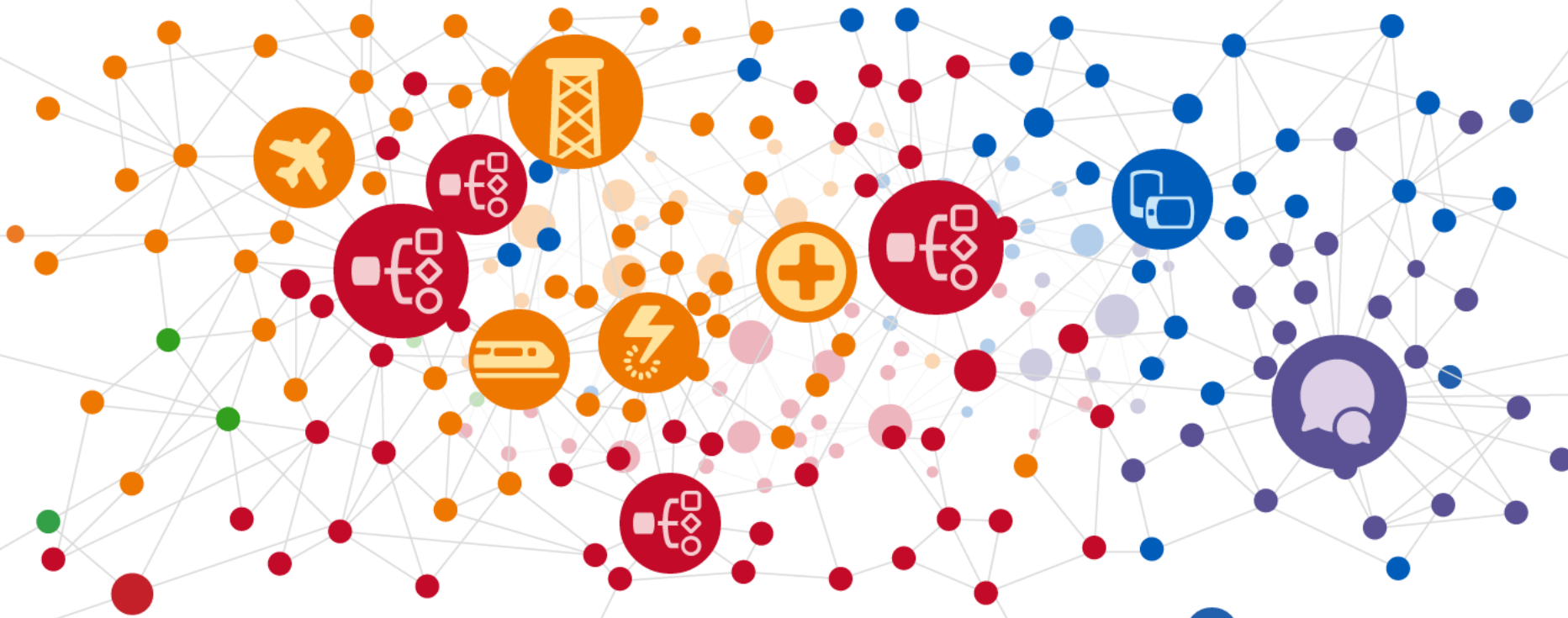
The Industrial Internet

5 2 0 0 0 0 0 0 0 0

Connected Devices



Devices Per Capita Worldwide



Past

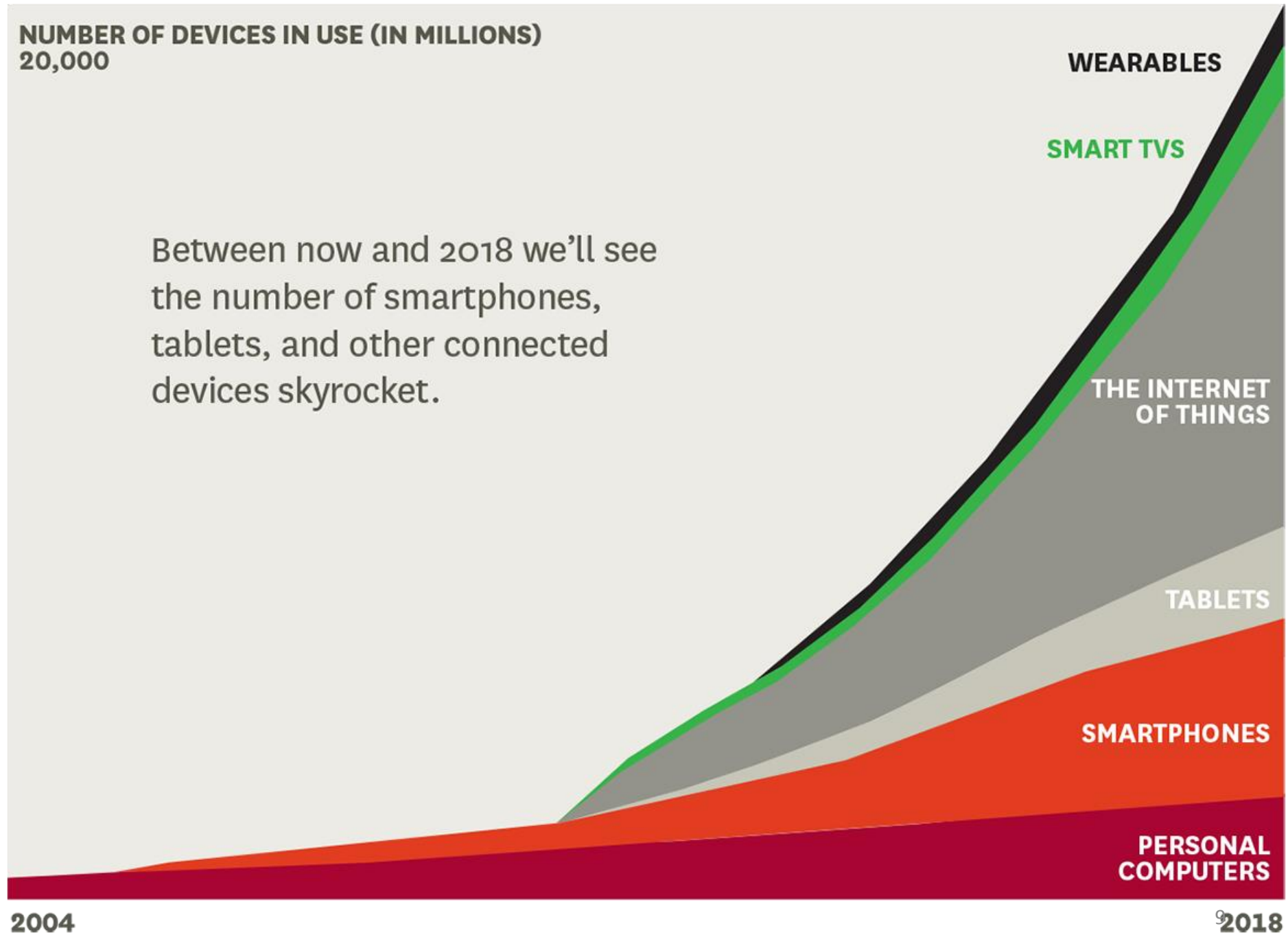


2020+

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.

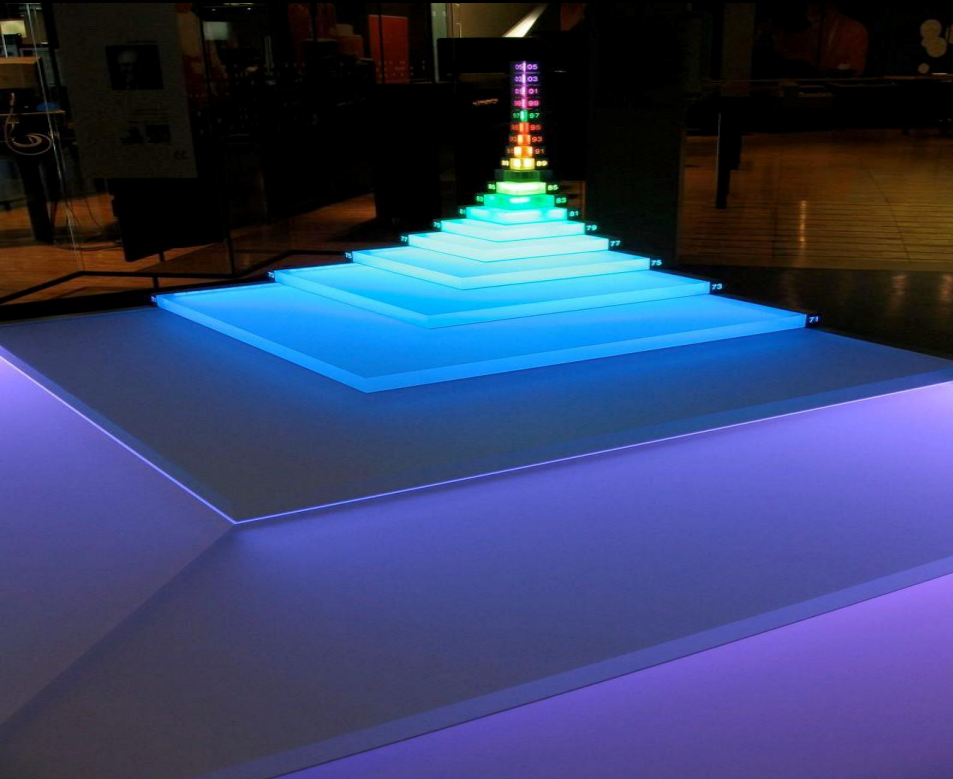
NUMBER OF DEVICES IN USE (IN MILLIONS)
20,000

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



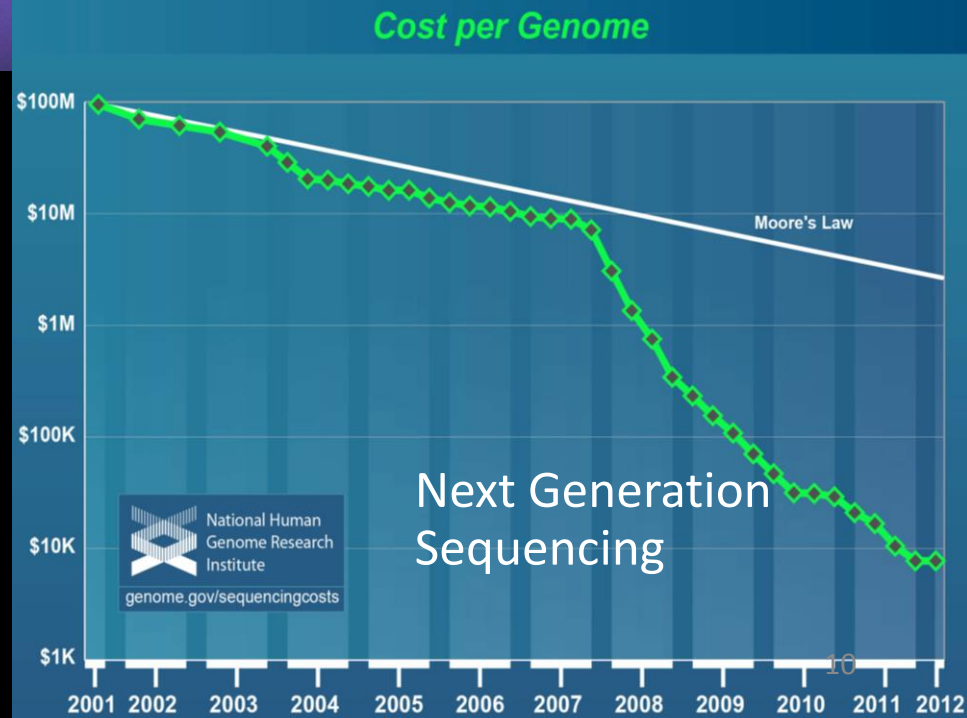
2004

2018



Moore's law:
 computing power
 doubles
 every 18 months

Carlson's law:
 complexity/cost
 evolves
 exponentially



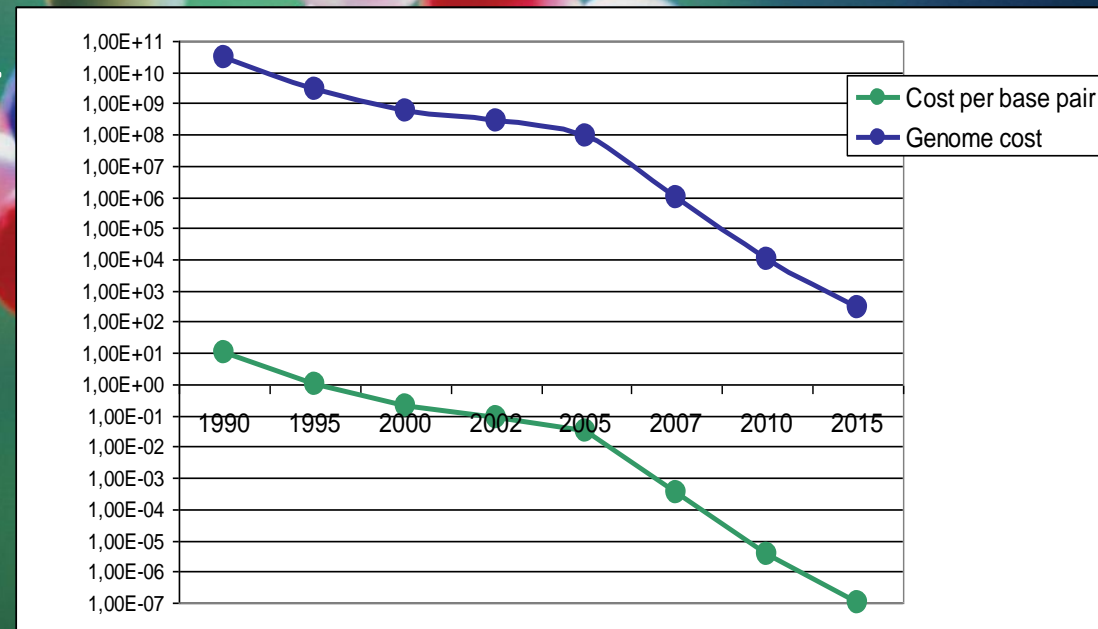
Genome data

- **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 μ m
resolution

81 GigaByte

1 CD-ROM

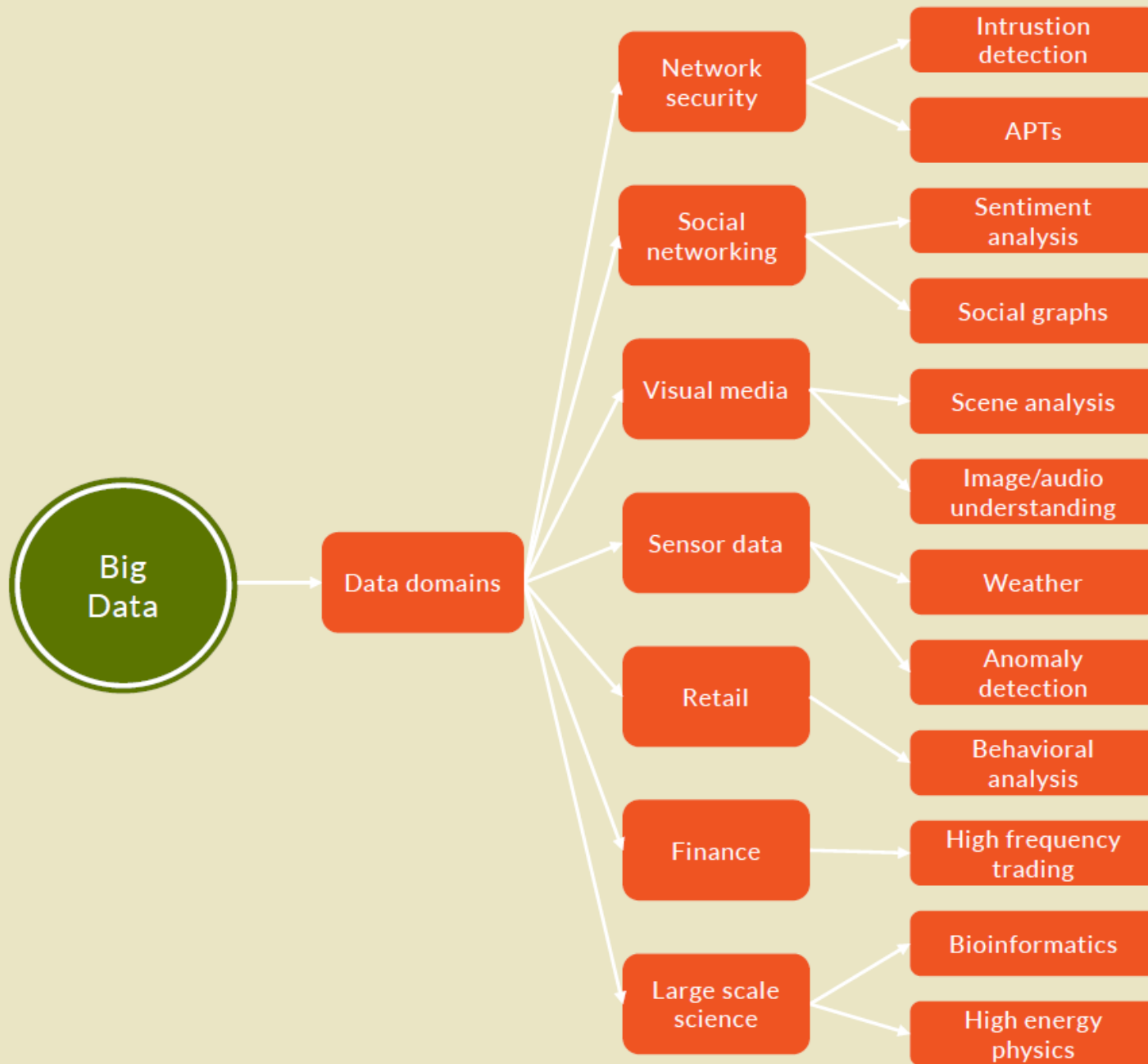
750
MegaByte

Data explosion in finance



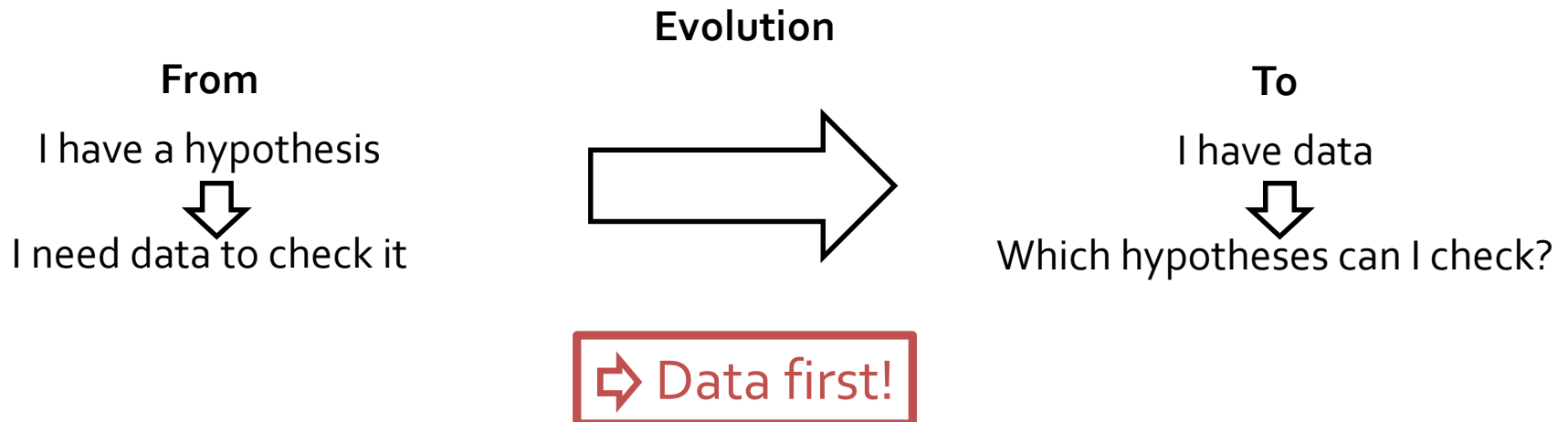
**Growing ~ 30-50% every year,
half of this is unstructured!**

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



The Fourth Paradigm

Paradigm	Time Ago	Method
First	A millenium	Empirical
Second	A few centuries	Theoretical
Third	A few decades	Computational
Fourth	Today	Data-driven





Big Data

What

Who

Six dimensions

Data

Compute Infrastructure

Storage Infrastructure

Analytics

Visualization

Security & Privacy

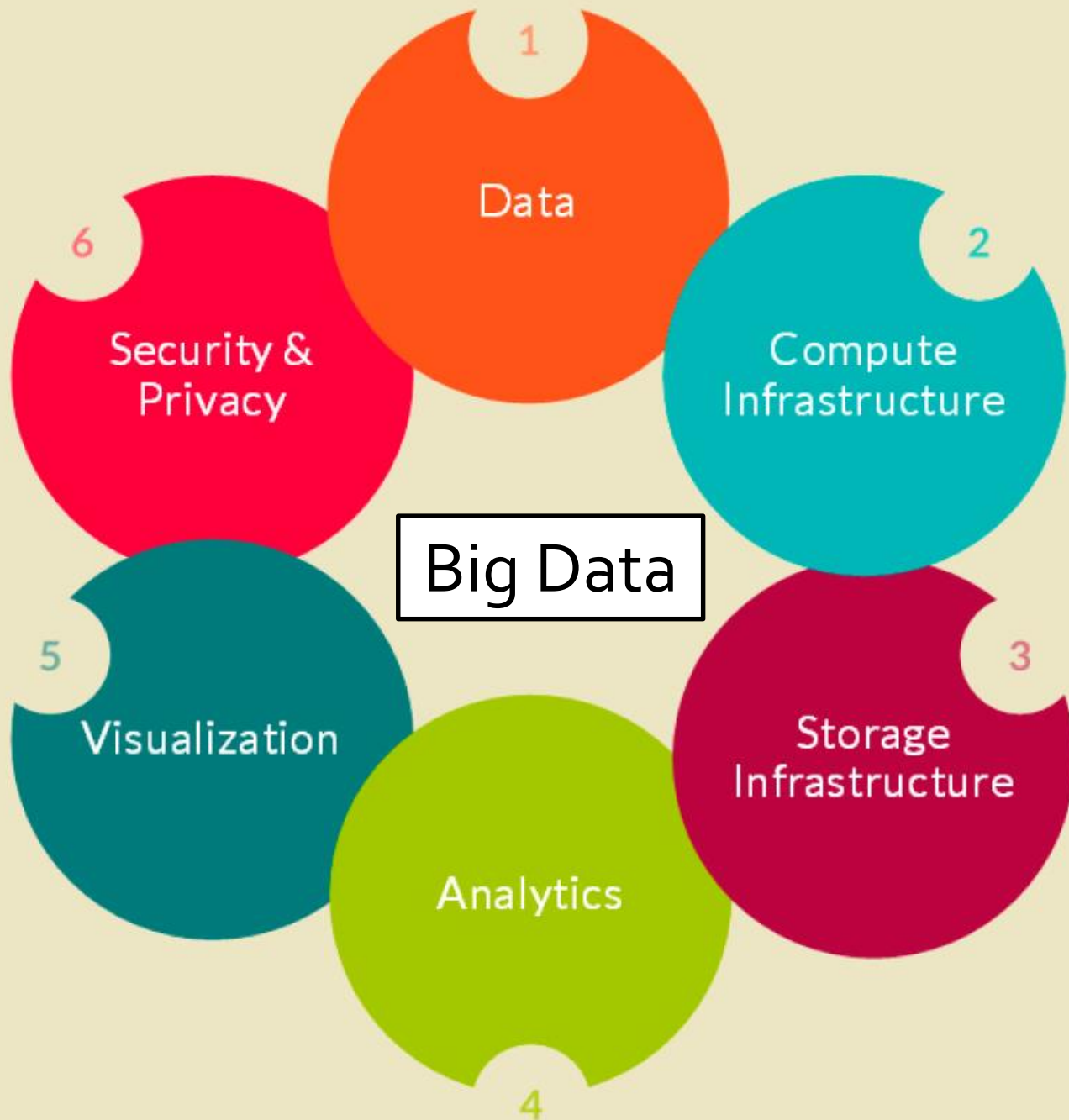
Machine learning as a commodity

Expertise

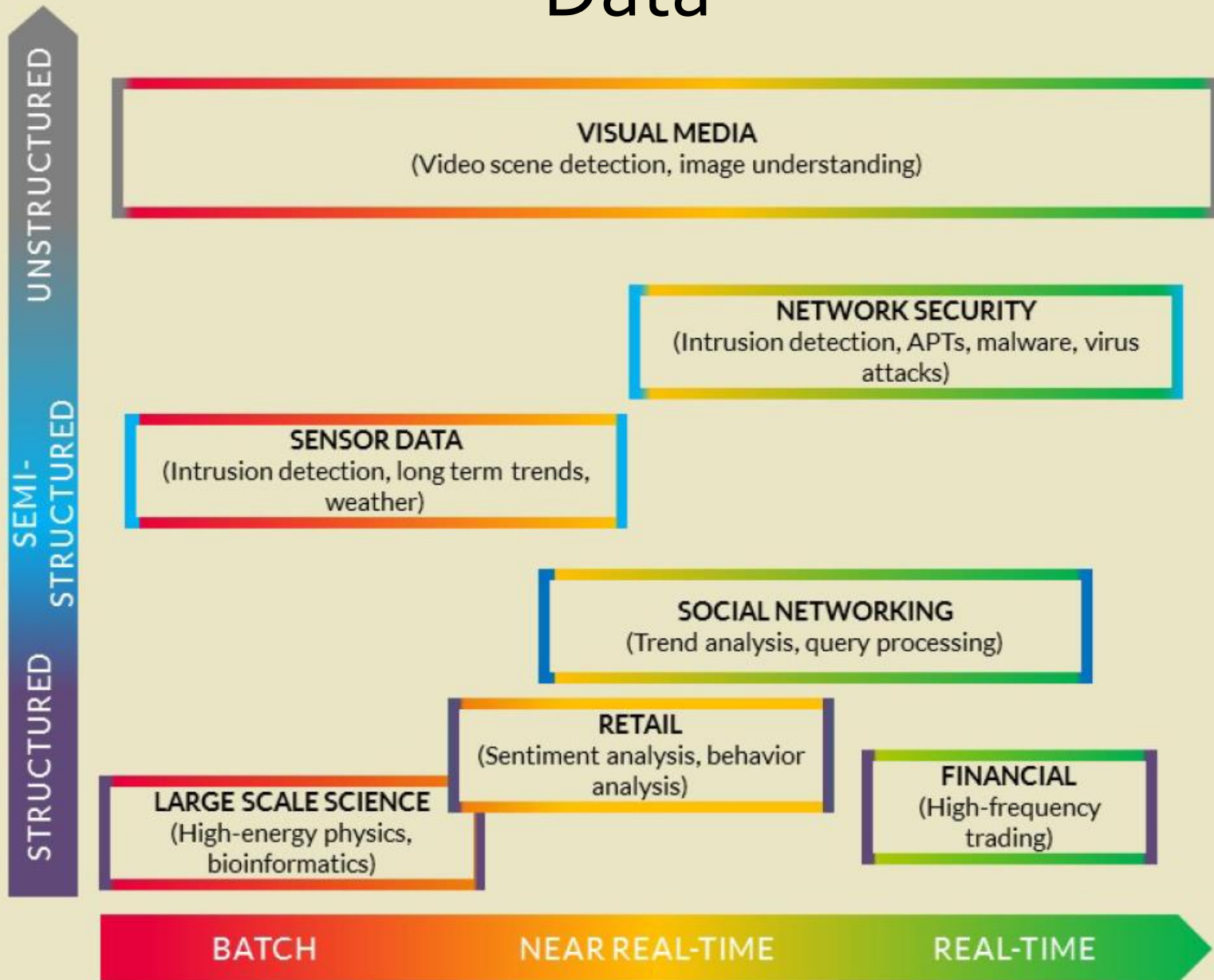
Books & Spin-offs

Algorithms

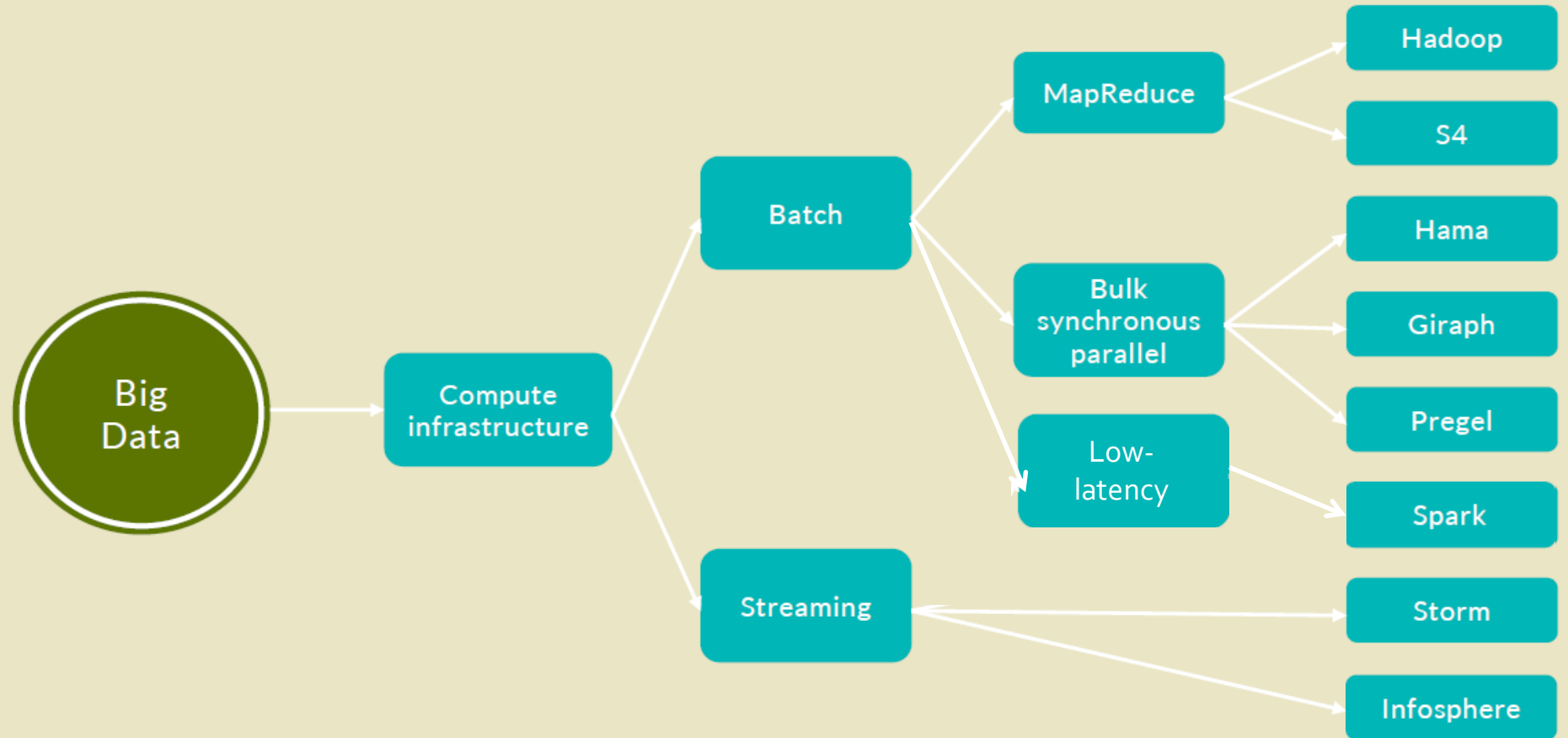
Applications



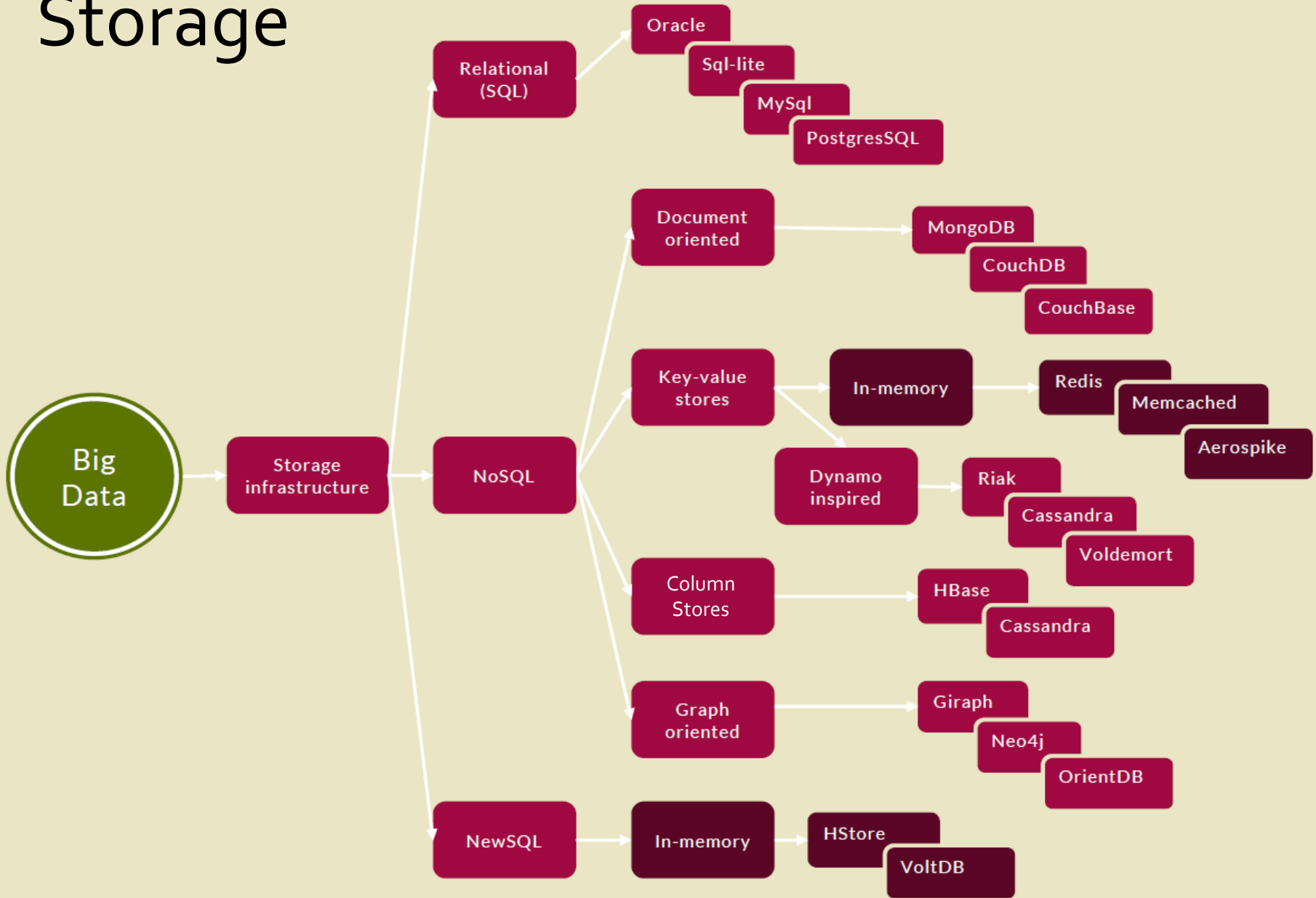
Data



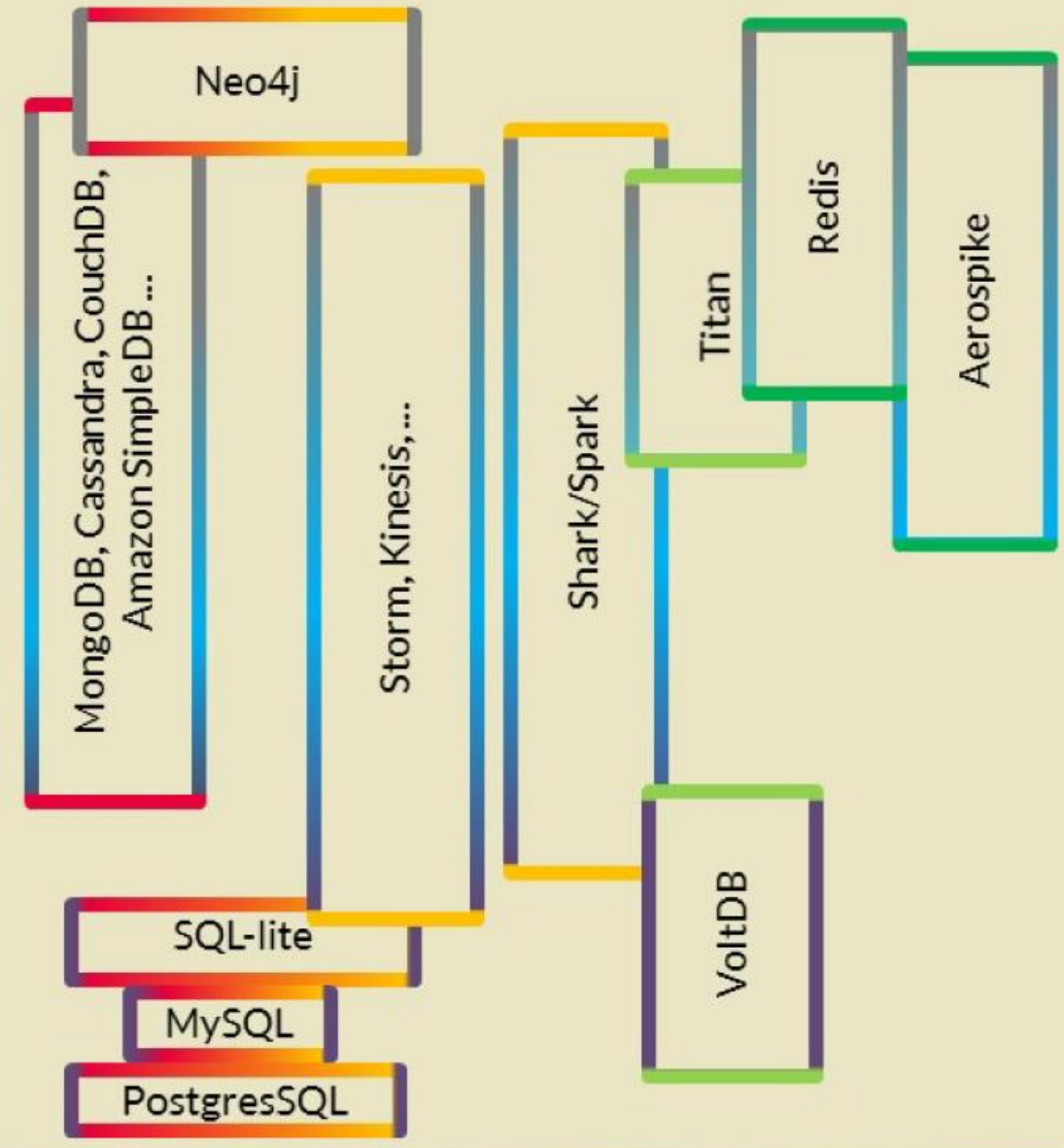
Compute infrastructure



Storage

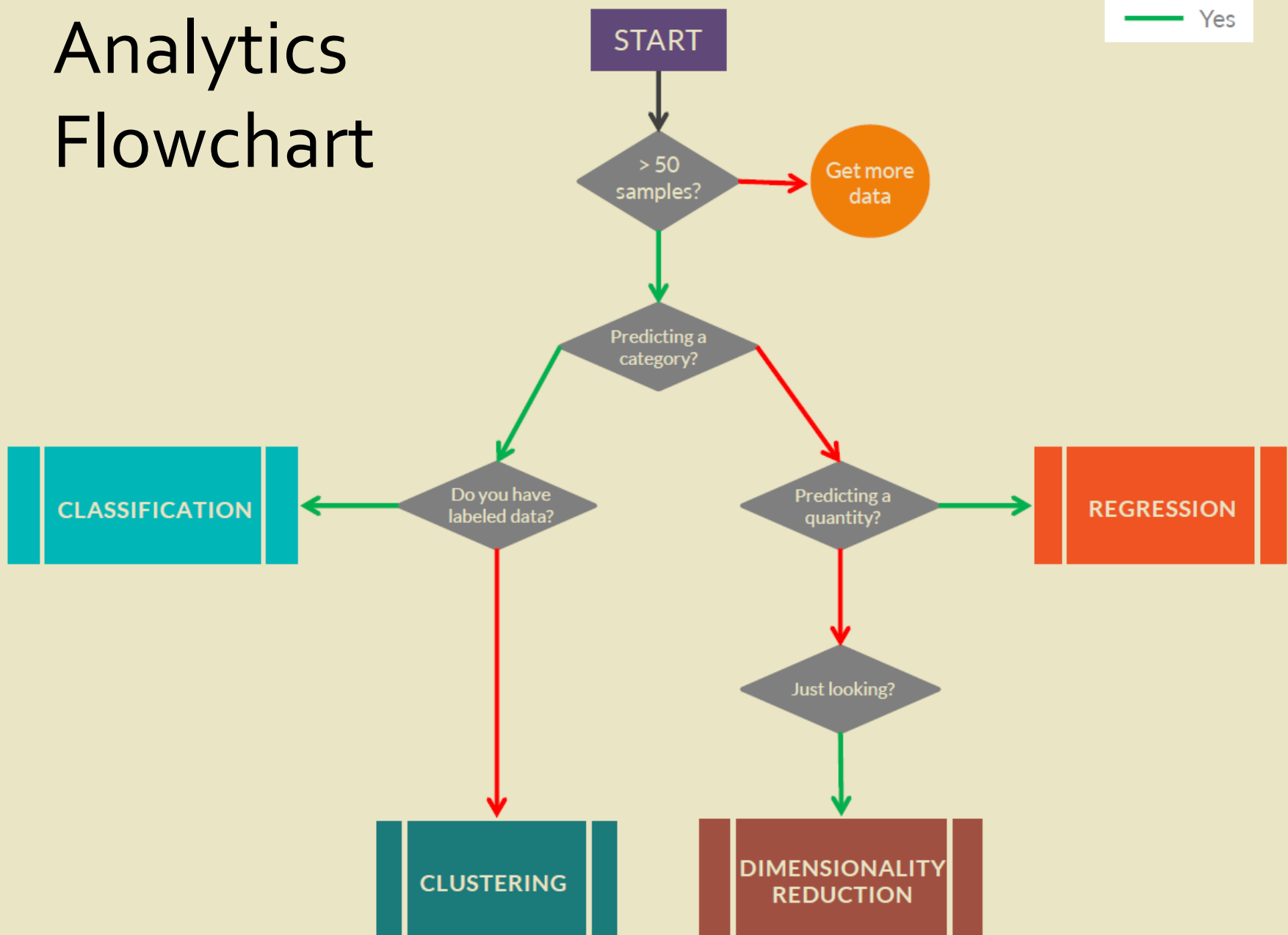


STRUCTURED SEMI-STRUCTURED UNSTRUCTURED

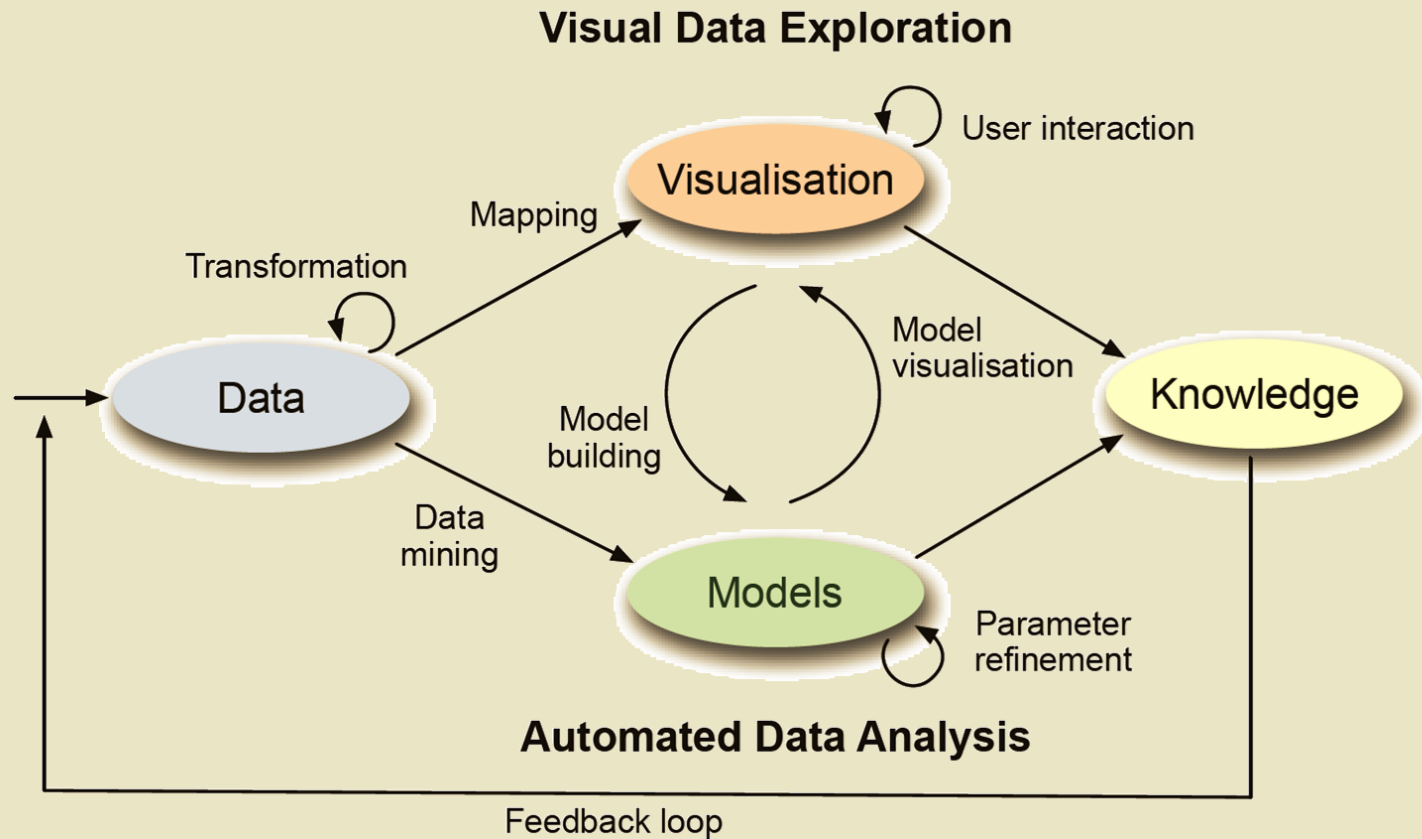


BATCH NEAR REAL-TIME REAL-TIME

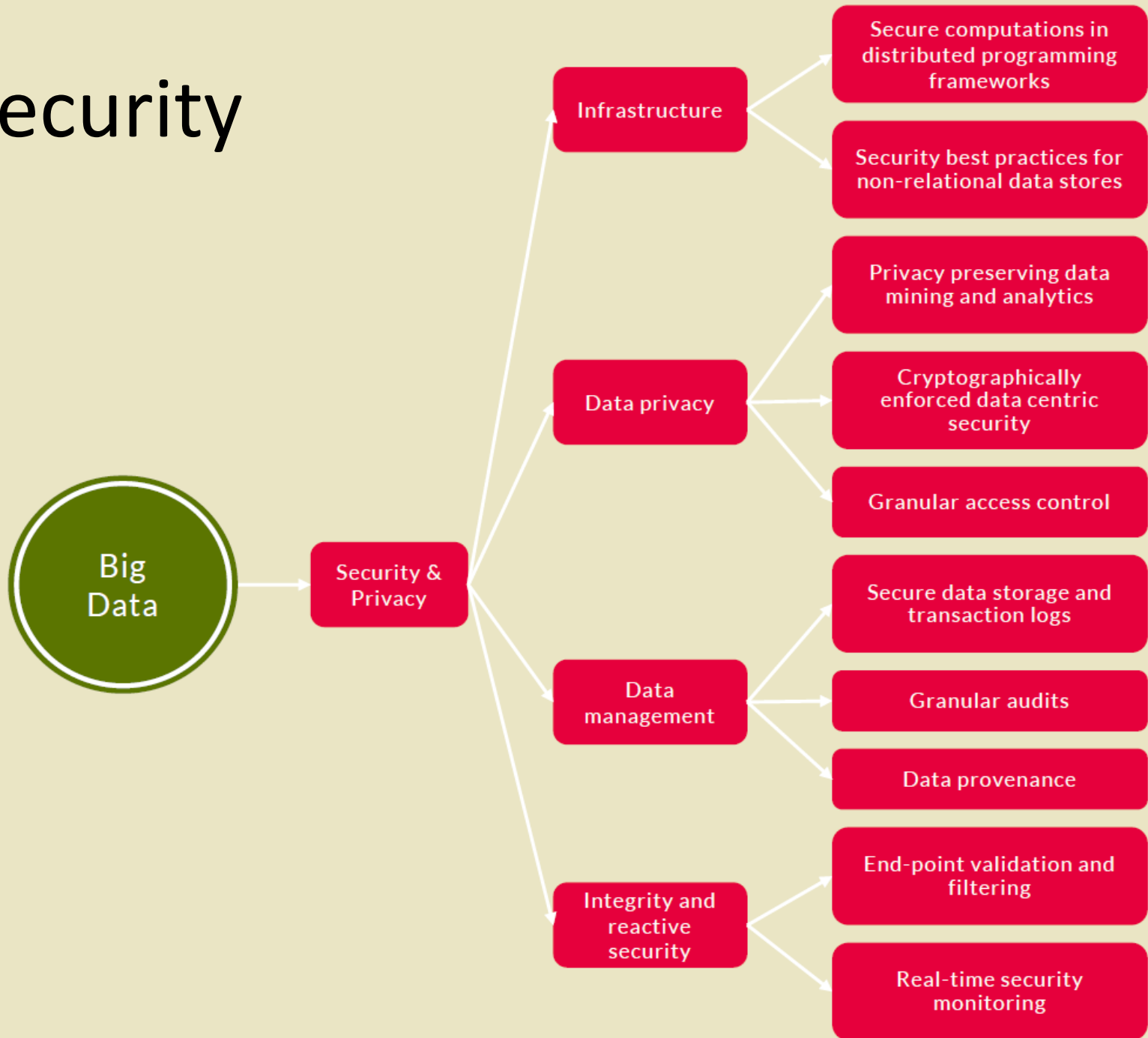
Analytics Flowchart



Visualization



Security





Big Data

What

Who

Six dimensions

Data

Compute Infrastructure

Storage Infrastructure

Analytics

Visualization

Security & Privacy

Machine learning as a commodity

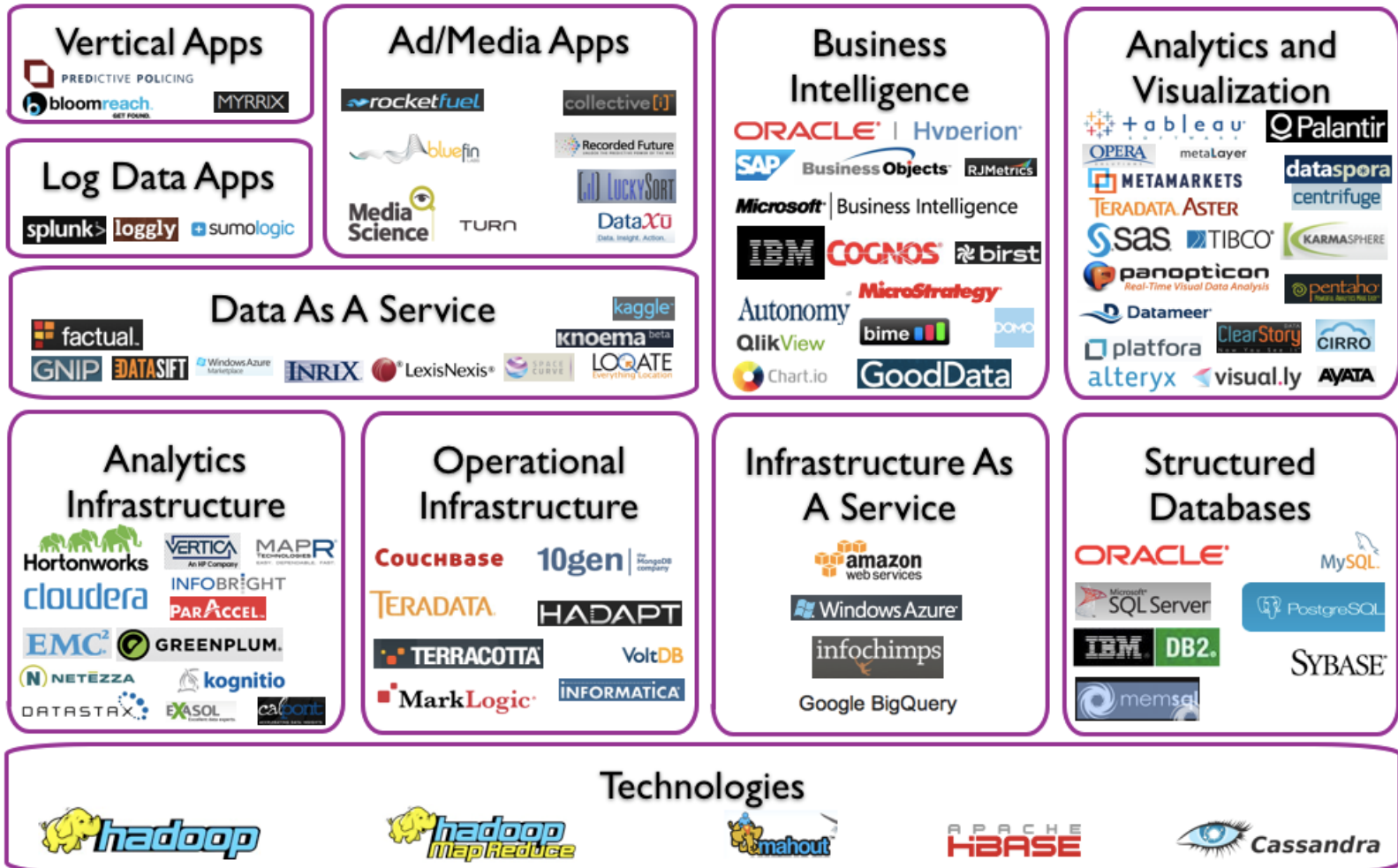
Expertise

Books & Spin-offs

Algorithms

Applications

Big Data Landscape



Big Data Landscape (Version 2.0)



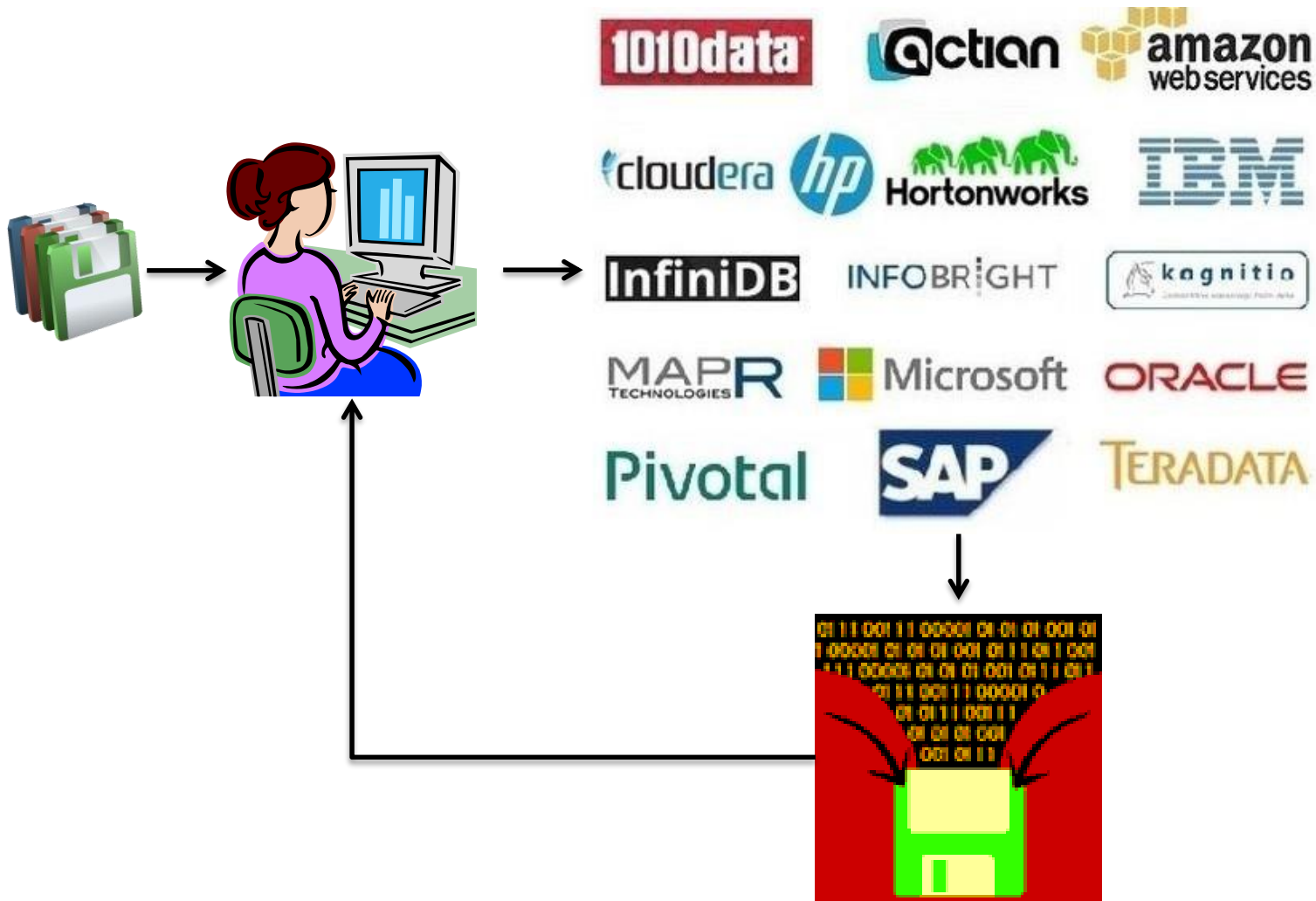
© Matt Turck (@mattturck) and ShivonZilis (@shivonz) Bloomberg Ventures

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



Big Data

What

Who

Six dimensions

Data

Compute Infrastructure

Storage Infrastructure

Analytics

Visualization

Security & Privacy

Machine learning as a commodity

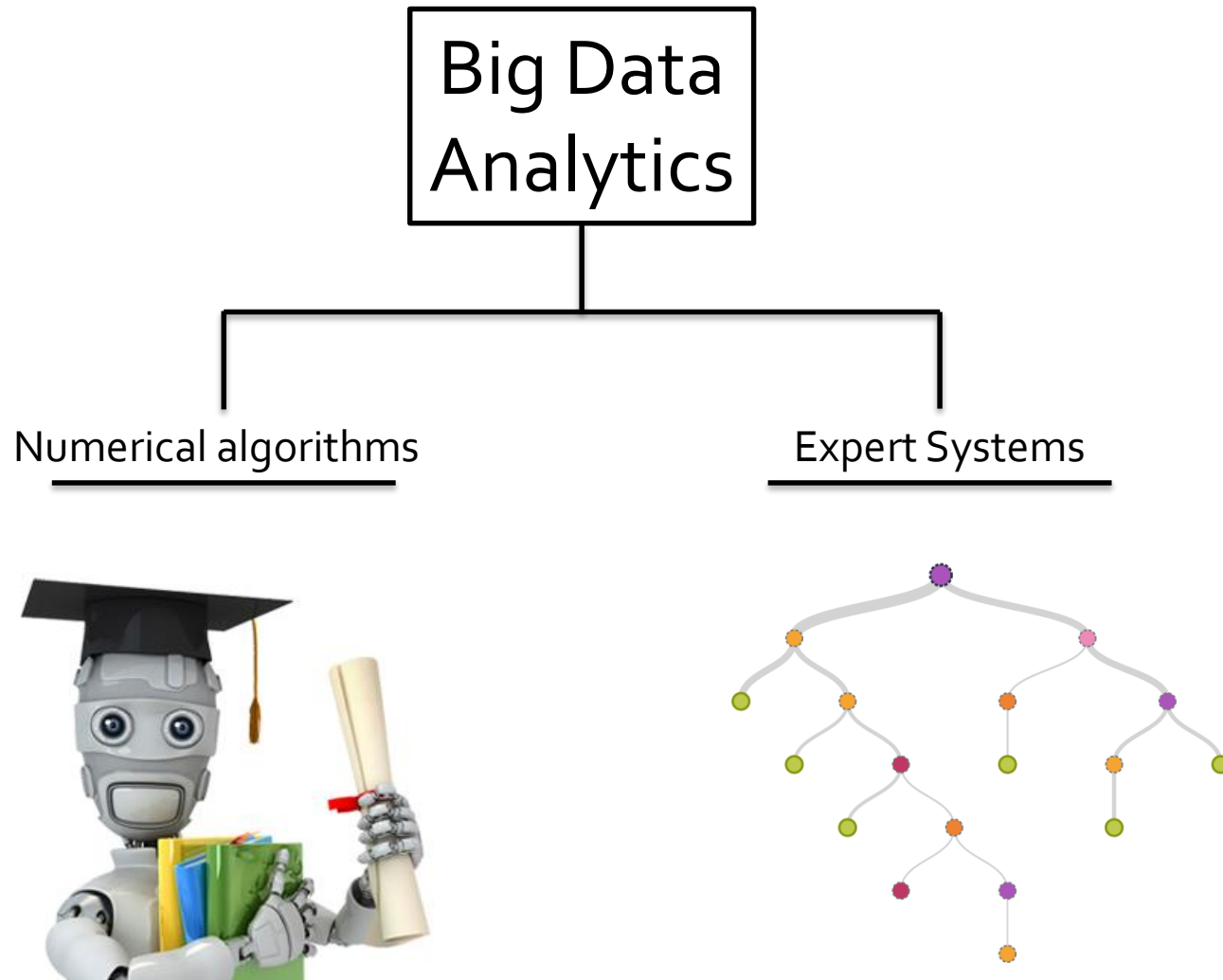
Expertise

Books & Spin-offs

Algorithms

Applications

Analytics



Objectives - ICT

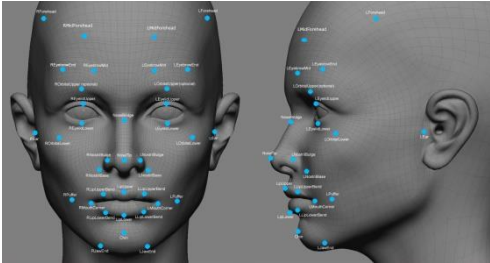
Communication networks



Home automation



Facial recognition



Digital signing



Data center optimization



Objectives - Finance

Fraud detection



Credit worthiness



Portfolio management

Order	Symbol	Company	Price	Change	% Change	Volume	Market Cap	PE Ratio	Dividend Yield	52-Week High	52-Week Low	Chart
*	HRP	7.01	0.10	1.4%	1.17M	\$400.70	0.00%	27.08%				
*	SLE	14.66	0.14	1.0%	2.20M	\$291.20	-0.30%	14.42%				
*	NWS	19.41	0.10	0.5%	1.10M	\$100.00	0.00%	14.42%				
*	MO	20.43	0.10	0.5%	1.10M	\$246.30	0.10%	14.42%				
*	HRB	22.55	0.10	0.4%	1.10M	\$411.00	0.00%	14.42%				
*	CAG	23.51	0.10	0.4%	1.10M	\$25.10	-0.30%	14.42%				
*	FRE	27.09	0.10	0.4%	1.10M	\$79.90	0.00%	14.42%				
*	HAL	45.23	0.10	0.2%	1.10M	\$94.00	0.00%	14.42%				
*	HUM	47.77	0.10	0.2%	1.10M	\$95.40	0.00%	14.42%				
*	DGX	49.79	0.10	0.2%	1.10M	\$95.80	0.00%	14.42%				
*	K	52.40	0.10	0.2%	1.10M	\$24.00	0.00%	14.42%				

Risk assessment

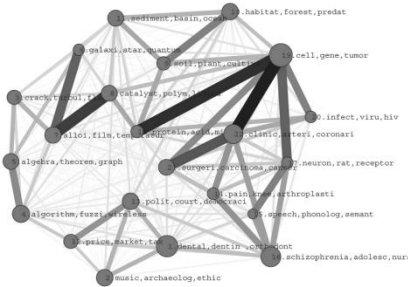


Just-in-time production



Objectives - Education

Scientometrics



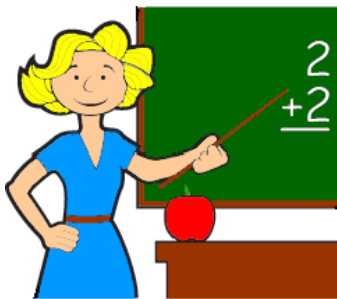
Detecting plagiarism



Grading



Teacher performance



Student performance



Objectives – Smart Cities

Predictive maintenance



Flood prediction



Smart lighting



Traffic management



Electricity Demand

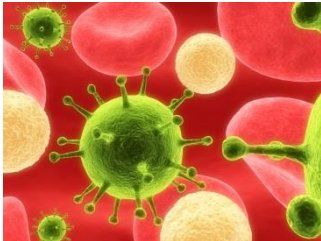


Objectives – Health

Diagnostics



Disease spreading



Genome sequencing



Tumour detection

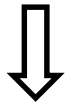


Medical fraud detection



Main tasks

Prediction



Regression

Segmentation



Clustering

Classification

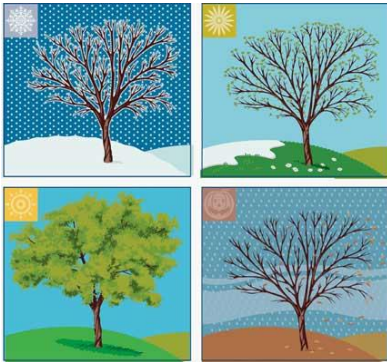
Anomalies



Detect outliers

Main tasks

Filtering effects



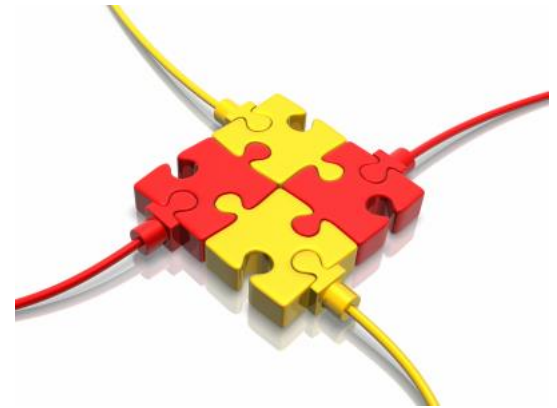
Normalization

Assess relevance



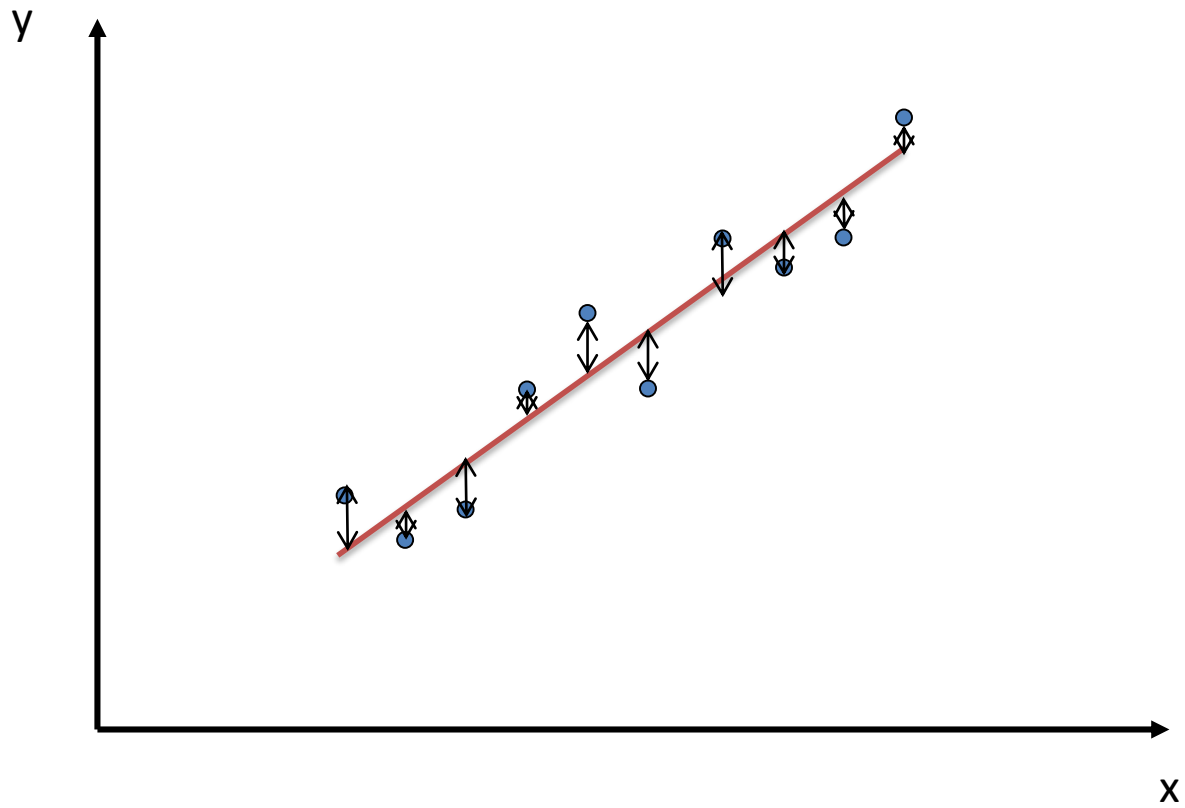
Ranking

Combining info

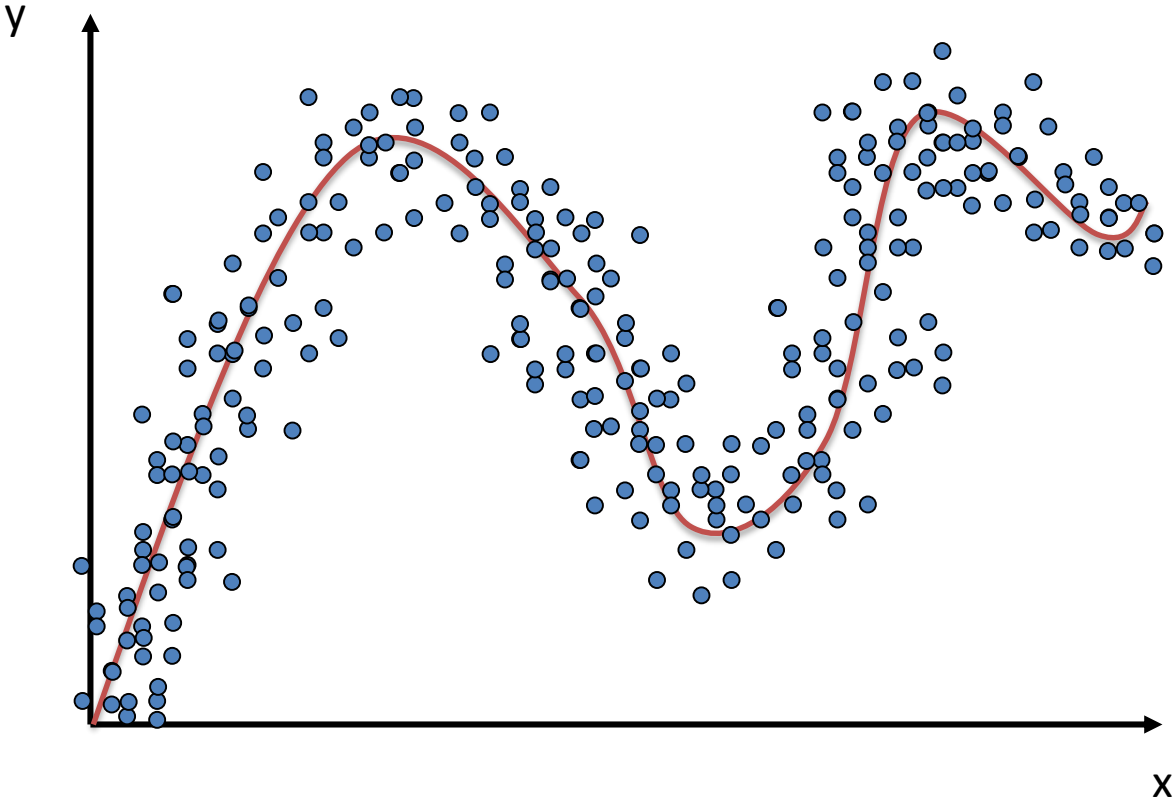


Data fusion

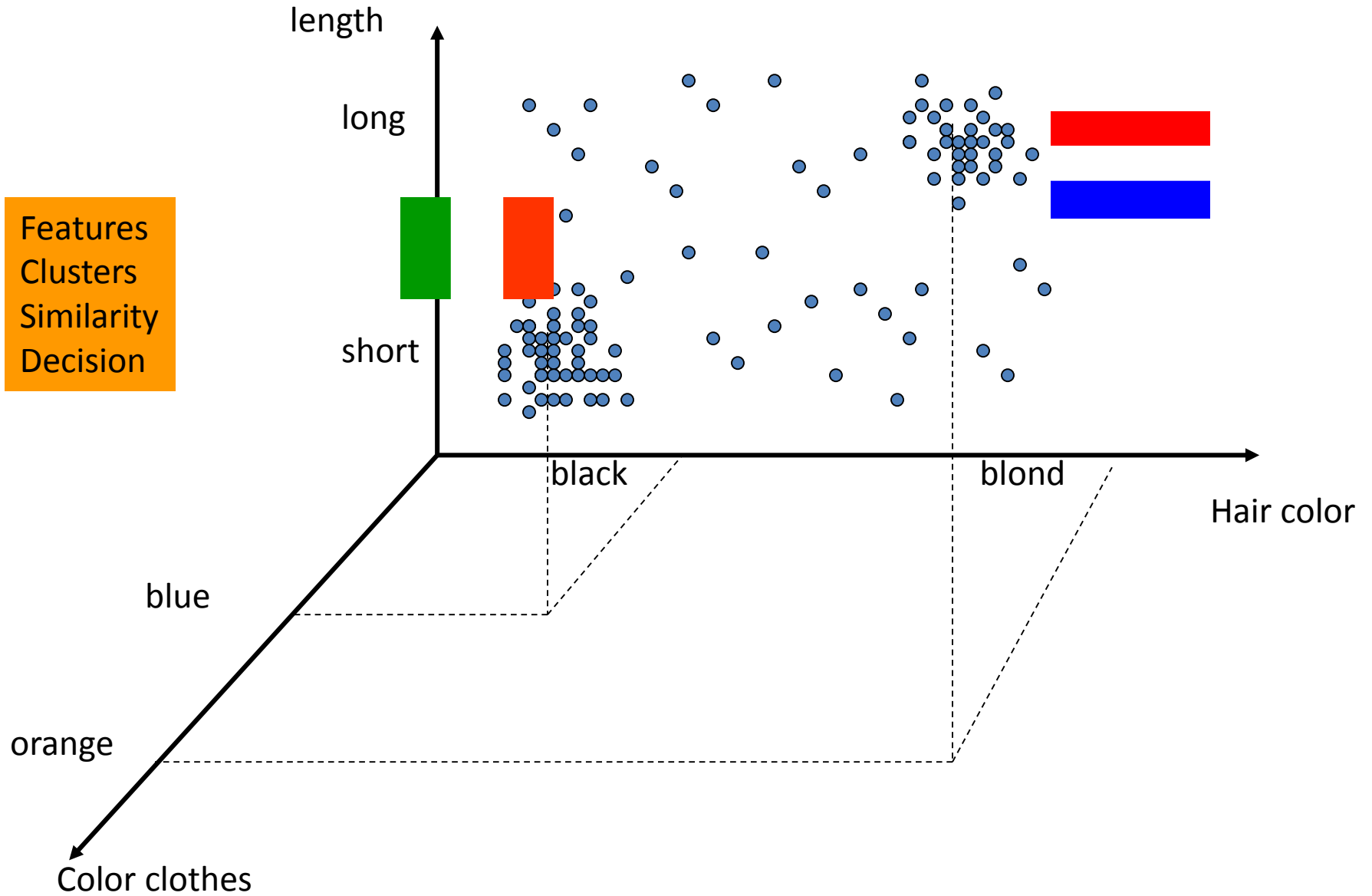
Regression



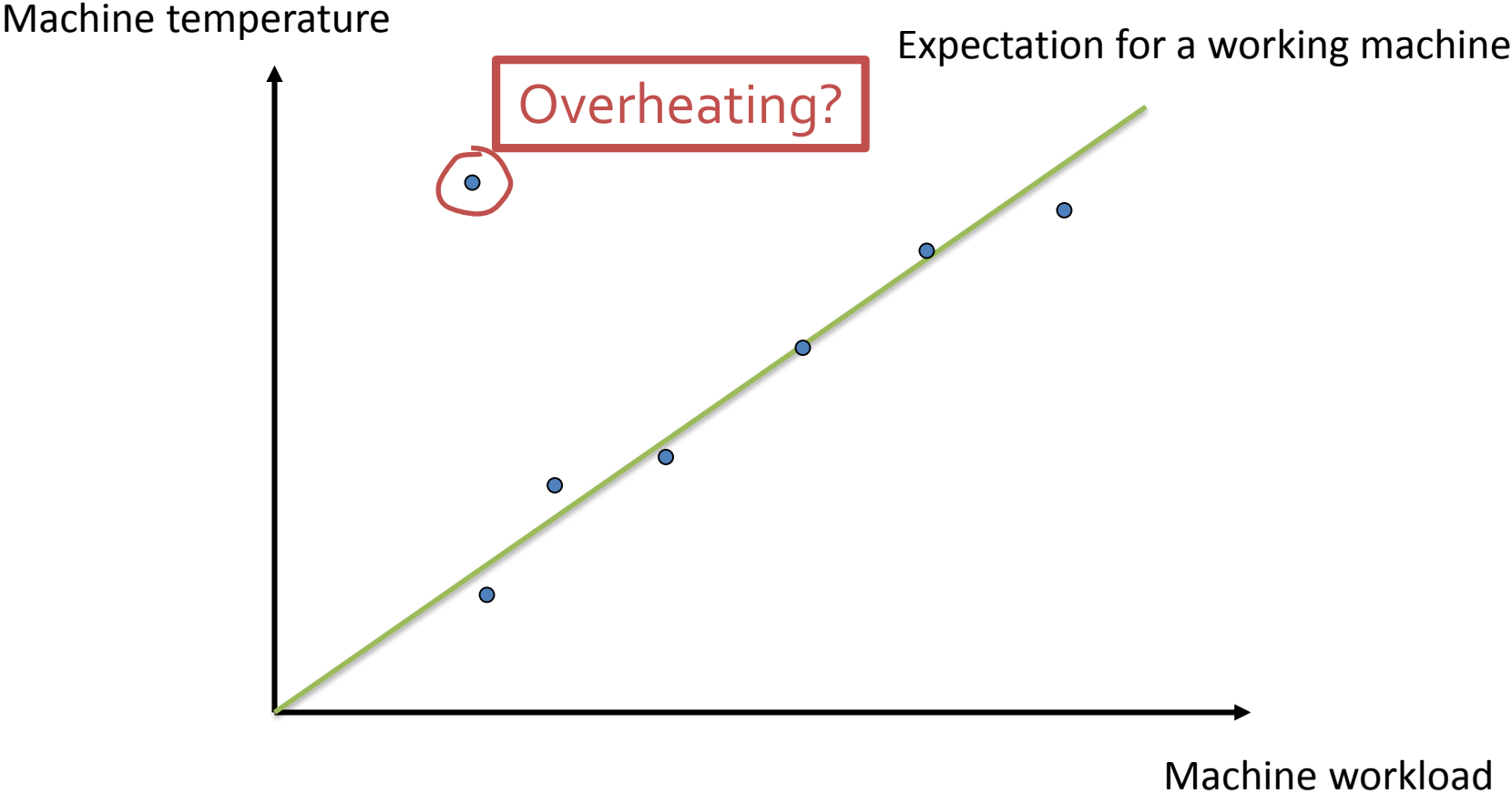
Regression



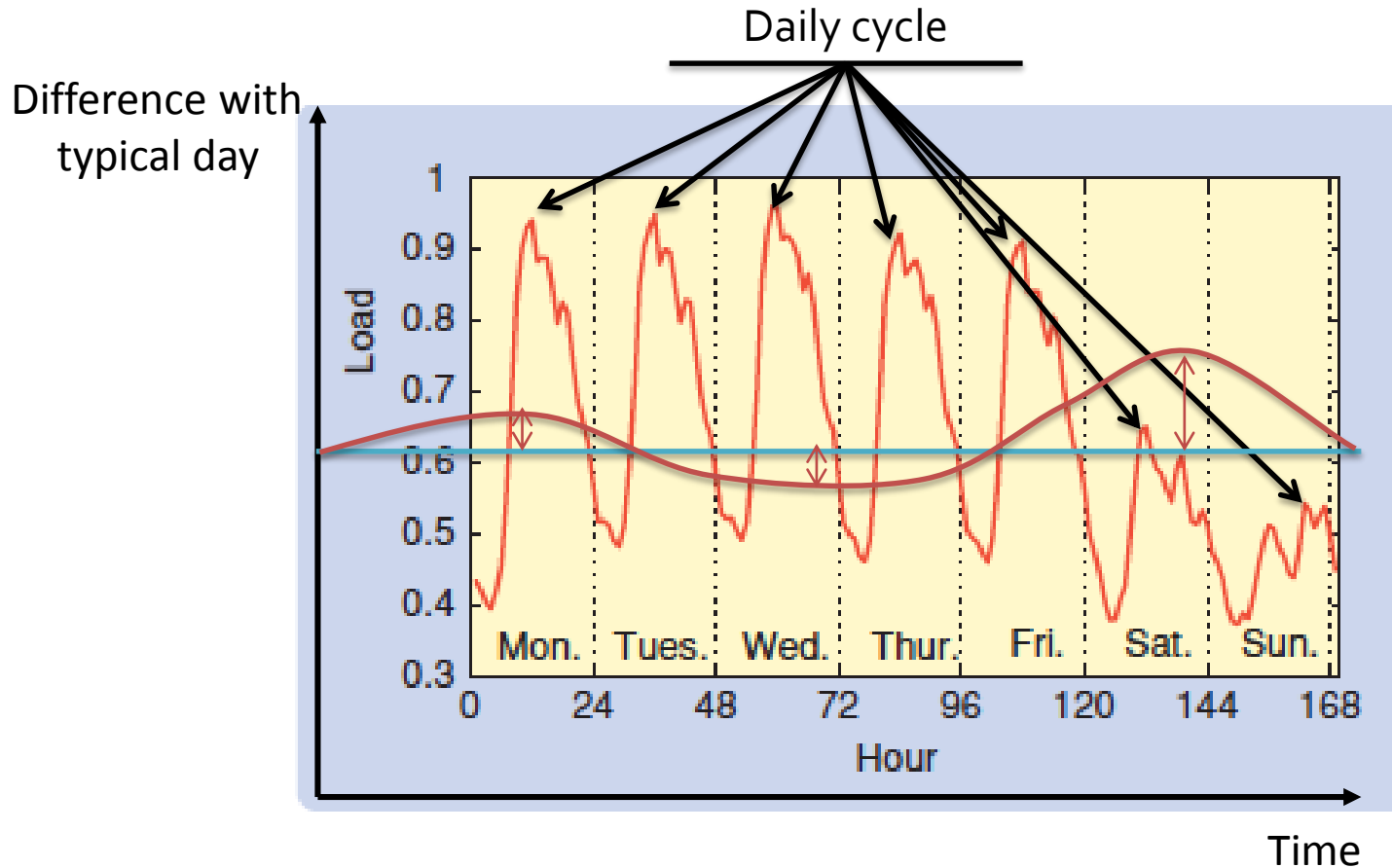
Clustering/Classification



Outlier Detection



Normalization



Look at divergence from
⇒ "normal" behaviour

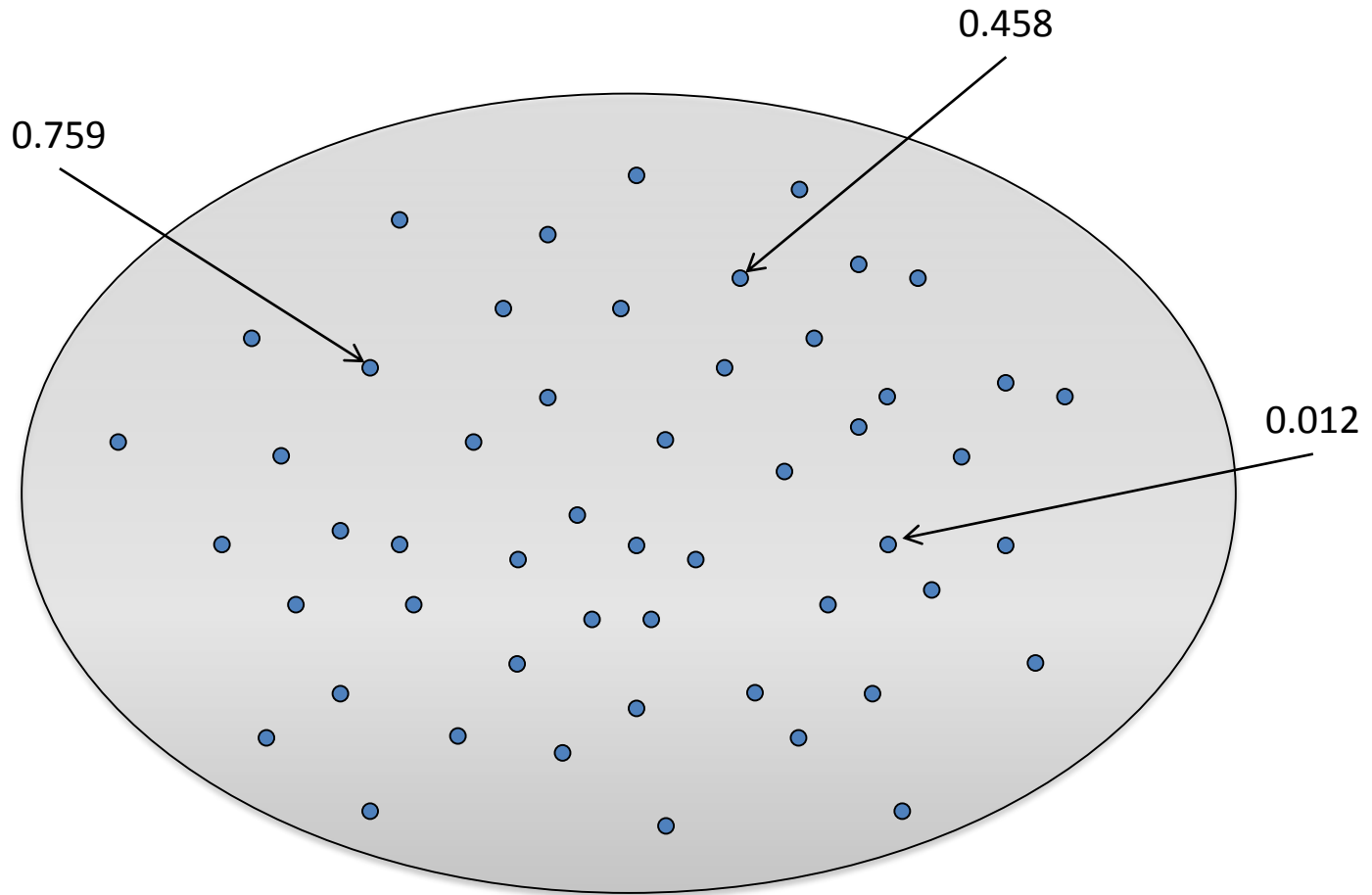
Ranking

A search input field with a light grey border. The text "Query a database" is centered within the field. On the right side of the field, there is a small microphone icon.

Google Search

I'm Feeling Lucky

Ranking



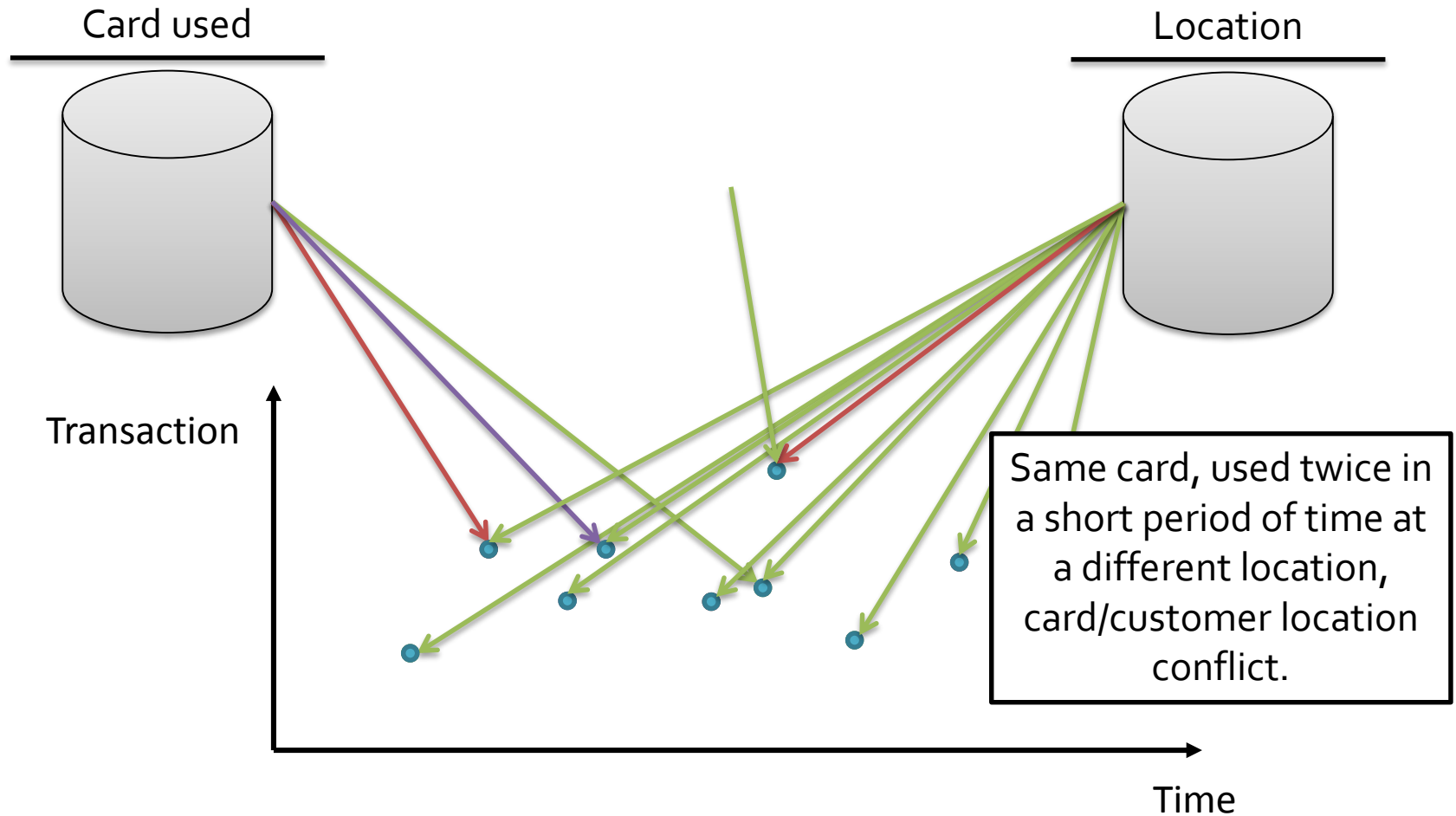
Score each point as a possible answer

Ranking



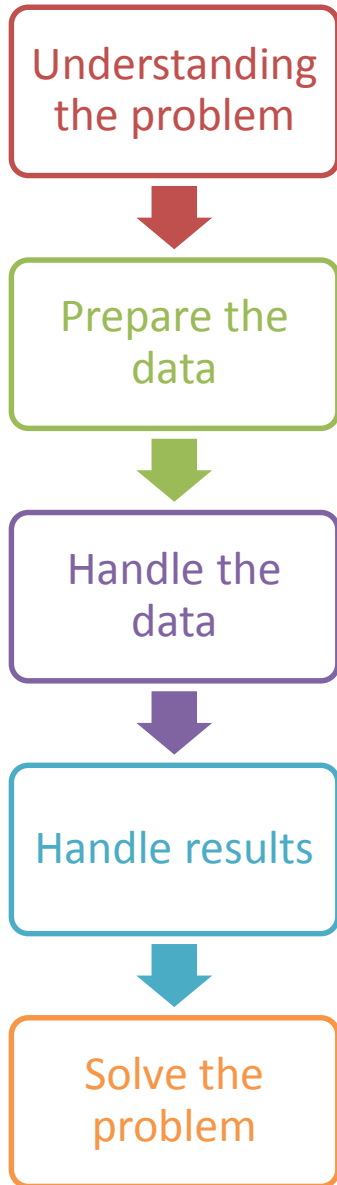
Rank all points from highest to lowest score and return this list

Data fusion

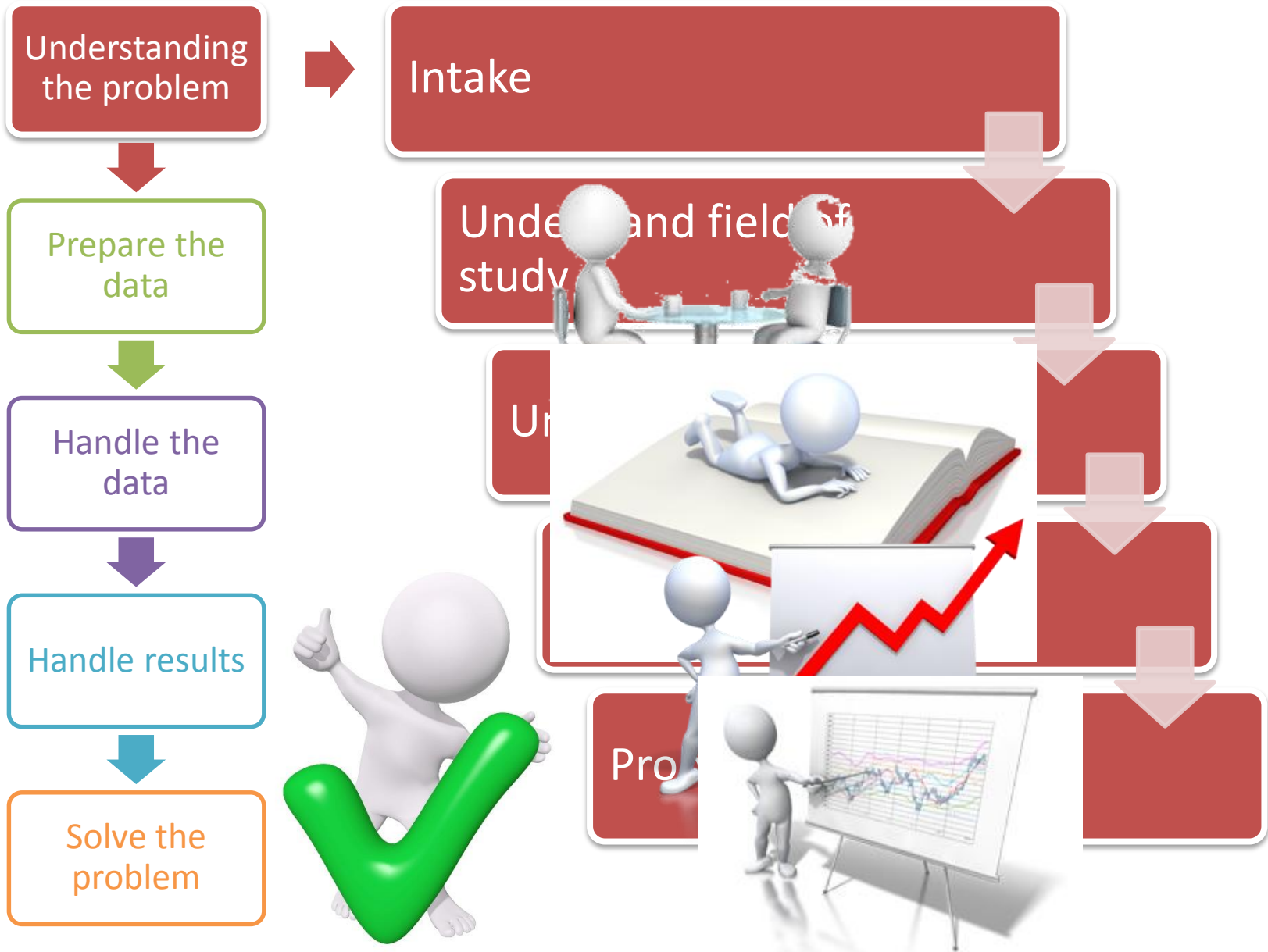


At first sight, **Card theft?** looks like normal behaviour

Taming a Data Project



Taming a Data Project



Taming a Data Project

Understanding the problem



Prepare the data



Handle the data

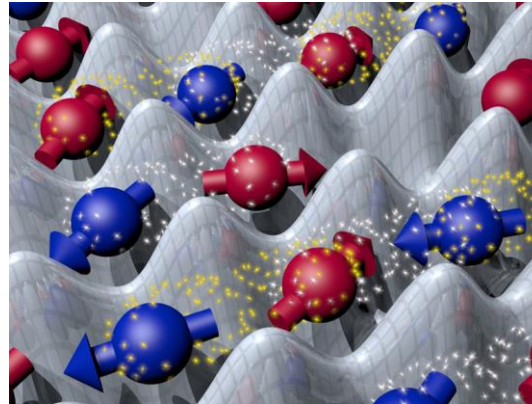


Handle results



Solve the problem

Change this

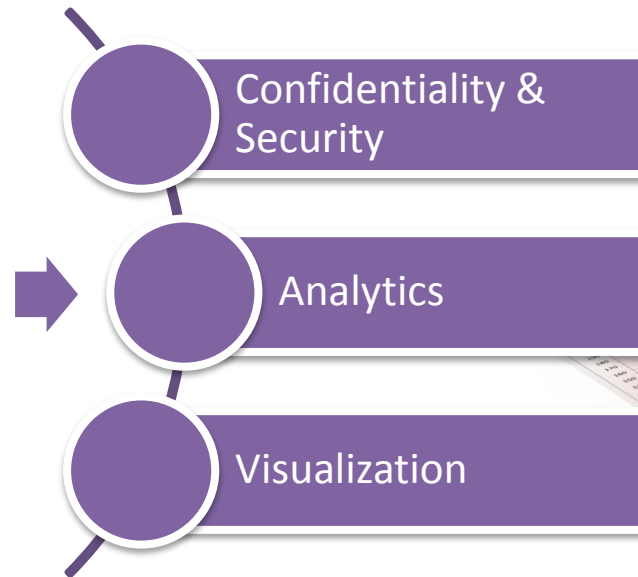
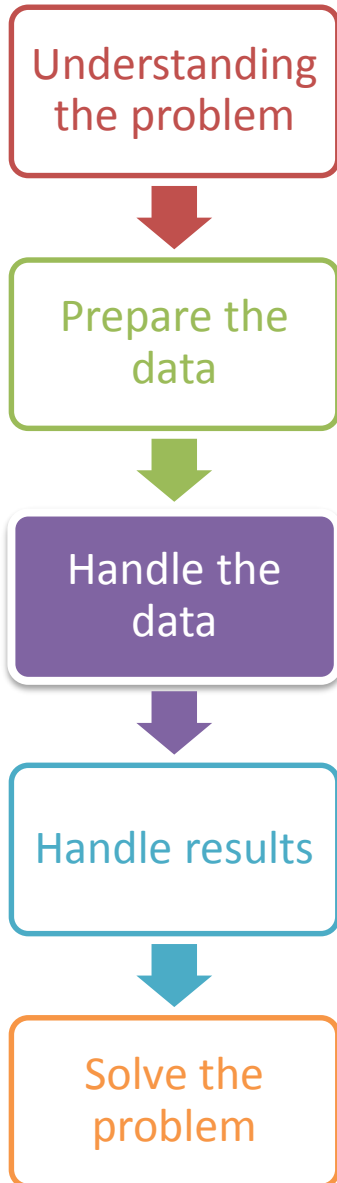


Into this

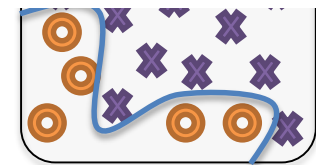
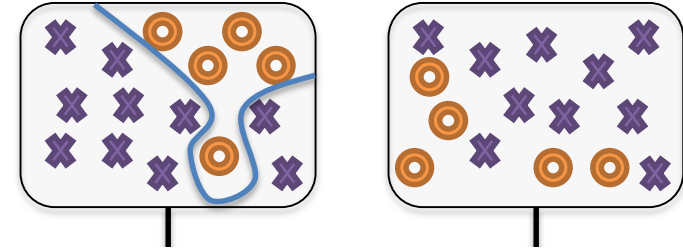
Dim	Attribute	K^{11}	K^{12}	K^{22}
1	Angle (degree)	0.0	42.9	0.0
2	Vertex distance (Å)	0.0	9.7	0.0
3	Square-root of area	4.0	4.5	5.0
4	Aspect ratio	1.0	0.8	1.0
5	Mean distance (Å)	4.2	14.0	3.9
6	Std. Dev. of distance (Å)	3.2	3.8	2.1
7	Contact pattern type	3	5	1

Data
Input module

Taming a Data Project



Is there enough and adequate data?



Taming a Data Project

Understanding the problem



Prepare the data



Handle the data



Handle results

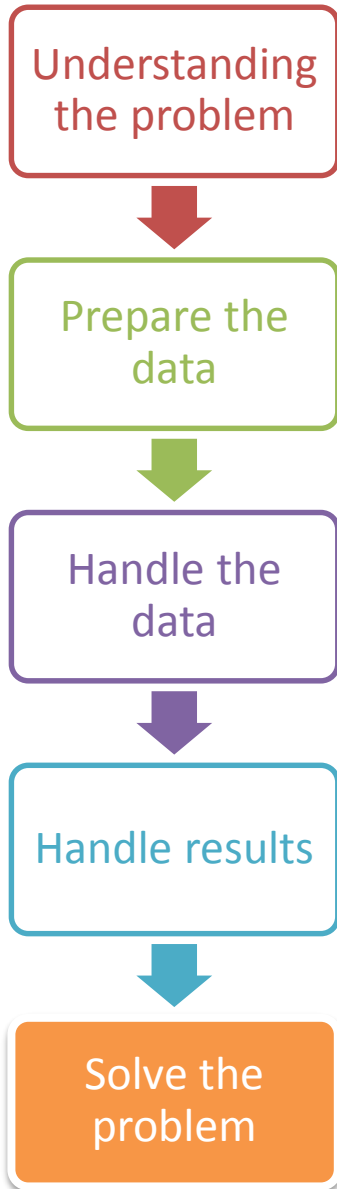


Solve the problem

Corrective objectives



Taming a Data Project



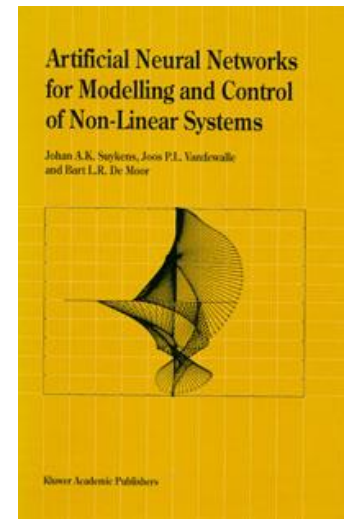
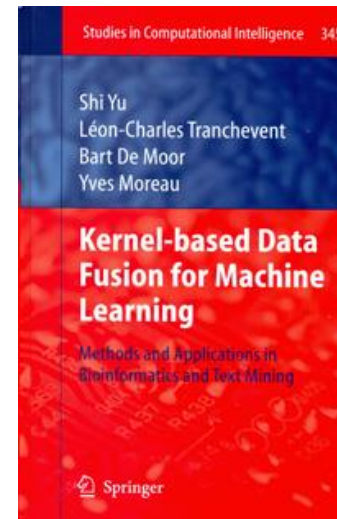
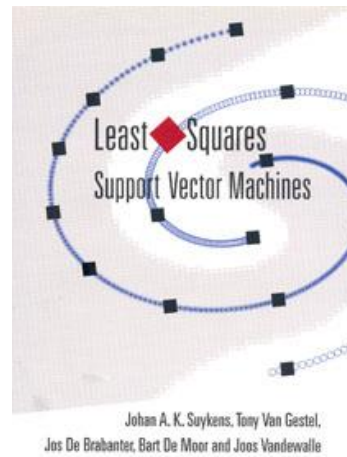
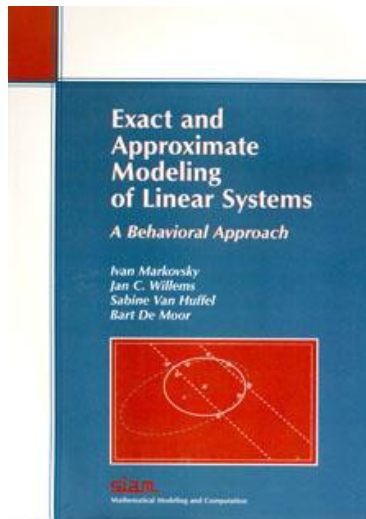
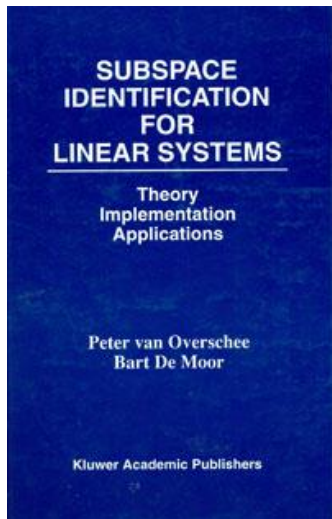
Stadius Expertise

Serious Data
Serious Mining

Energy
Industry
Environment
Social networks
Fraud and predictive analysis
Health
...

Real Quantifiable
Return on Investment

Stadius - Books



Stadius - Software



- EnsembleSVM
- TensorLab
- Clinical Data Miner
- Endeavour
- Beegle
- ...

Stadius - Spin-offs



IPCOS

- K.U.Leuven spin-off company, 1995
- www.ipcos.be
- specialized in modelling and control of multivariable industrial processes (chemical and power plants, oil exploration)
- CEO : Dr. Peter Van Overschee
- Awarded with :
 - Starters award Foundation King Baudouin (1995)
 - 3-annual starters award Flanders Technoland (1996)

 **BASF**
The Chemical Company



 **OCI**
NITROGEN

 **ARKEMA**

 大唐国际发电股份有限公司
BAZANG INTERNATIONAL POWER GENERATION CO., LTD.

 سابك
sabic



TOTAL



 **KOCH**
KOCH FERTILIZER, LLC

 **INEOS**
THE WORD FOR CHEMICALS

 قافكو
GAFCO
GULF ARABIAN FERTILIZER CO. LTD.

- K.U.Leuven, Spin-off company, 2005
- In silico drug discovery and screening
- CEO : Wilfried Langenaeker
- Silicos in the press:
 - 41ste Van Cauterenleerstoel : Maken ingenieurs het verschil ?
 - Silicos expands biology lab and announces move to new facilities (20/07/2009)
 - Silicos, a computational drug design biotech, appoints Dr. Jack Elands as Chief Business Officer (29/06/2009)

- K.U.Leuven spin-off company, 2002

- www.tmleuven.be

- Traffic and mobility consulting

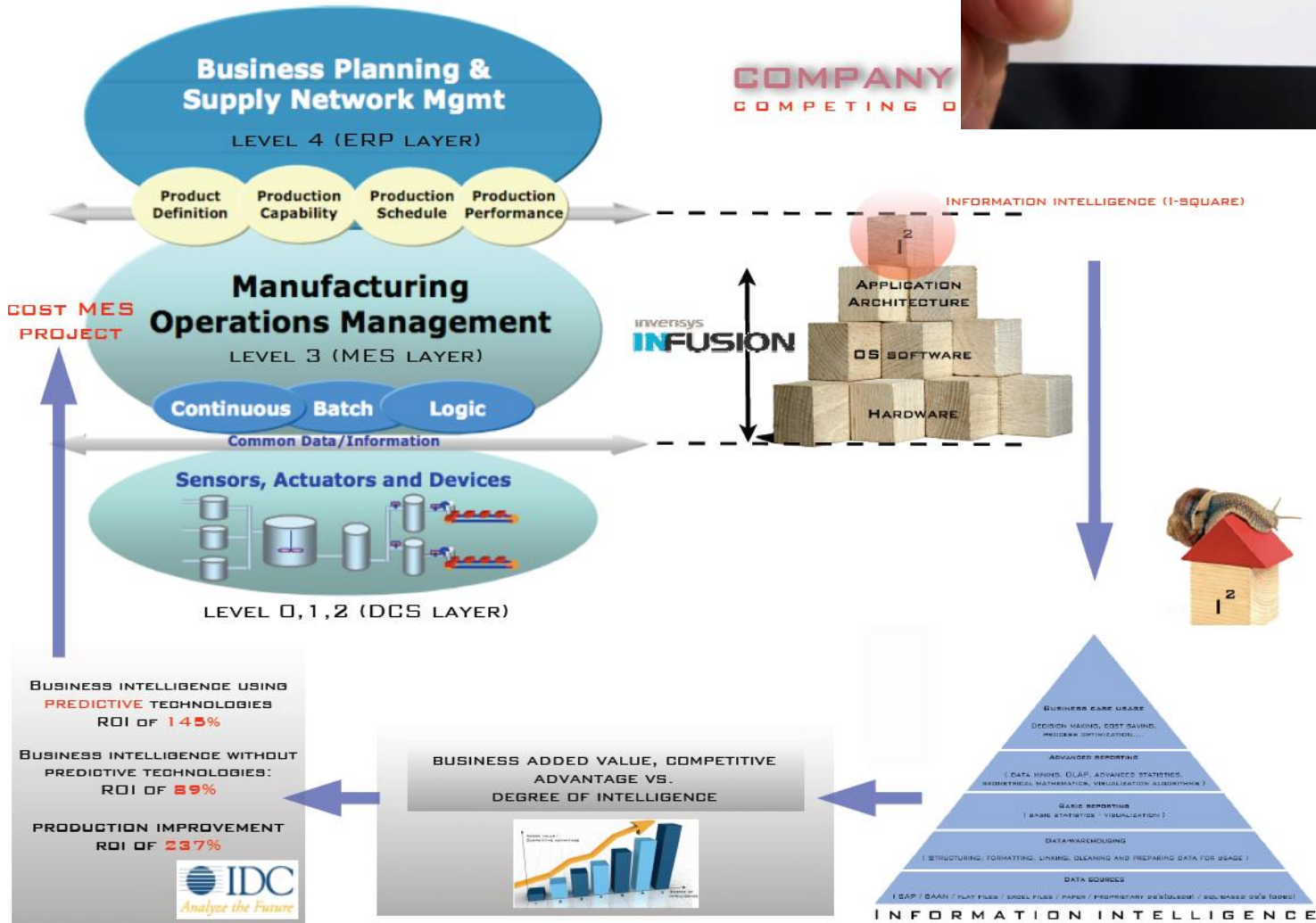
- CEO : Griet De Ceuster



- TML in the press:
 - Vlaamse regering houdt meccano-tracé buitenspel
 - Meccano wil alsnog meespelen
 - Knack : Meccano superieur aan BAM-tracé
 - De Standaard : Ik mijd dure spits in de stad
 - De Tijd : Slimme kilometerheffing is mogelijk
 - Rekeningrijden : 30 % minder files
 - De Standaard : Verdubbeling Brusselse Ring kost miljard
 - AHA-wetenstappen : Ontdek Leuven via de AHA ! wetenstappen route
 - Wiskundige modellen in het verkeer



COMPANY
COMPETING O



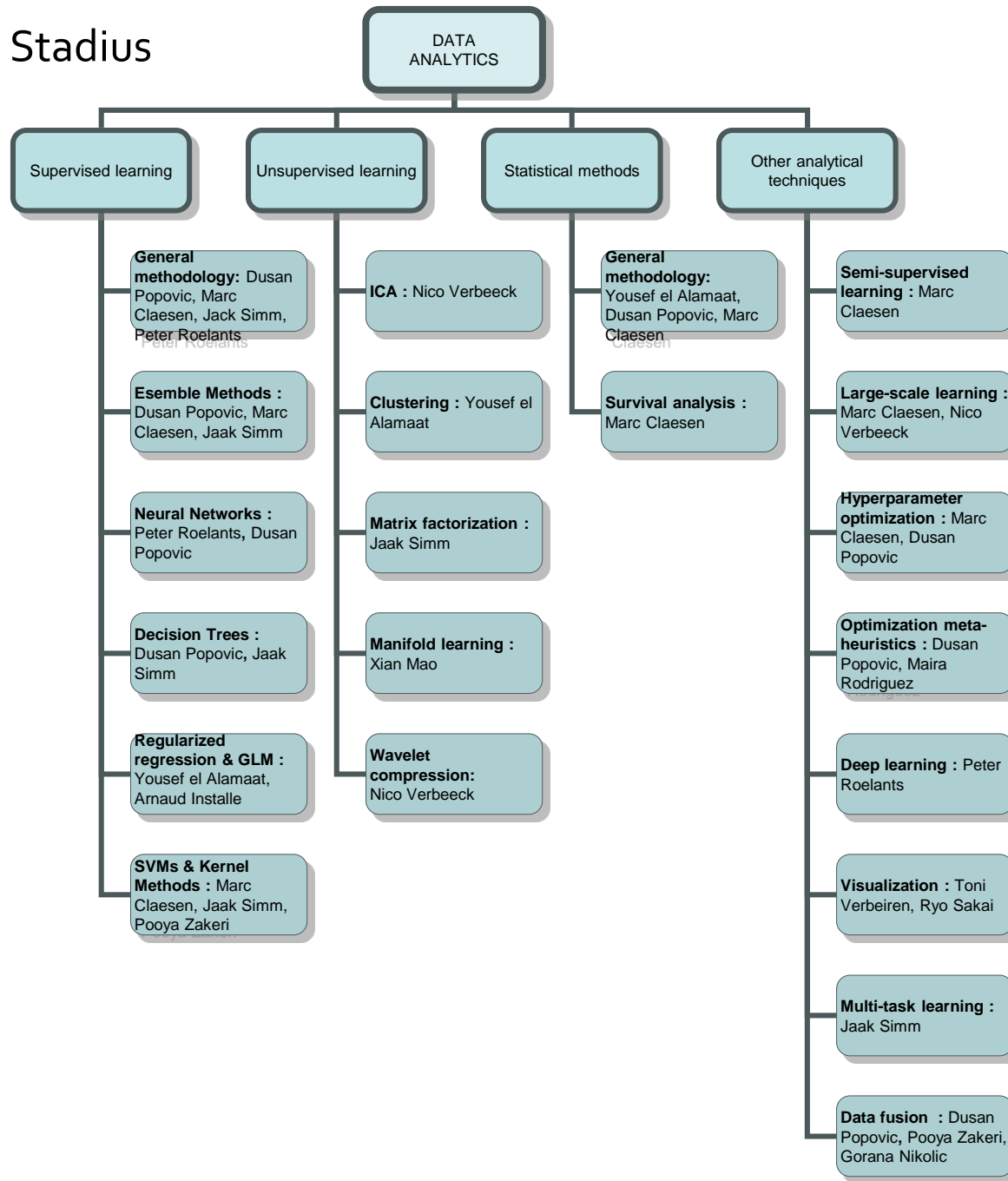


- K.U.Leuven spin-off company, 2000
- www.norkom.com (acquired by Norkom Technologies, 2006)
- Fraud detection in telecom and finance
- CEO : Steven Verhoeven
- Data4s in the press:
 - Campuskrant : "Gemma Frisiusfonds voor prille spin-offs bestaat 10 jaar" (Dec. 2007)
 - LeuvenInc : Norkom technologies acquires European risk management firm data4s (Nov. 2004)
 - Norkom Technologies neemt Europese Risk Management firma DATA4s over (Nov. 2004)
 - Norkom Technologies acquires European risk management firm DATA4s (Nov. 2004)
 - Norkom and DATA4s join forces (Nov. 2004)
 - KBC bank selecteert anti-witwastechologie van data4s (Sep. 2004)
 - LeuvenInc : DATA4s secures 1.75 million euro in second financing round (Feb. 2002)
 - Rabobank bestrijdt witwassen met DATA4s technologie
 - Clear2pay integrates data4s's fraud money laundering prevention platform
 - Data4s, tussen data-mining en Rode Duivels (Campuskrant, Jun. 2000)
 - Technoparels



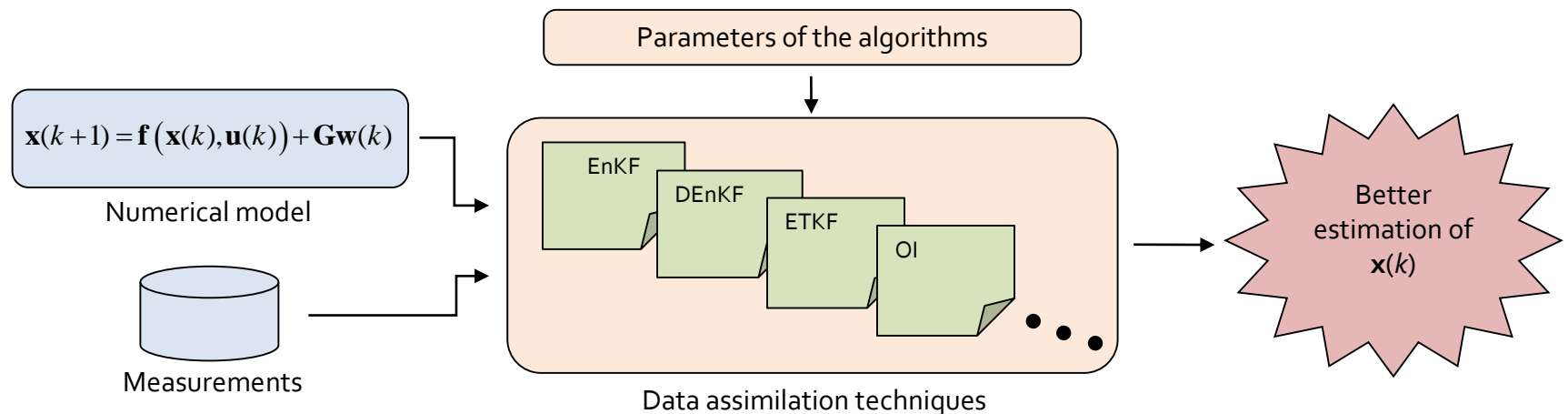
- K.U.Leuven Spin-off Company, 2008
- www.cartagenia.com
- CEO : Herman Verrelst
- Clinical applications of genetic analysis
- Cartagenia in the press:
 - Cartagenia partners major UK study on new prenatal diagnosis technology
 - Honderdste spin-off is kroon op succesverhaal LRD (Campuskrant 29/05/2013)

Algorithms in Stadius



Data Assimilation

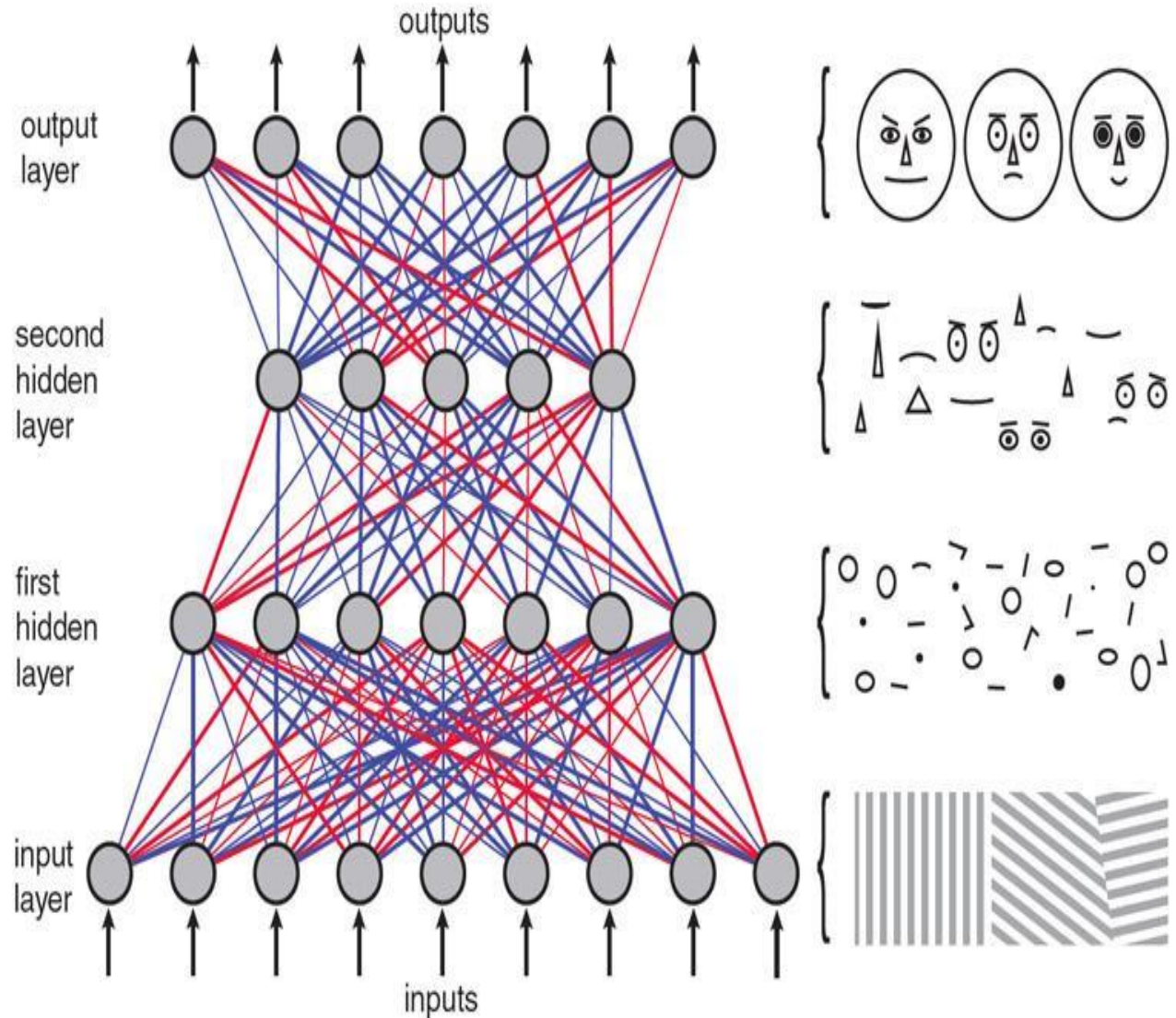
Data assimilation is the common name given to several numerical techniques that combine **the outputs of a numerical model** with **observational data** in order to improve the quality of the model predictions.



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

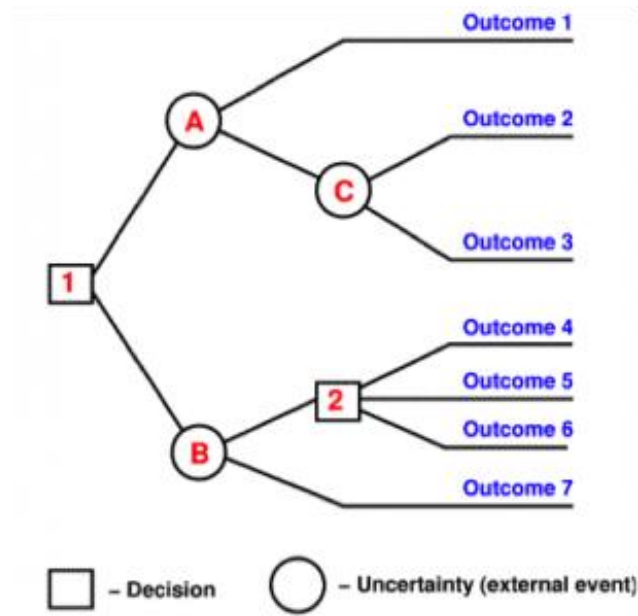
Deep Learning & Neural Networks

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



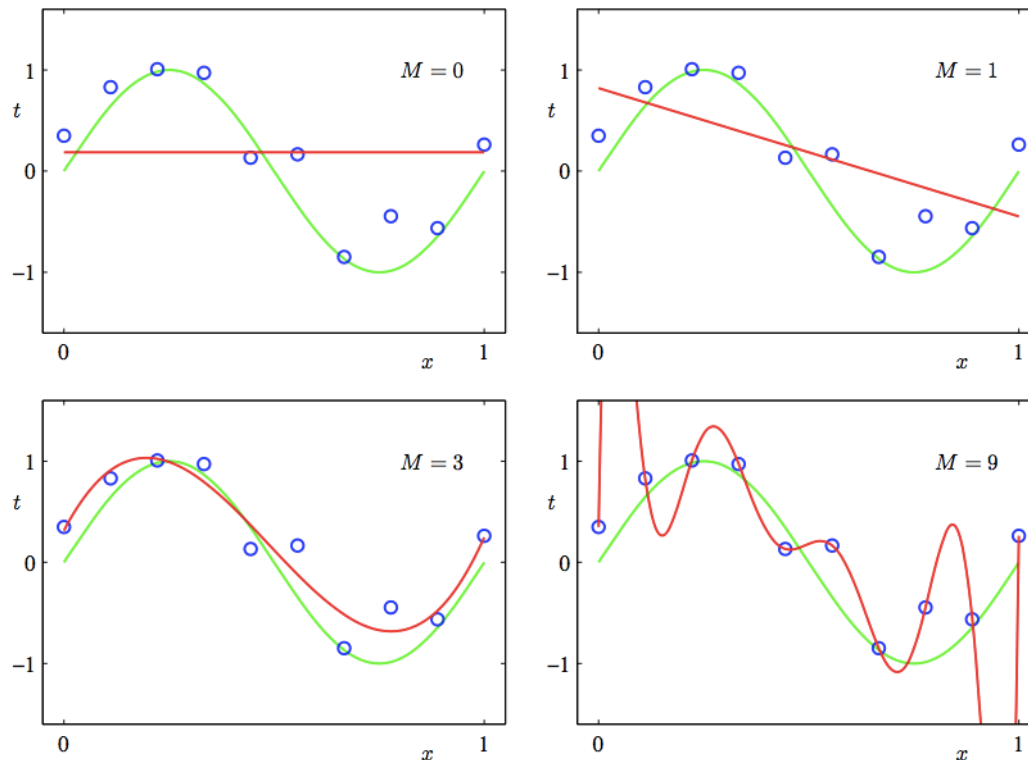
Decision trees

Decision nodes are trained according to a labeled set of data points. A new instance is given as an input and run through the tree, which then produces the most likely output.



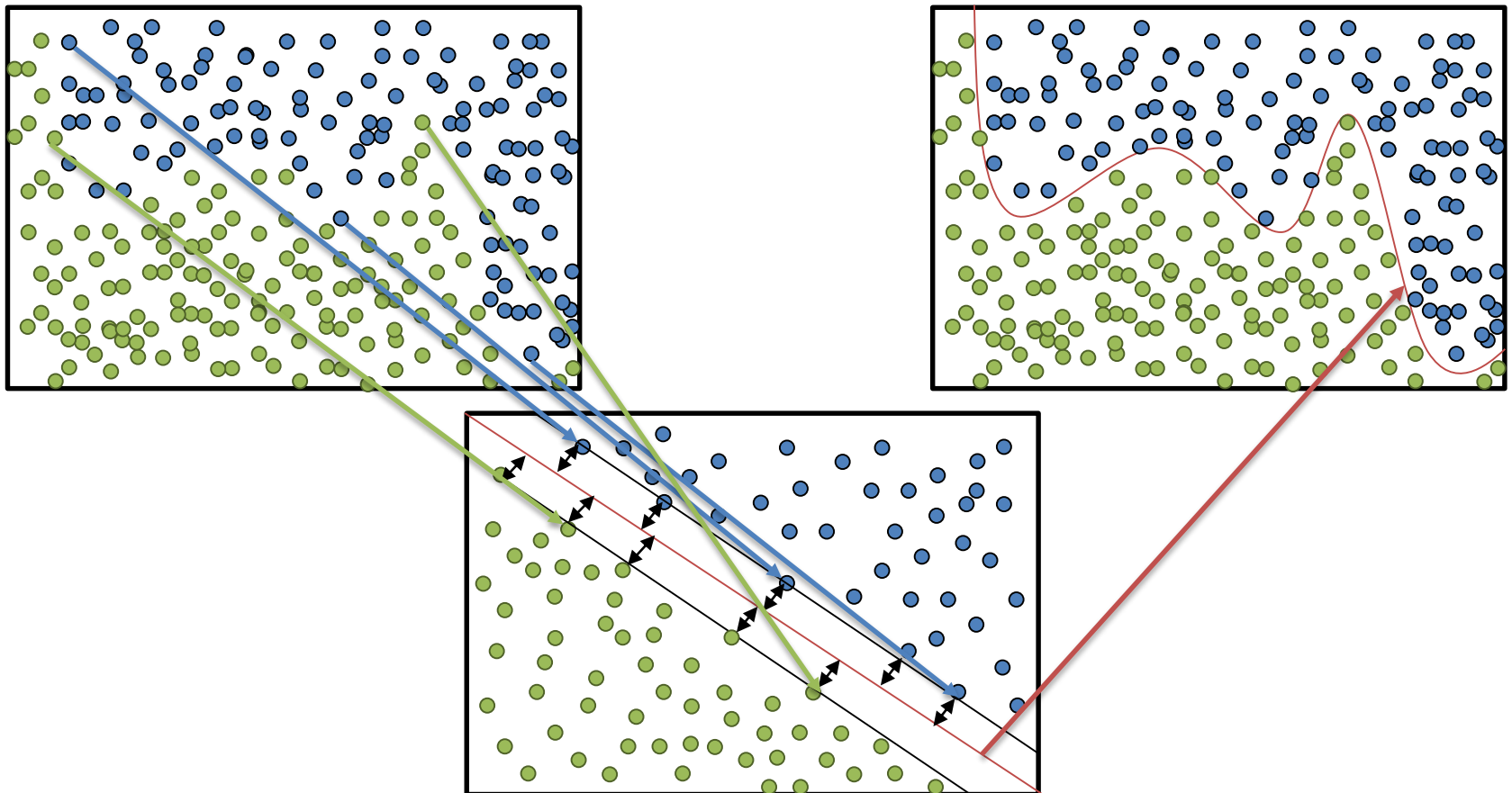
Regularized Regression

Fitting a regression function on a data set can result in overfitting: the regression fits to the data, but not to the general trend. The regression is thus not generalizable! A solution is to punish the learner for creating a model with high complexity.

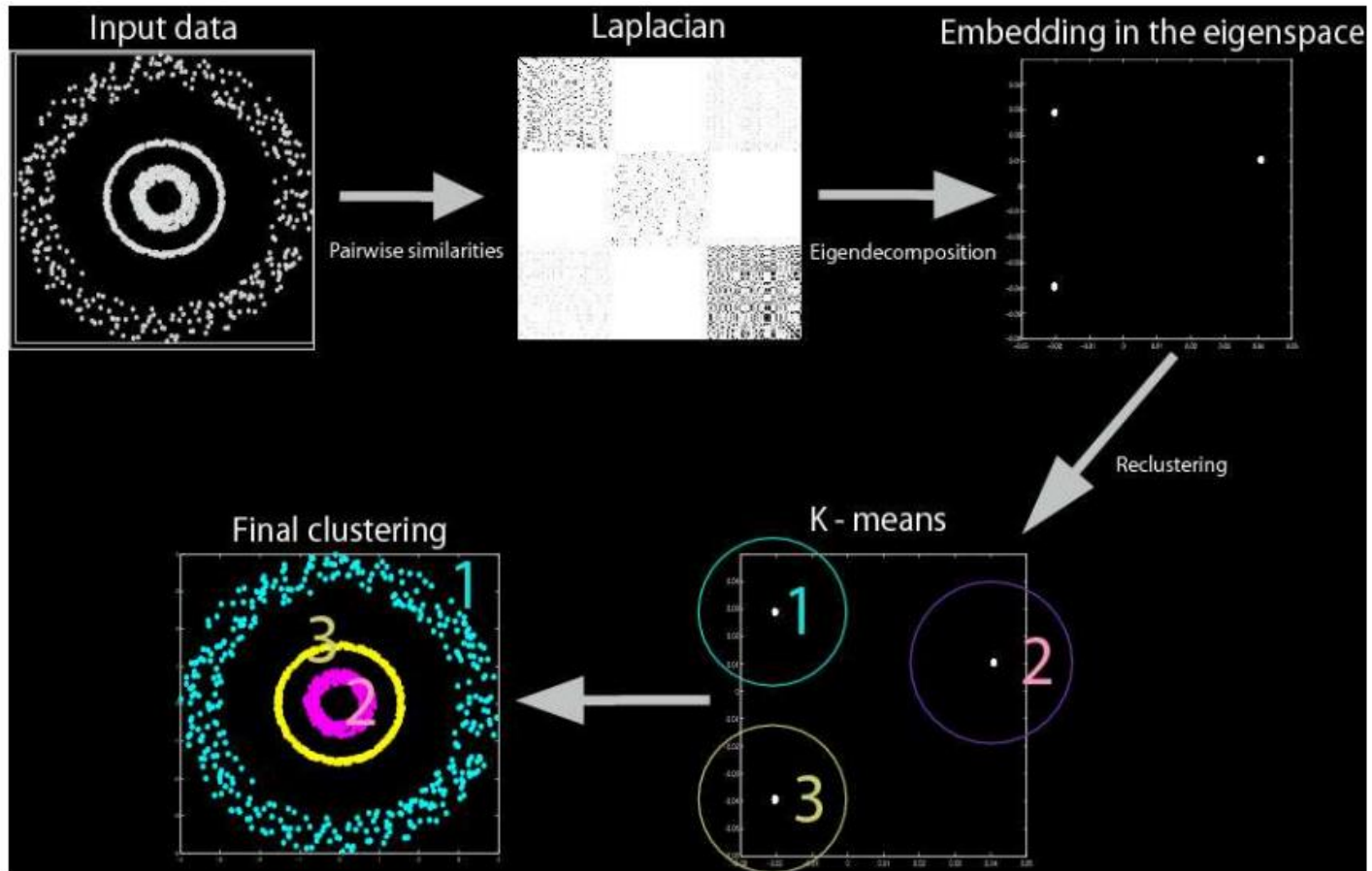


Support Vector Machine

First transform the problem to a high-dimensional form, where the solution is easily found, through the so-called 'kernel trick'. Then, transform the decision boundary back to the original form.

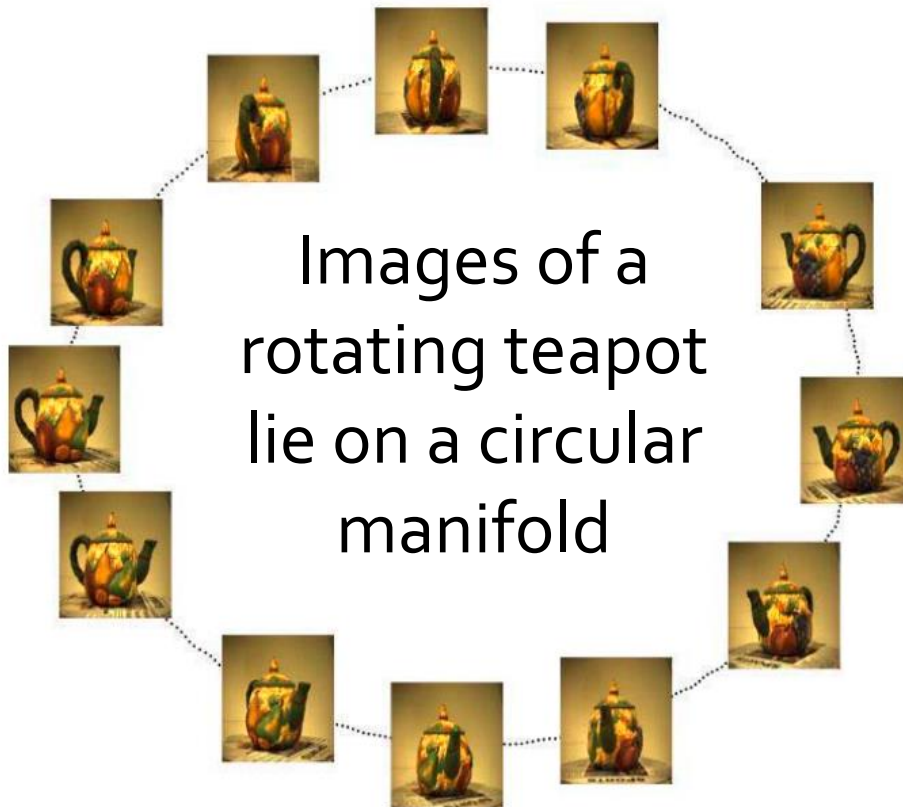
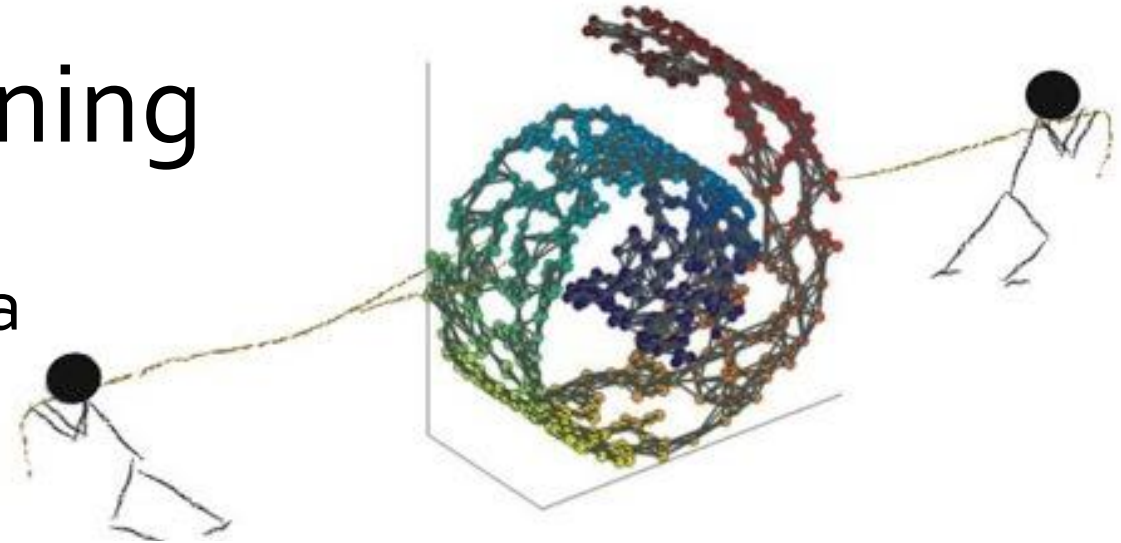


Spectral clustering



Manifold learning

A lot of datasets live on a low dimensional manifold.

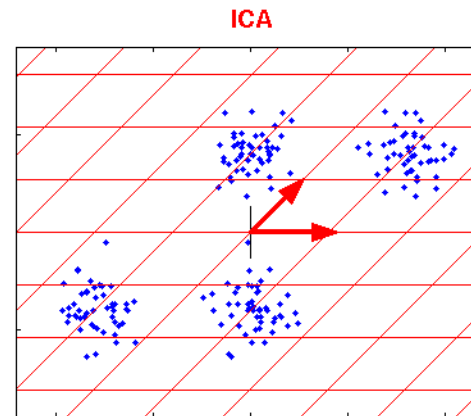
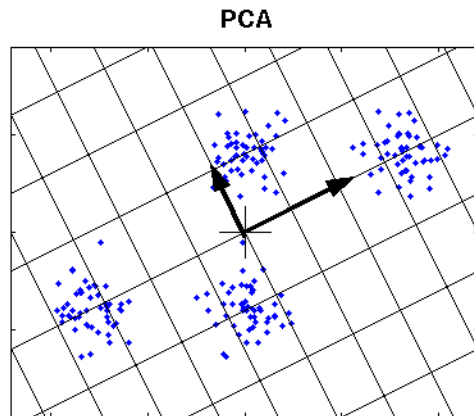


Images of a rotating teapot lie on a circular manifold

Goal: Find a low-dimensional basis for describing the high-dimensional data

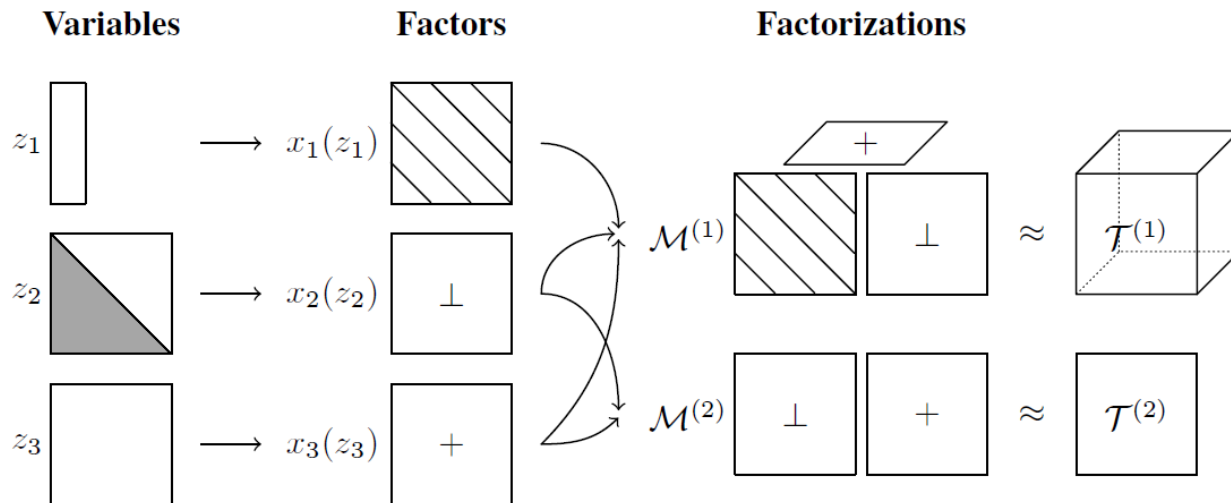
Component analysis

The data dimensionality is reduced by dividing the data set into smaller, relevant components. This can be done by maximizing the variance (principal component analysis), or by finding independent sources of data (independent component analysis).



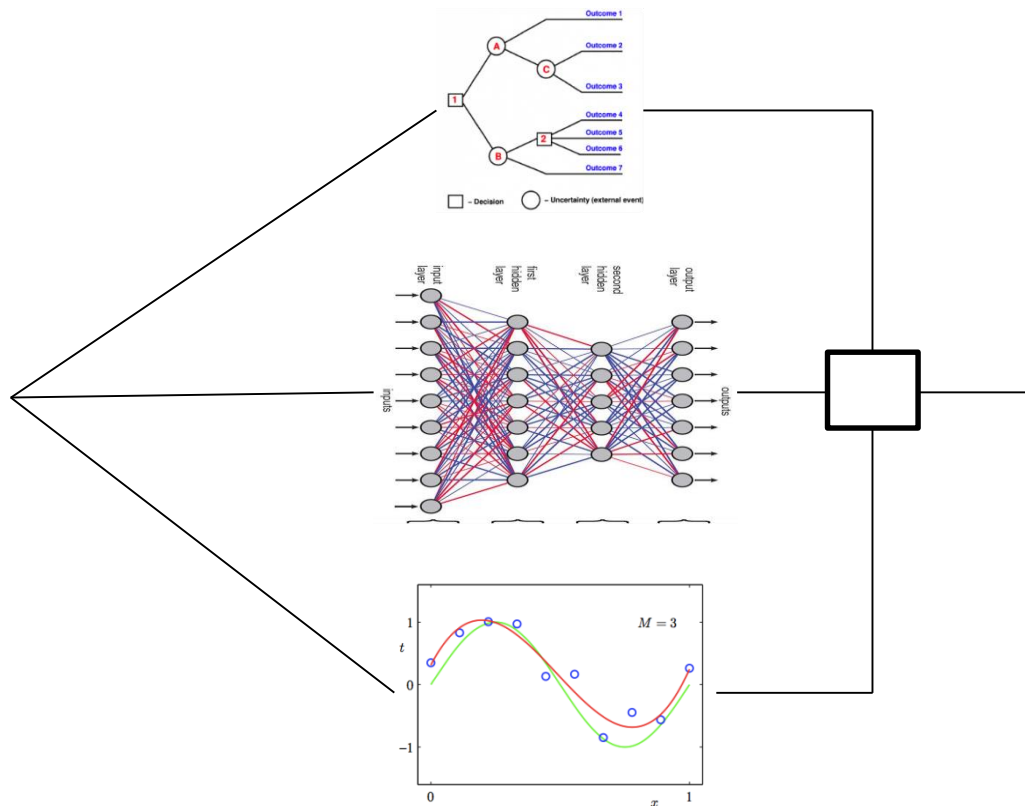
Tensor methods

The data is represented in a so-called tensor, which is a higher-dimensional extension to a vector or a list (i.e., a tensor in three dimensions is a list of lists of lists). This can be used to represent highly-dimensional data in a very concise way. We use very efficient methods to decompose these tensors, that allow us to extract common influential factors among different features of a data set. These techniques are highly scalable and work well on Big Data.



Ensemble methods

Several machine learning algorithms are implemented in parallel to each other. A decision on the outcome is then made, based on some decision rule (e.g., majority voting).





Big Data

What

Who

Six dimensions

Data

Compute Infrastructure

Storage Infrastructure

Analytics

Visualization

Security & Privacy

Machine learning as a commodity

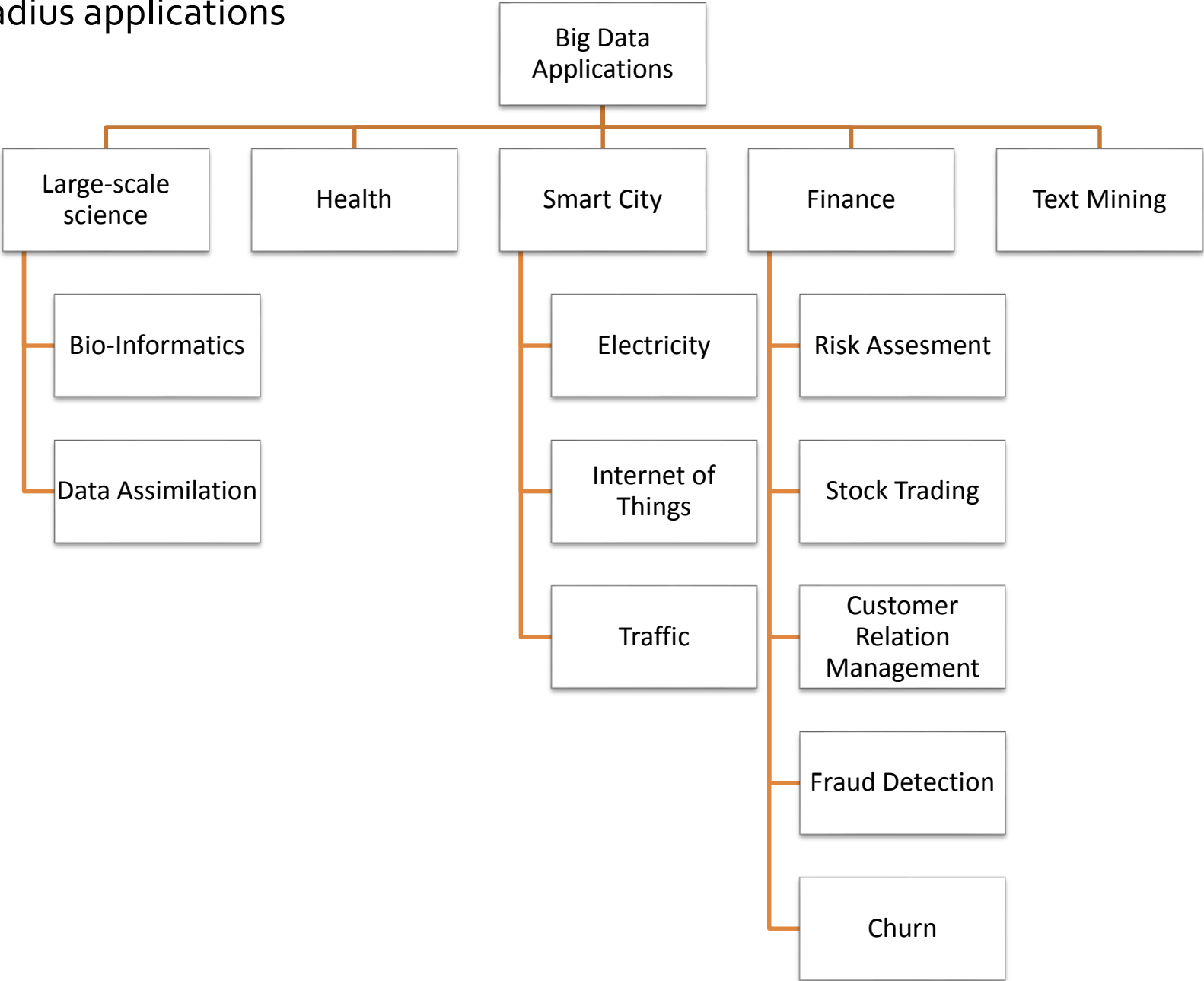
Expertise

Books & Spin-offs

Algorithms

Applications

Stadius applications



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

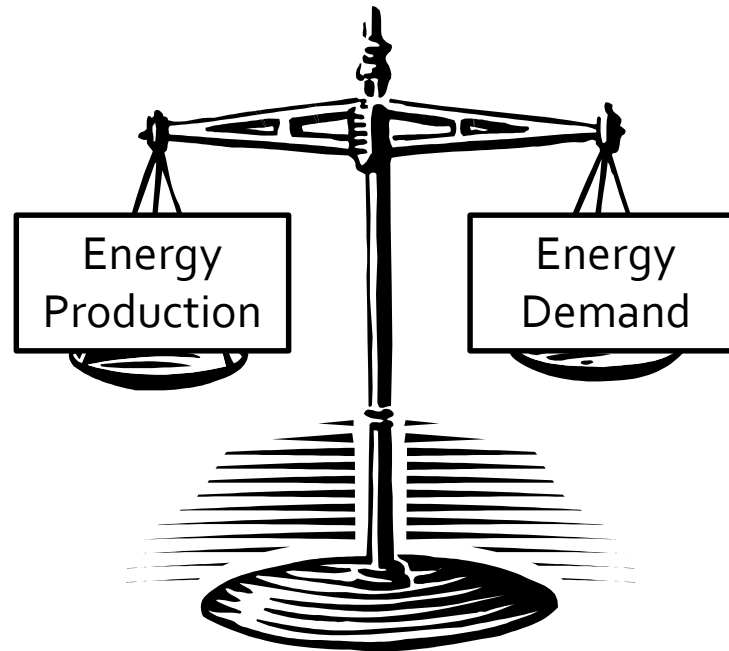
Social networks

Finance and Fraud

Health

Electric load forecasting

Problem



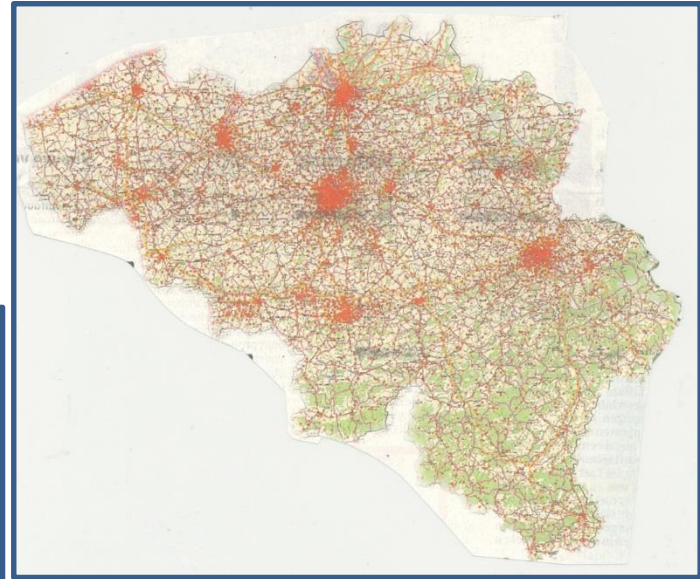
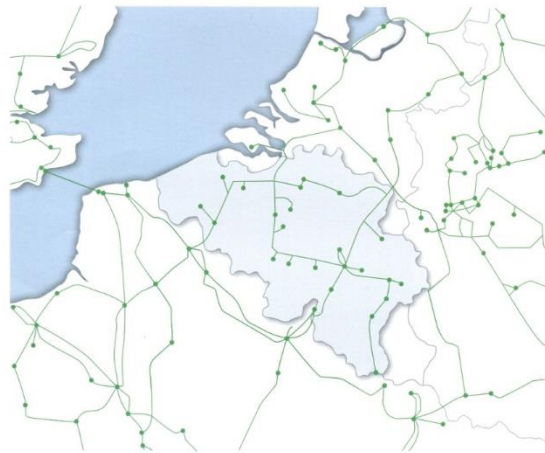
How to forecast
the demand?

Power grid

België en Europa

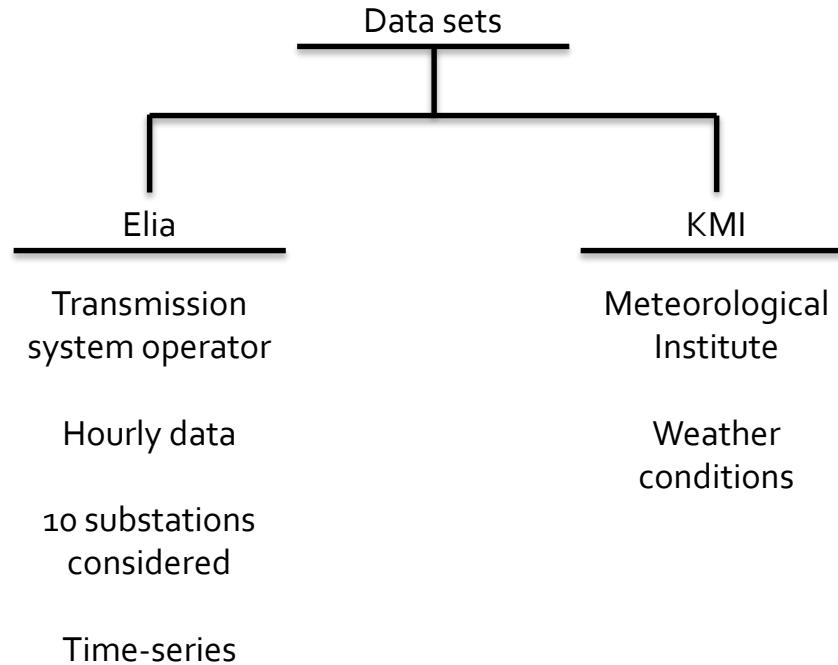
Het Elia-net:
knooppunt van elektriciteitsverkeer in Europa

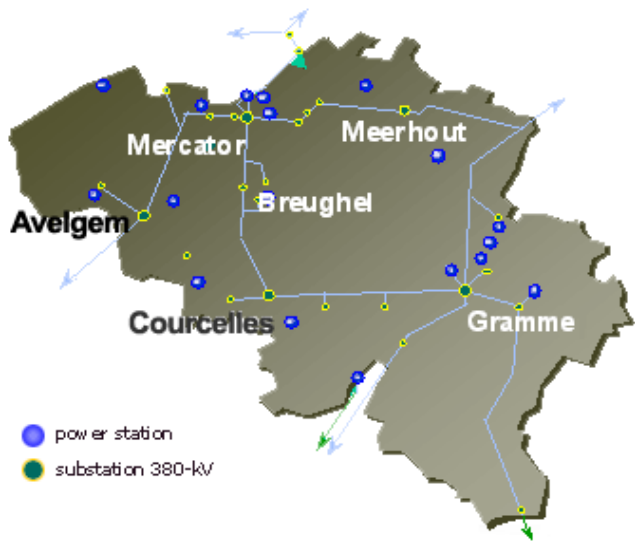
380 kV interconnectienet met hoogspanningsstations



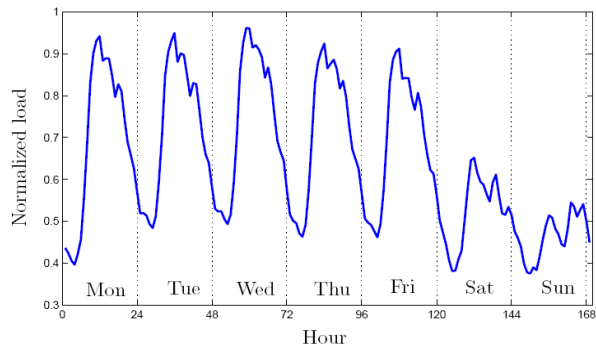
Electric load forecasting

Data

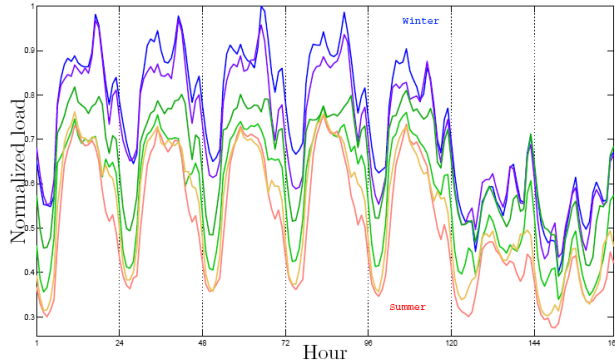




**250 transformer substations
Every 15 min, 5 years**



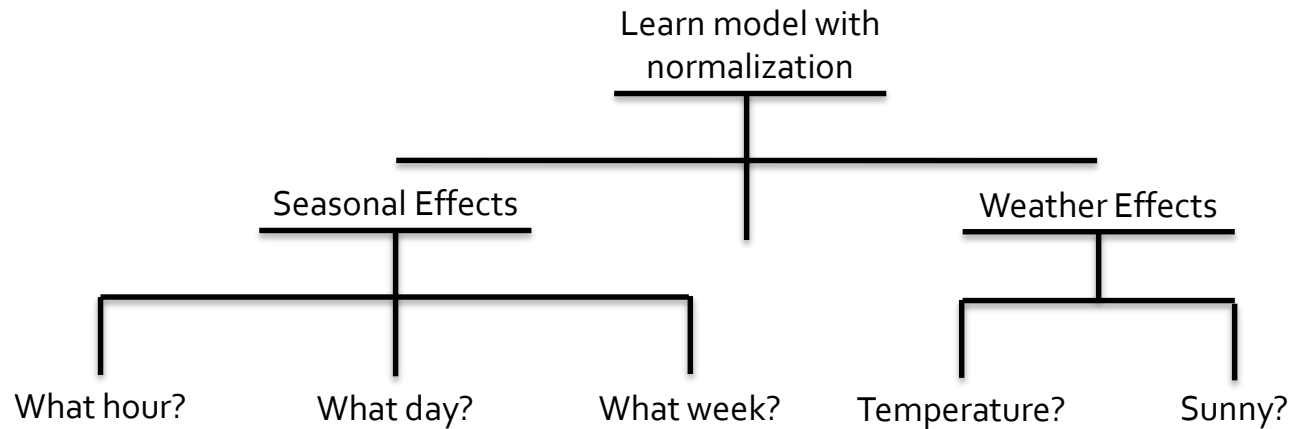
83
1 post, 1 week



1 post, four seasons

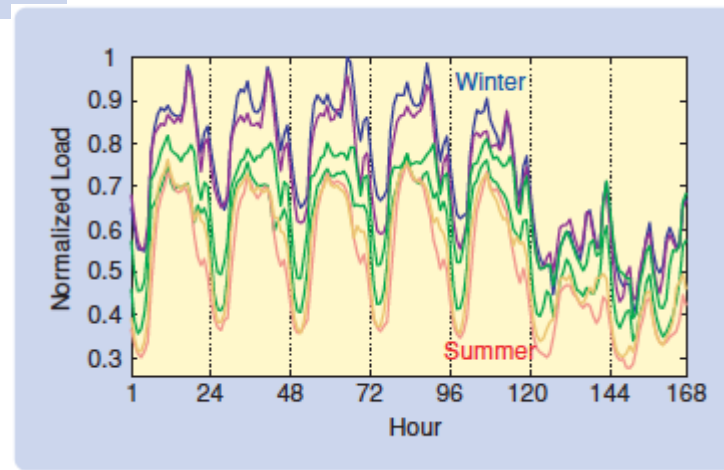
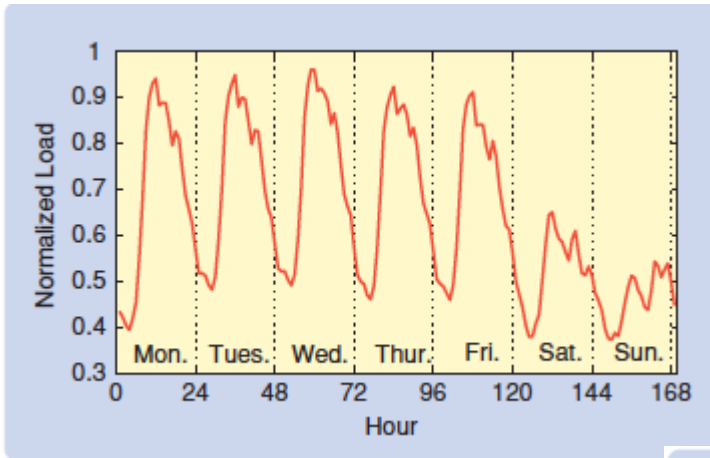
Electric load forecasting

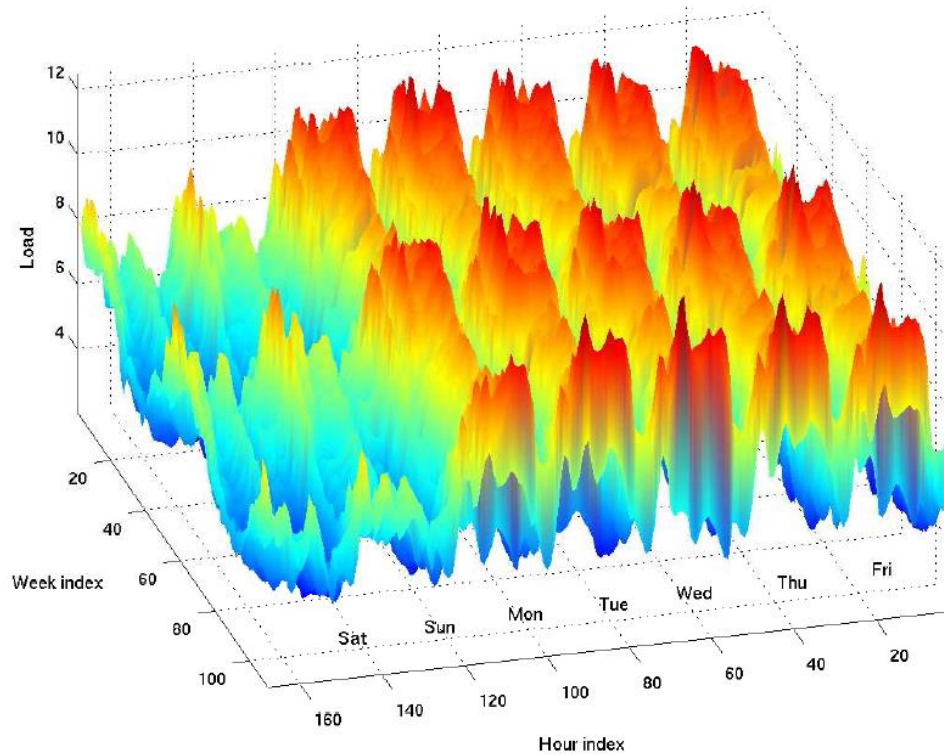
Expertise in Action



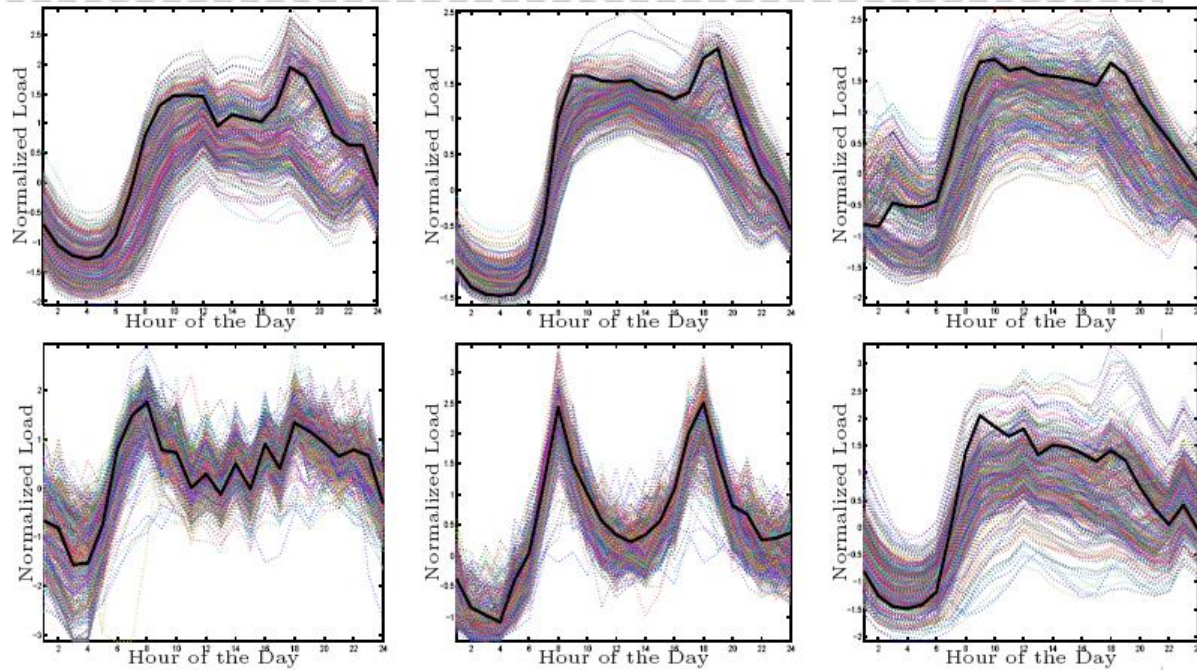
Model update: Every week!

Electric load forecasting





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



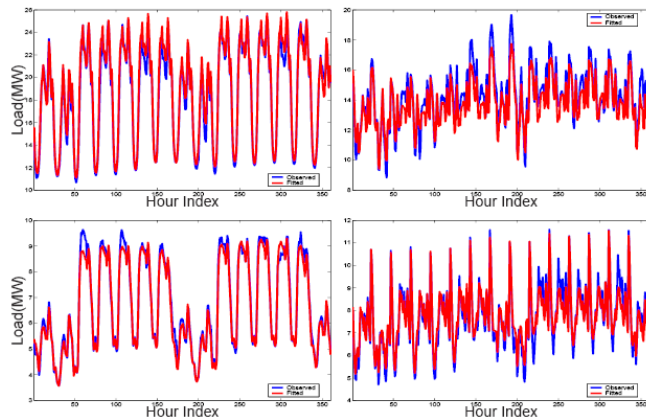
6 posts, 1 year
Seasonalities, calendar holidays !

Electric load forecasting

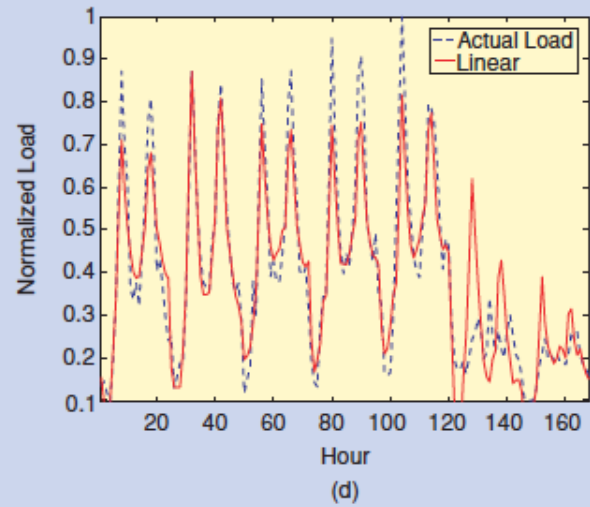
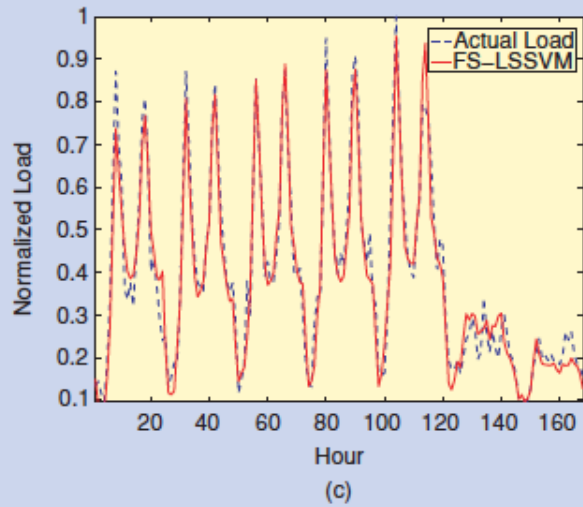
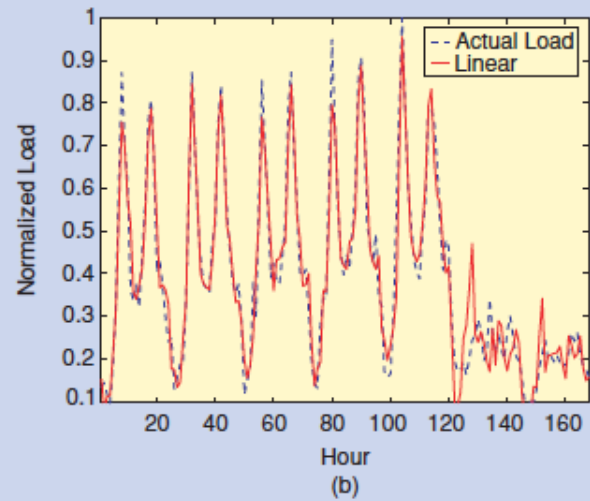
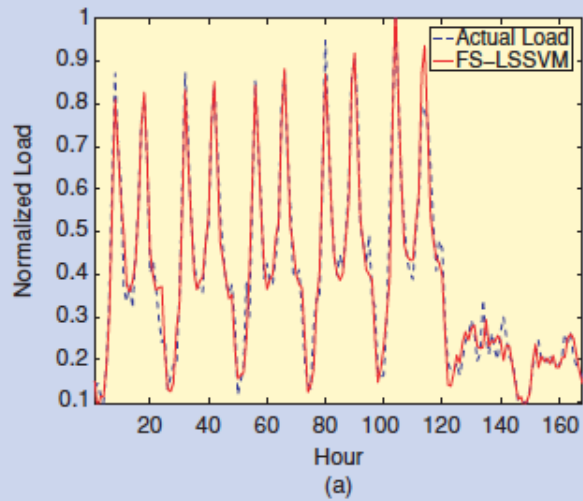
Problem Solved

- 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:



⇒ Accurate forecast



Energy

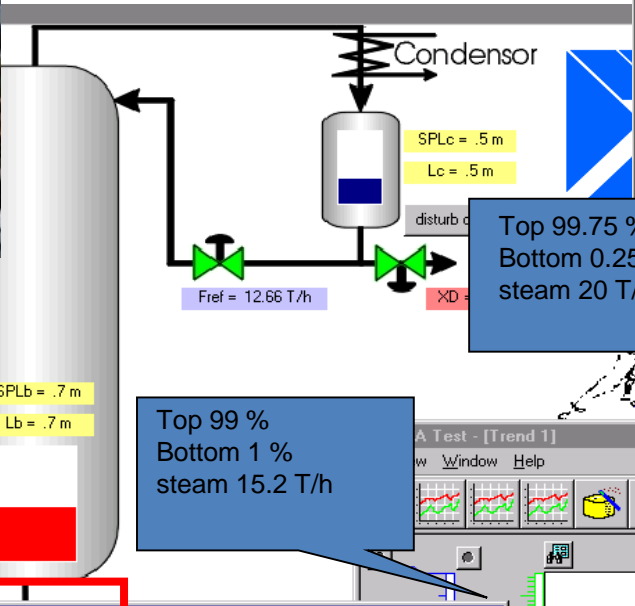
Industry

Environment

Social

Finance and Fraud

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

A Test - [Trend 1]
 View Window Help

[Time] : 3/8/0 2:58

- X_bottom.pv: 0.32
- X_destil.pv: 99.800
- F_steam.sp: 20.00

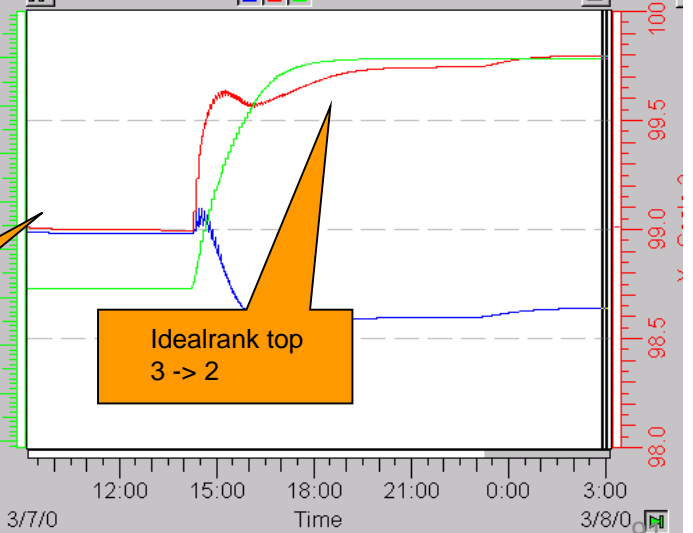
Status: Idle Test Status: Idle

INCAView - [Overview]
 View Window Help

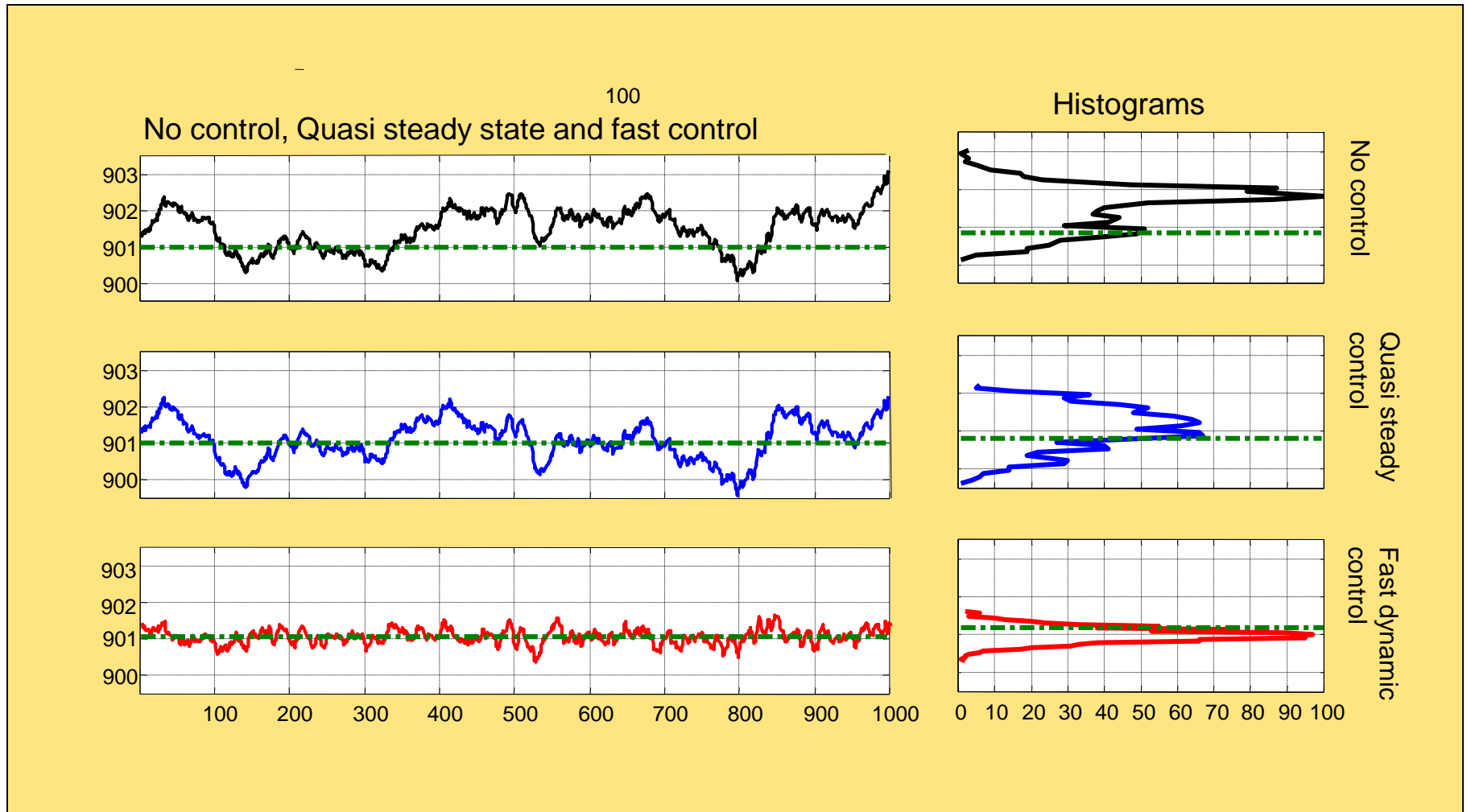
Log Status: Idle Controller Status, Reason: Turned on by operator request

All	CVs	MVs	DVs	Miscellaneous					
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"/>

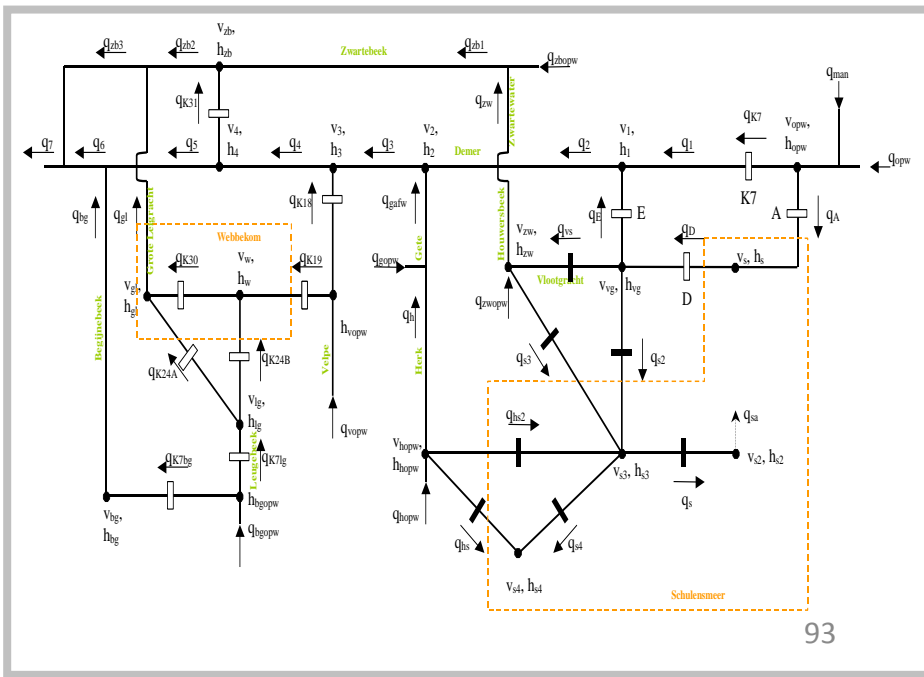
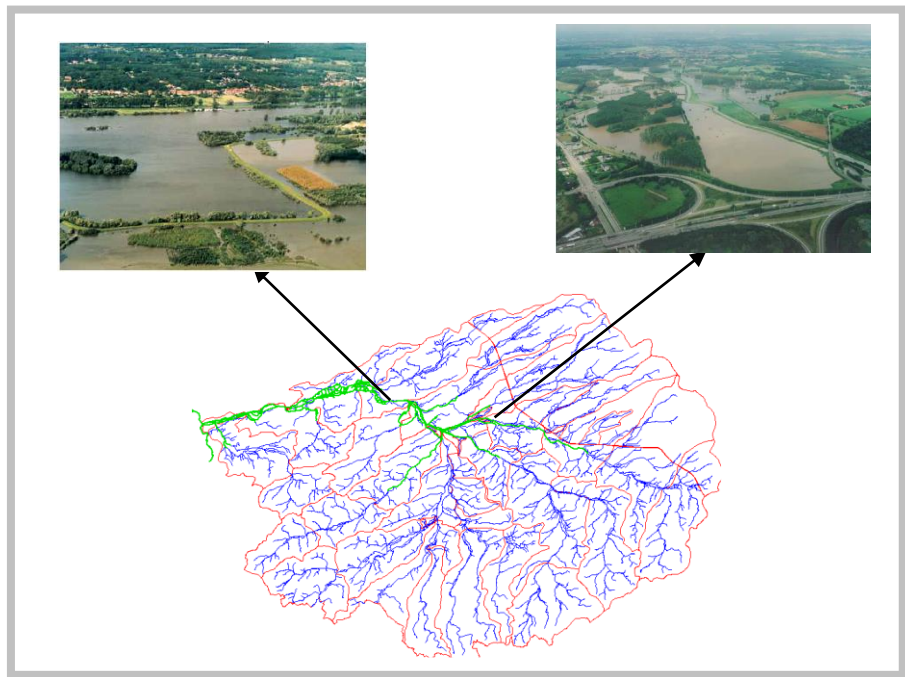
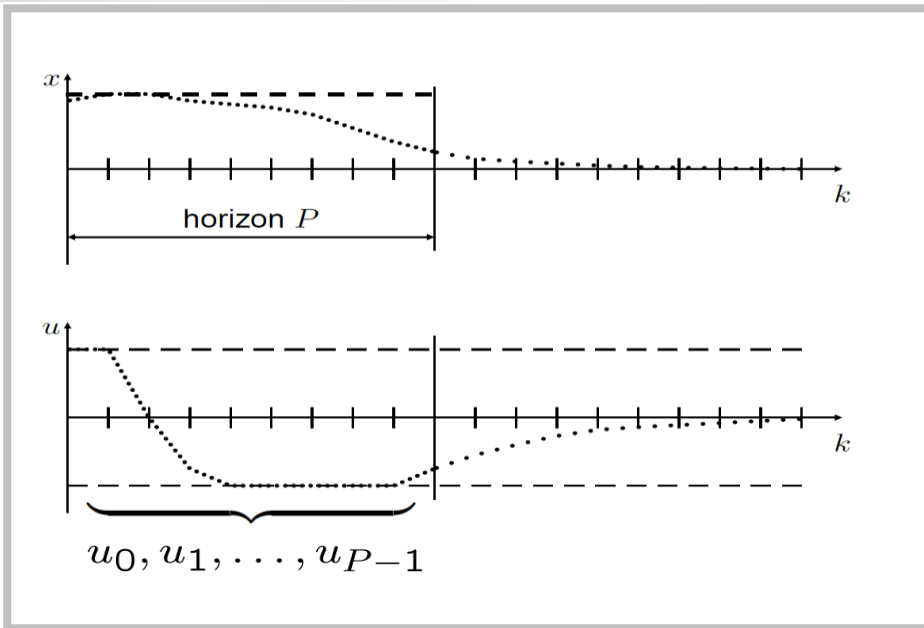
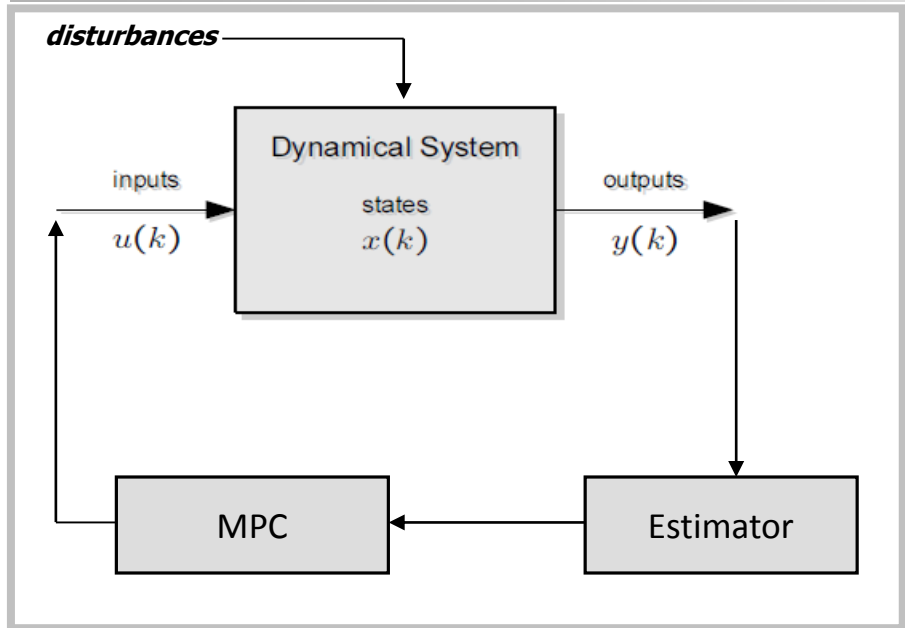
Setpoint changes



Modelling for control

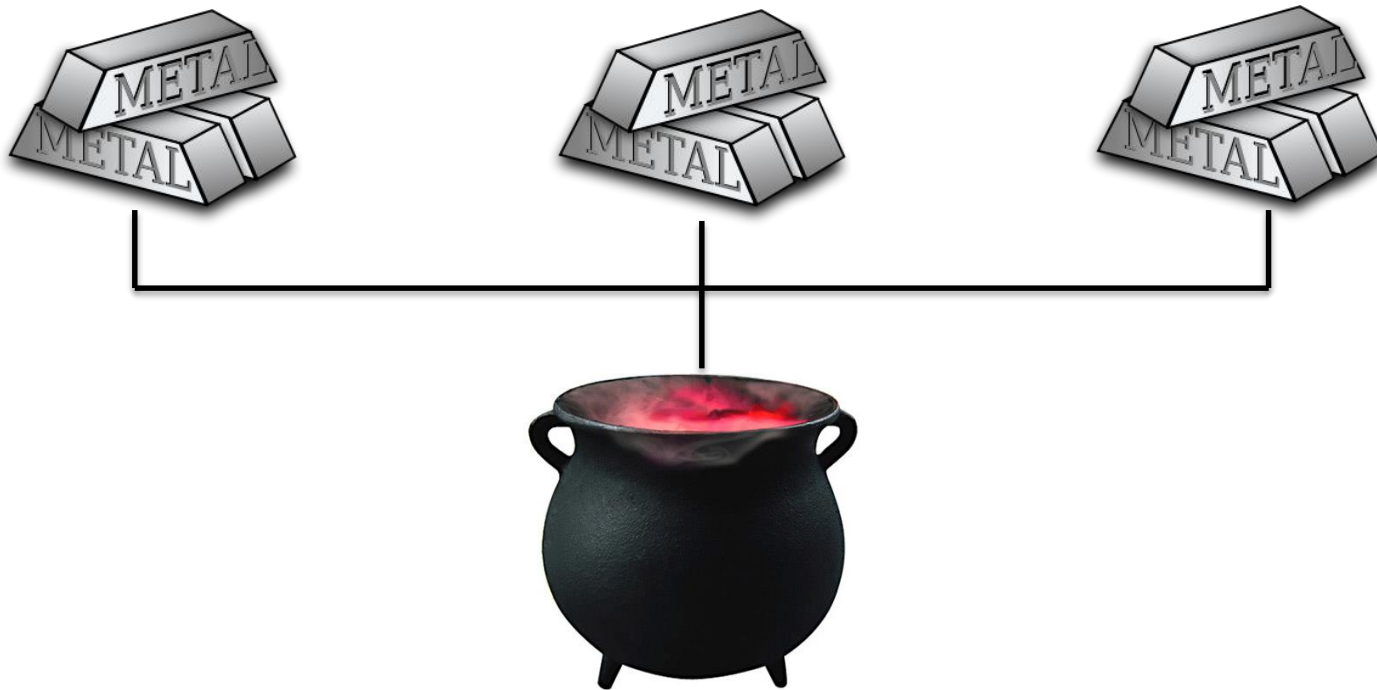


Model Based Predictive Control for Flood Regulation: Demer



Alloy melting point modelling

Problem & Objectives



How to predict material properties when designing a new alloy?

Alloy melting point modelling

Data

InsPyro NV

Random
measurements
of melting point
in function of
constituent
concentrations

10 different
materials!

130 000 points
sampled!

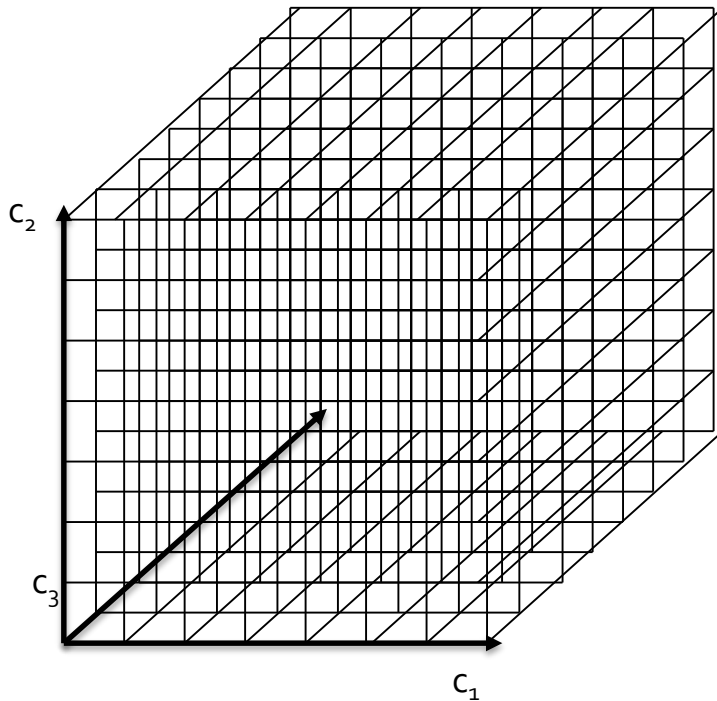
Many missing
values.

Curse of Dimensionality!

Alloy melting point modelling

Expertise In Action

Represent as tensor



Decompose in factors



A diagram showing the decomposition of a tensor \mathcal{T} into a sum of rank-1 tensors. The tensor \mathcal{T} is represented as a cube. It is equal to the sum of rank-1 tensors, each represented as a vertical bar (matrix) and a horizontal bar (vector). The first term is $u_1^{(1)}$ (vertical bar) and $u_1^{(2)}$ (horizontal bar), with $u_1^{(3)}$ (diagonal bar) above it. This is followed by an ellipsis and then the last term $u_R^{(1)}$ (vertical bar) and $u_R^{(2)}$ (horizontal bar), with $u_R^{(3)}$ (diagonal bar) above it.

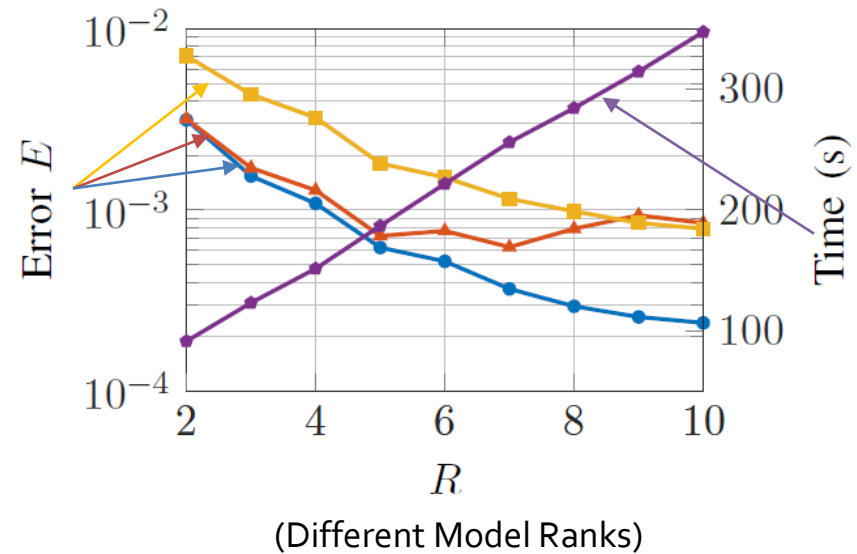
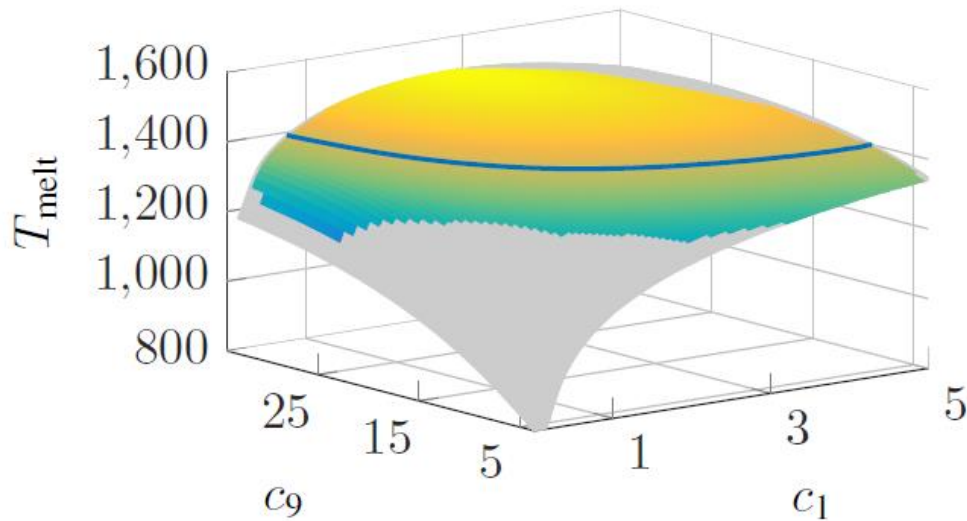
New model has only
4500 parameters

⇒ Curse of Dimensionality broken!

Alloy melting point modelling

Problem Solved

Economical prediction of new material properties



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

Energy

Industry

Environment

Social networks

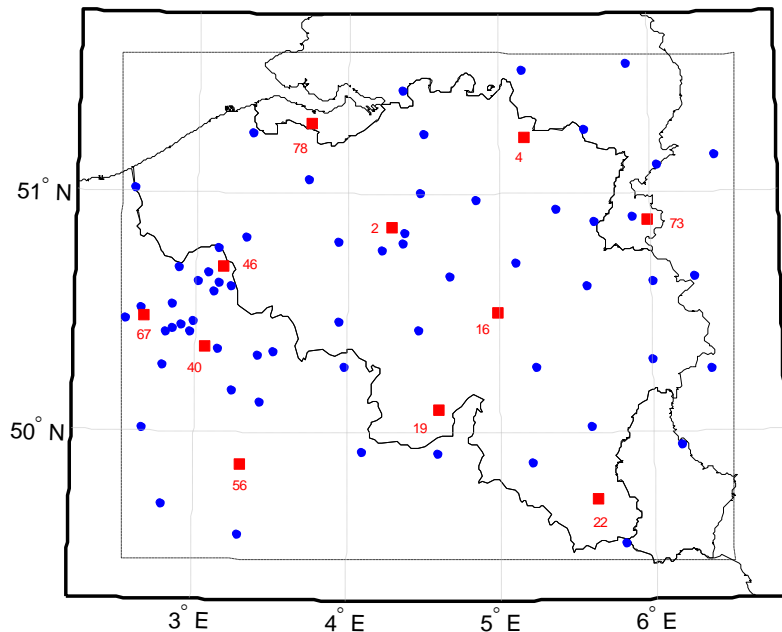
Finance and Fraud

Health

Data Assimilation

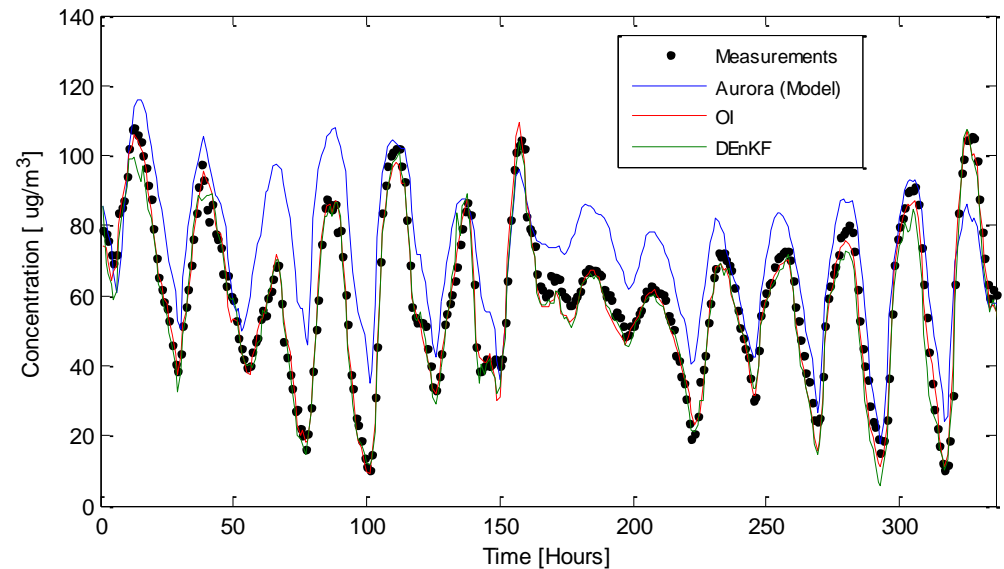
The Deterministic Ensemble Kalman Filter (DEnKF) and the OI technique have been used to improve the O₃ estimates of the Air-quality model AURORA.

O₃ air-quality stations



- - Assimilation stations
- - Validation stations

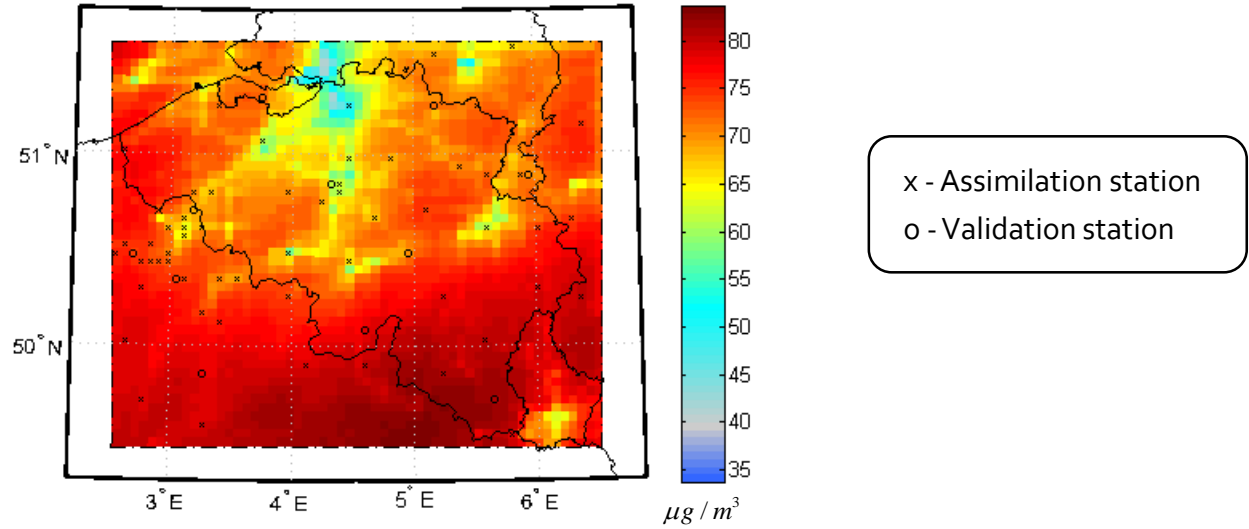
Average of the O₃ concentration over the validation stations



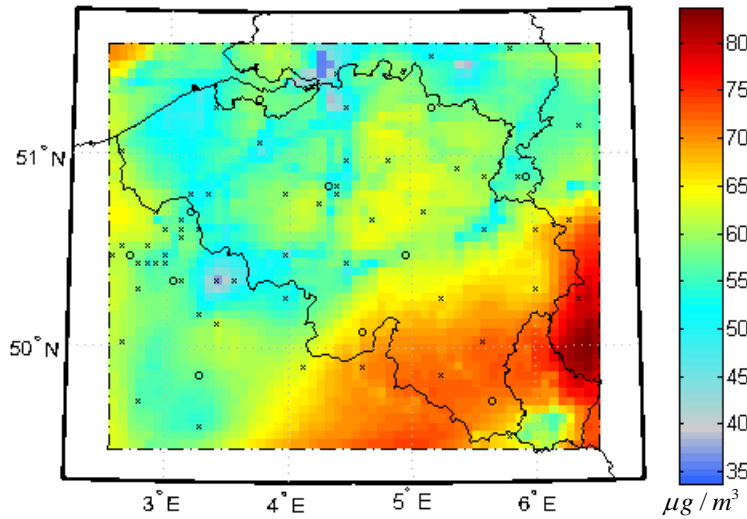
Starting date: May 28th, 2005 at midnight

Average of the O₃ concentration field over the 14 day period

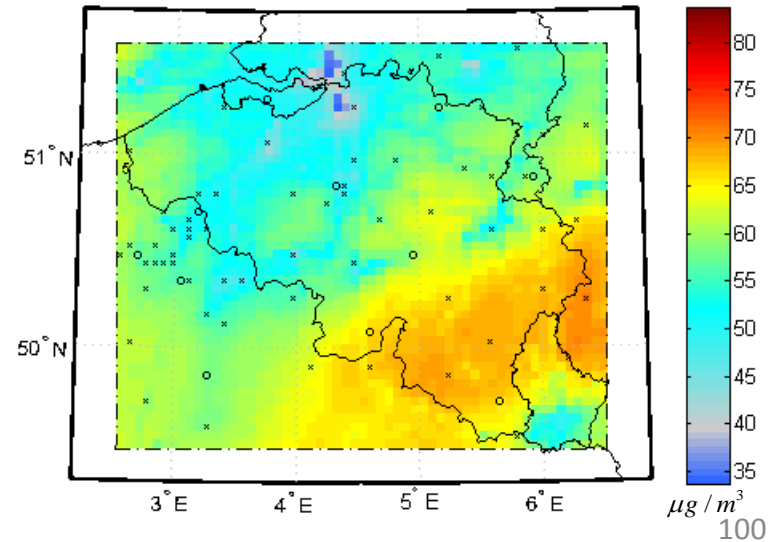
Free-run of Aurora



Optimal Interpolation



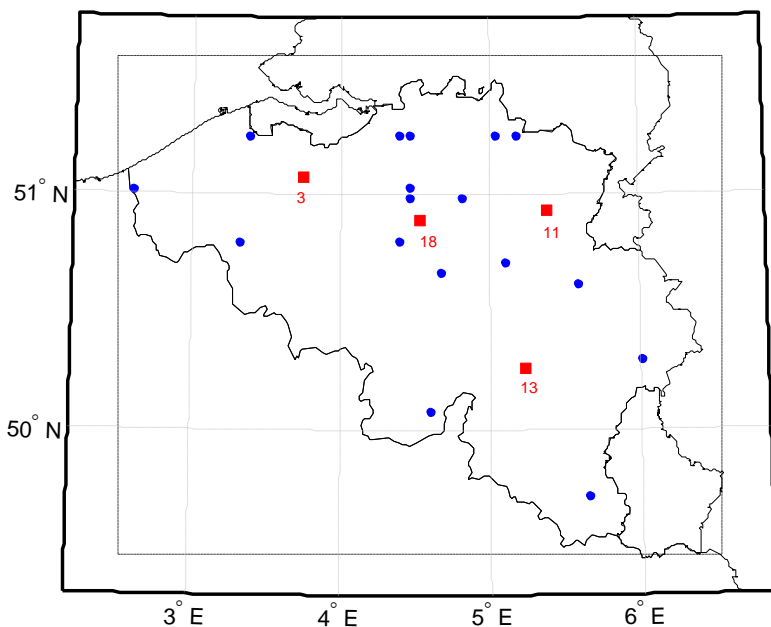
DEnKF



Data Assimilation

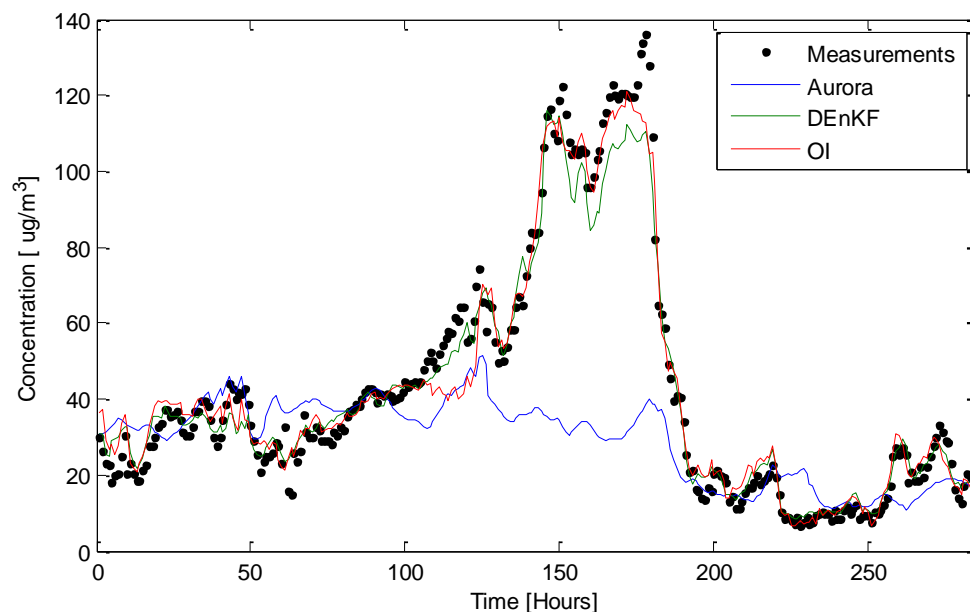
The Deterministic Ensemble Kalman Filter (DEnKF) and the OI technique have been used to improve the PM₁₀ estimates of the Air-quality model AURORA.

PM₁₀ air-quality stations



- - Assimilation stations
- - Validation stations

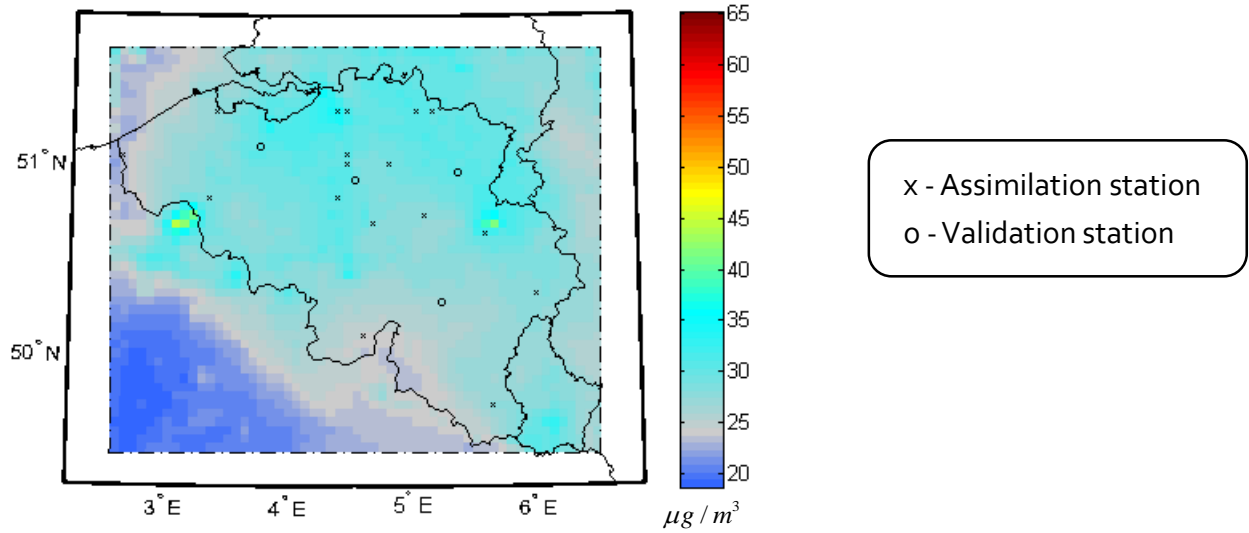
Average of the PM₁₀ concentration over the validation stations



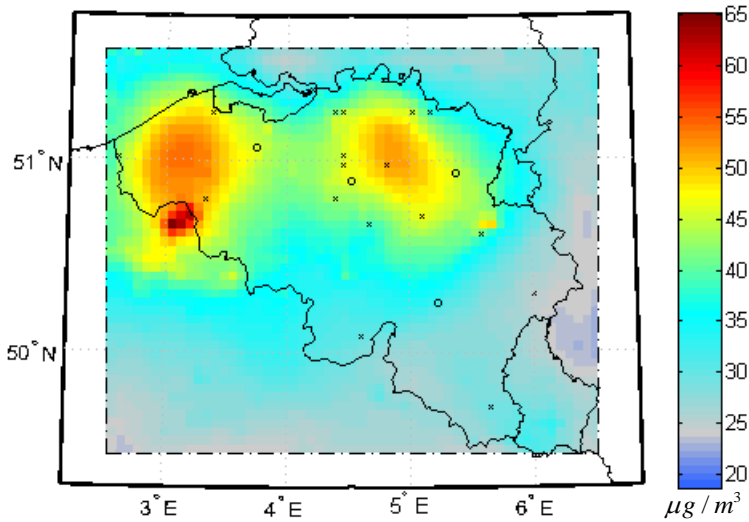
Starting date: January 20th, 2010 at midnight

Average of the PM₁₀ concentration field

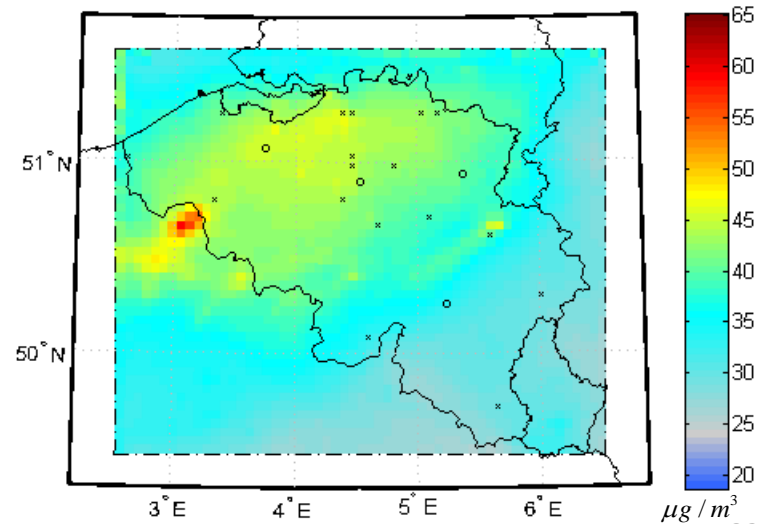
Free-run of Aurora



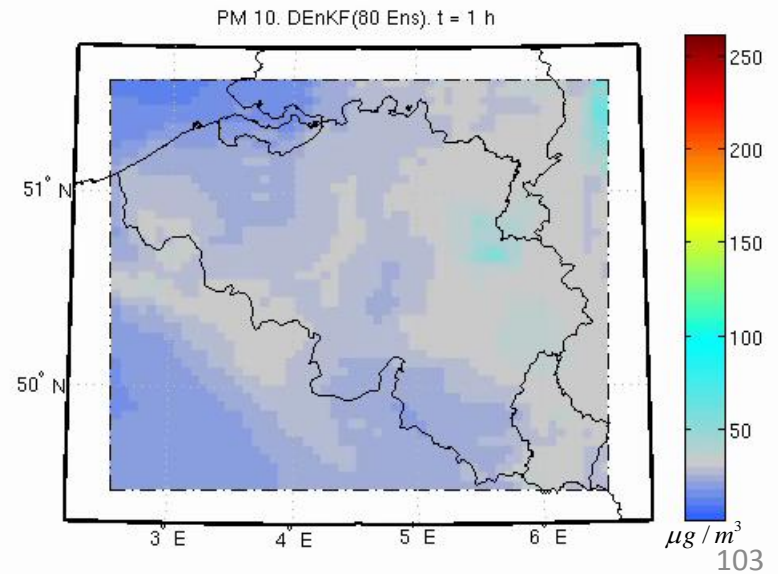
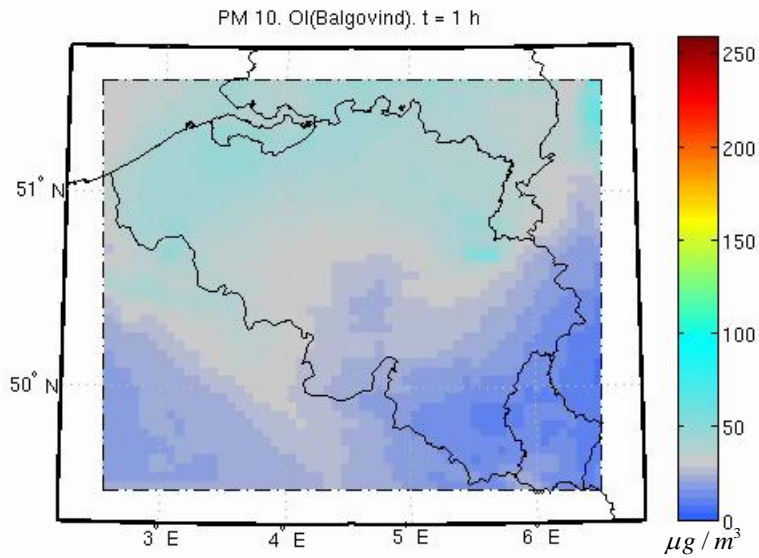
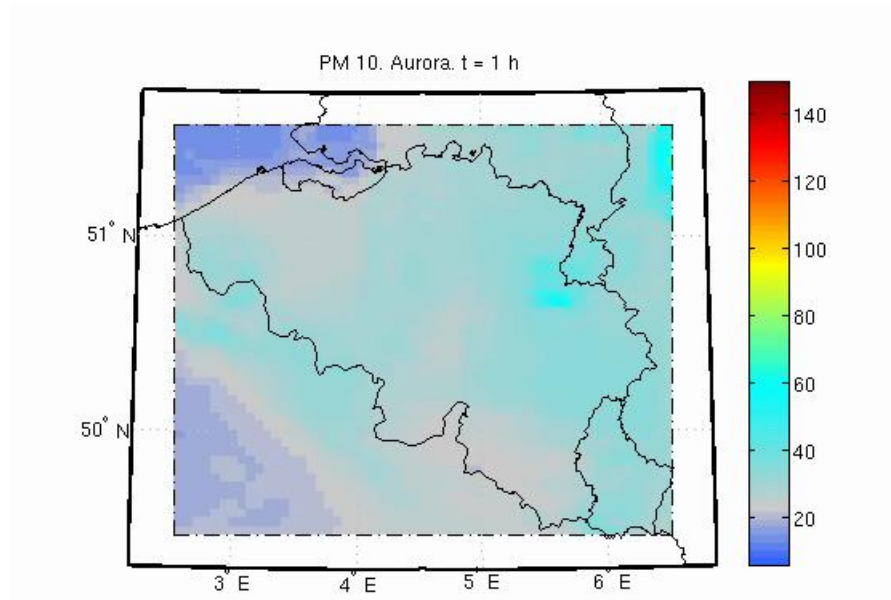
Optimal Interpolation



DEnKF



Average of the PM₁₀ concentration field



A woman with short brown hair and glasses, wearing a grey turtleneck, is looking at a tablet held by a healthcare professional in blue scrubs. The healthcare professional has her hair in a bun and is wearing a stethoscope. The scene is set in a clinical or office environment.

Energy

Industry

Environment

Social networks

Finance and Fraud

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

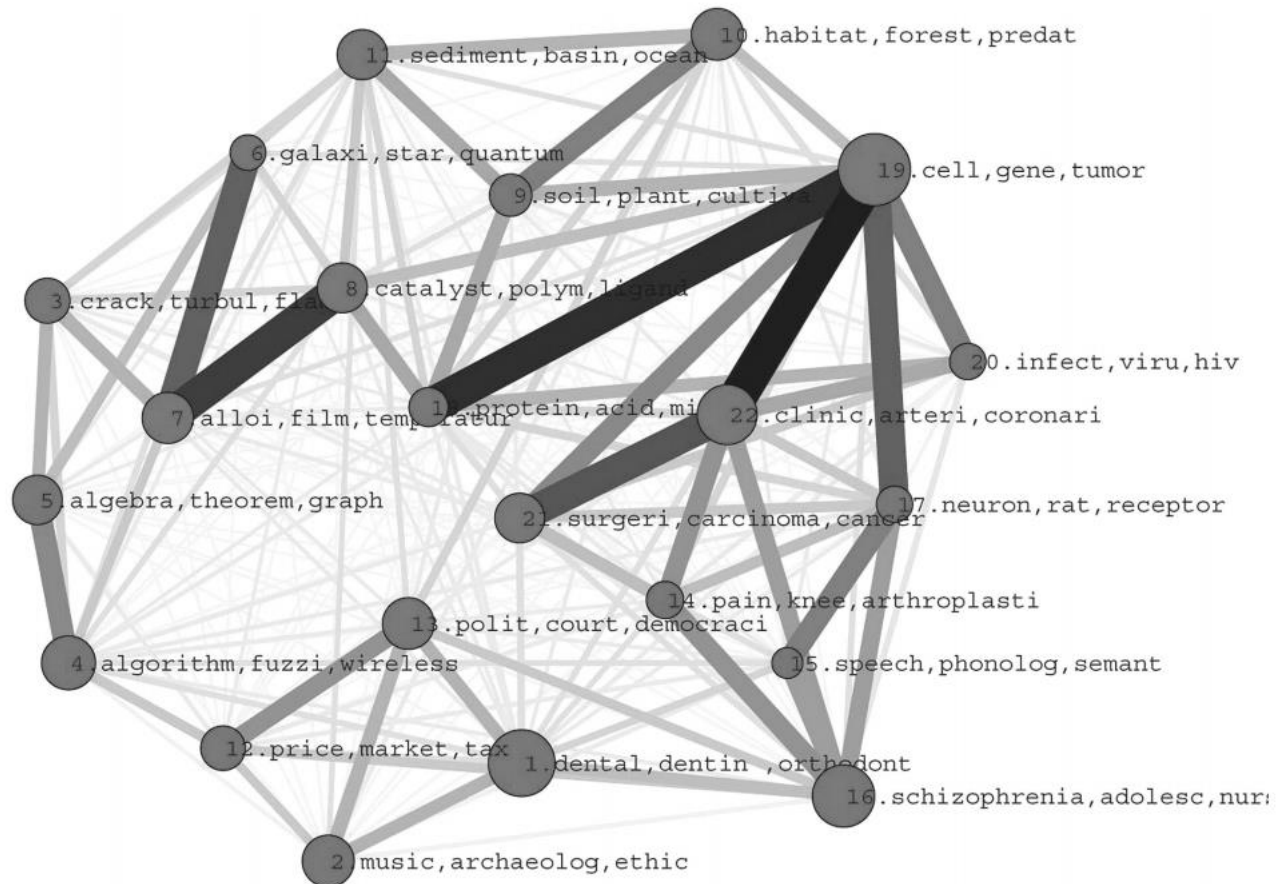
In scientometrics, information from journals can be categorized lexically or with citations. An important area of scientometric research is the clustering or mapping of scientific publications. The widely used method of cocitation clustering was introduced independently by Small (1973, 1978) and Marshakova (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

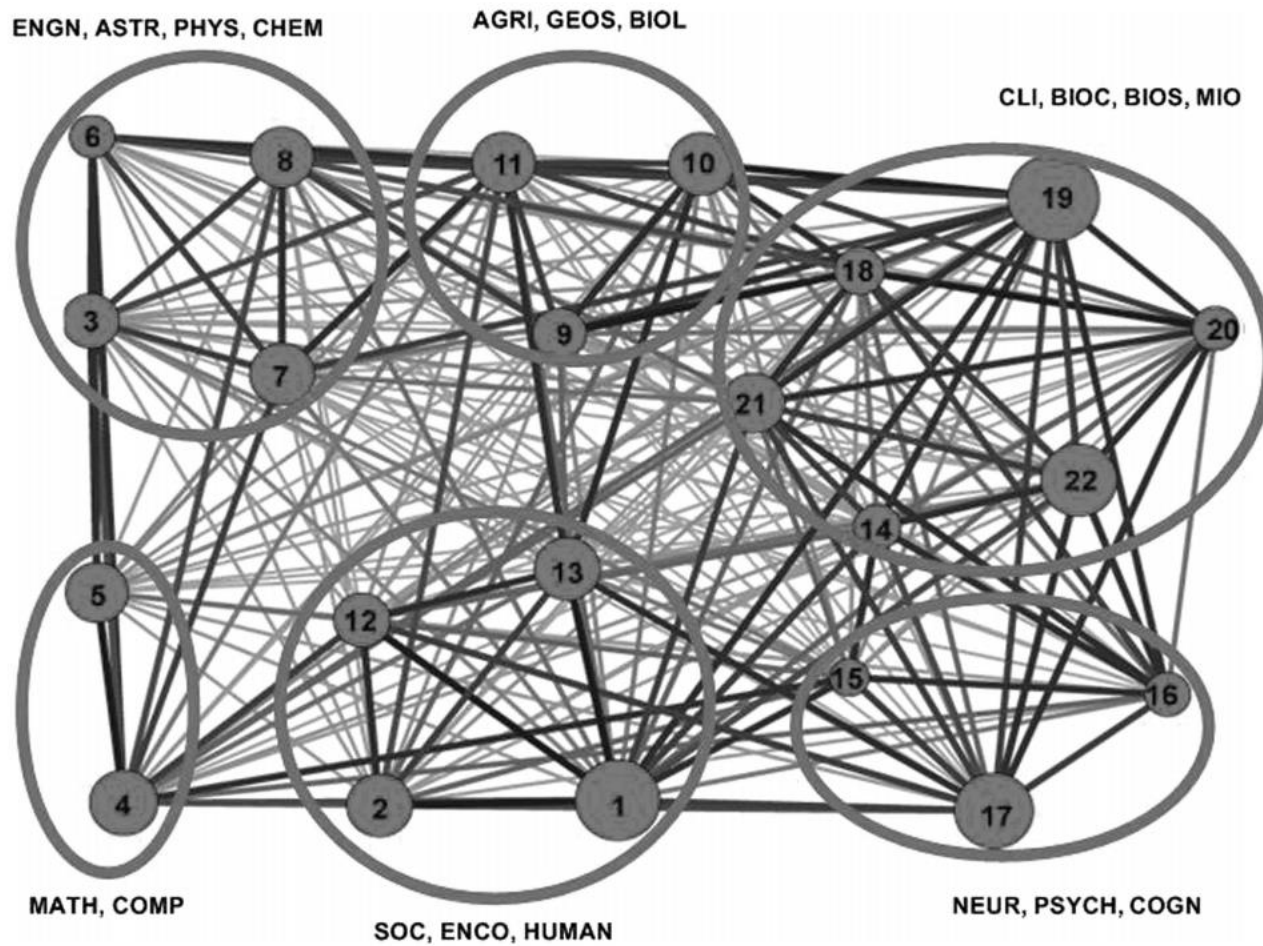
Received July 7, 2009; revised October 31, 2009; accepted December 30, 2009

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

Journal Clustering



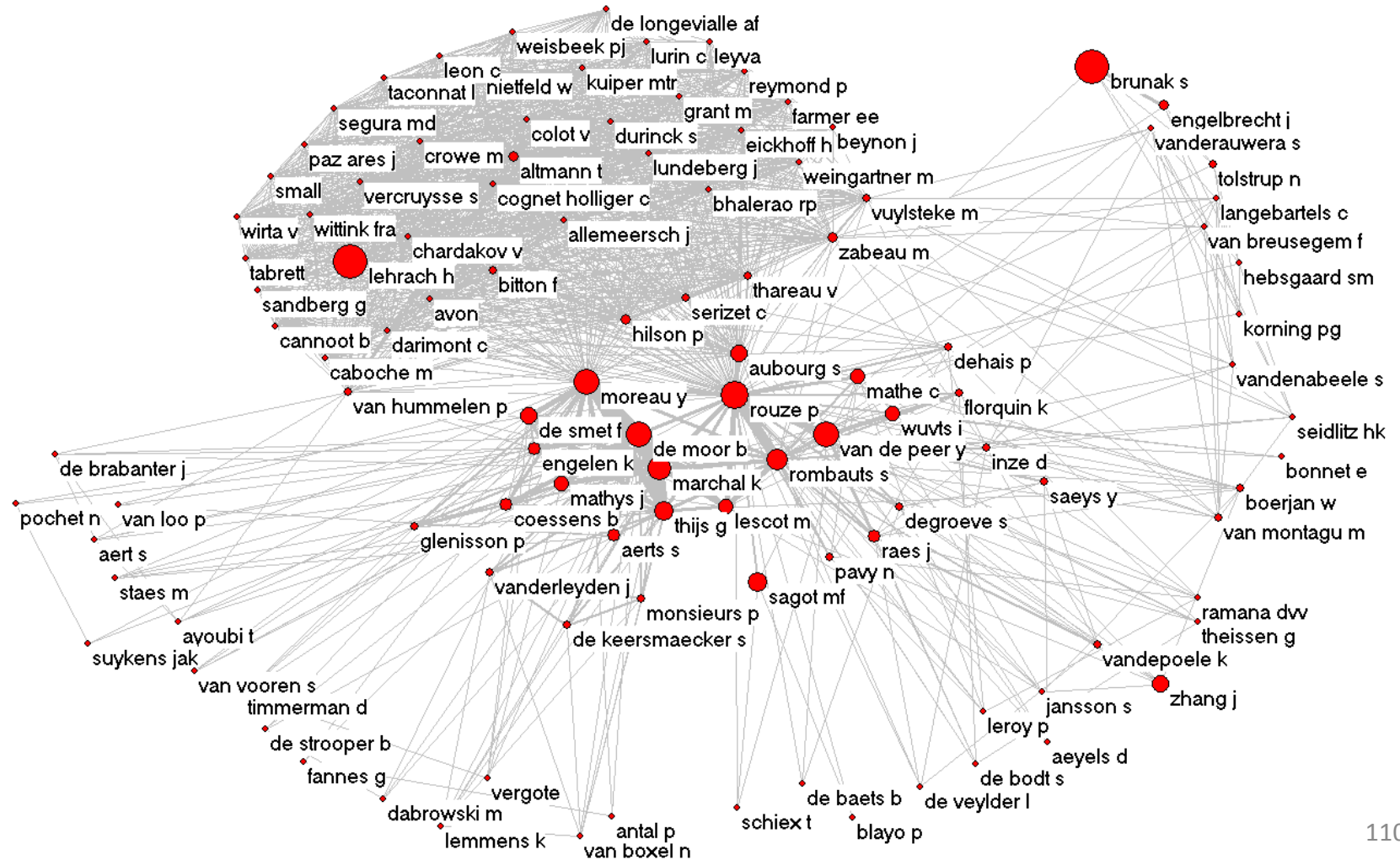
Journal Clustering



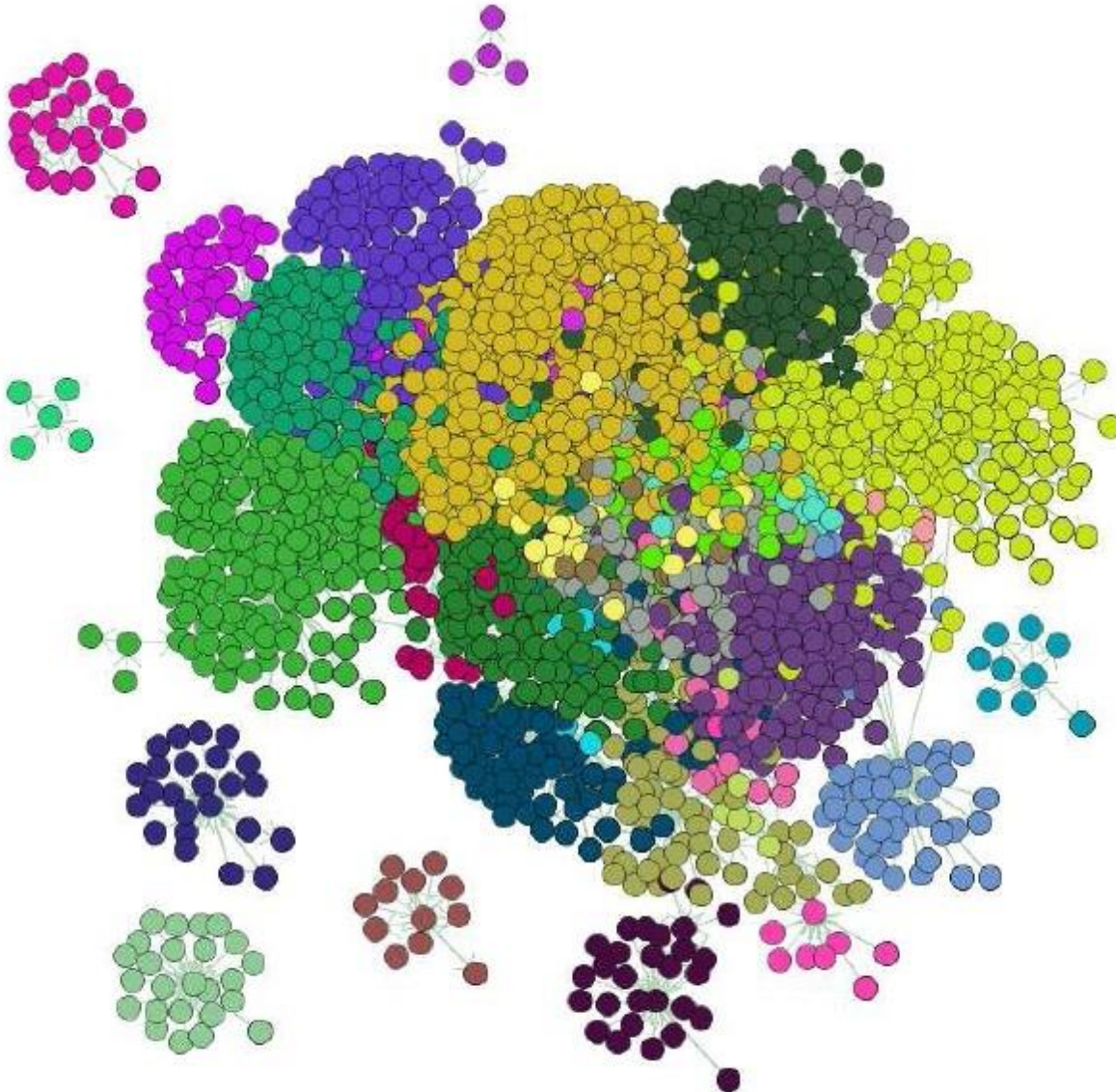
Author Collaboration Clustering



Author Collaboration Clustering



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.

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

Energy

Industry

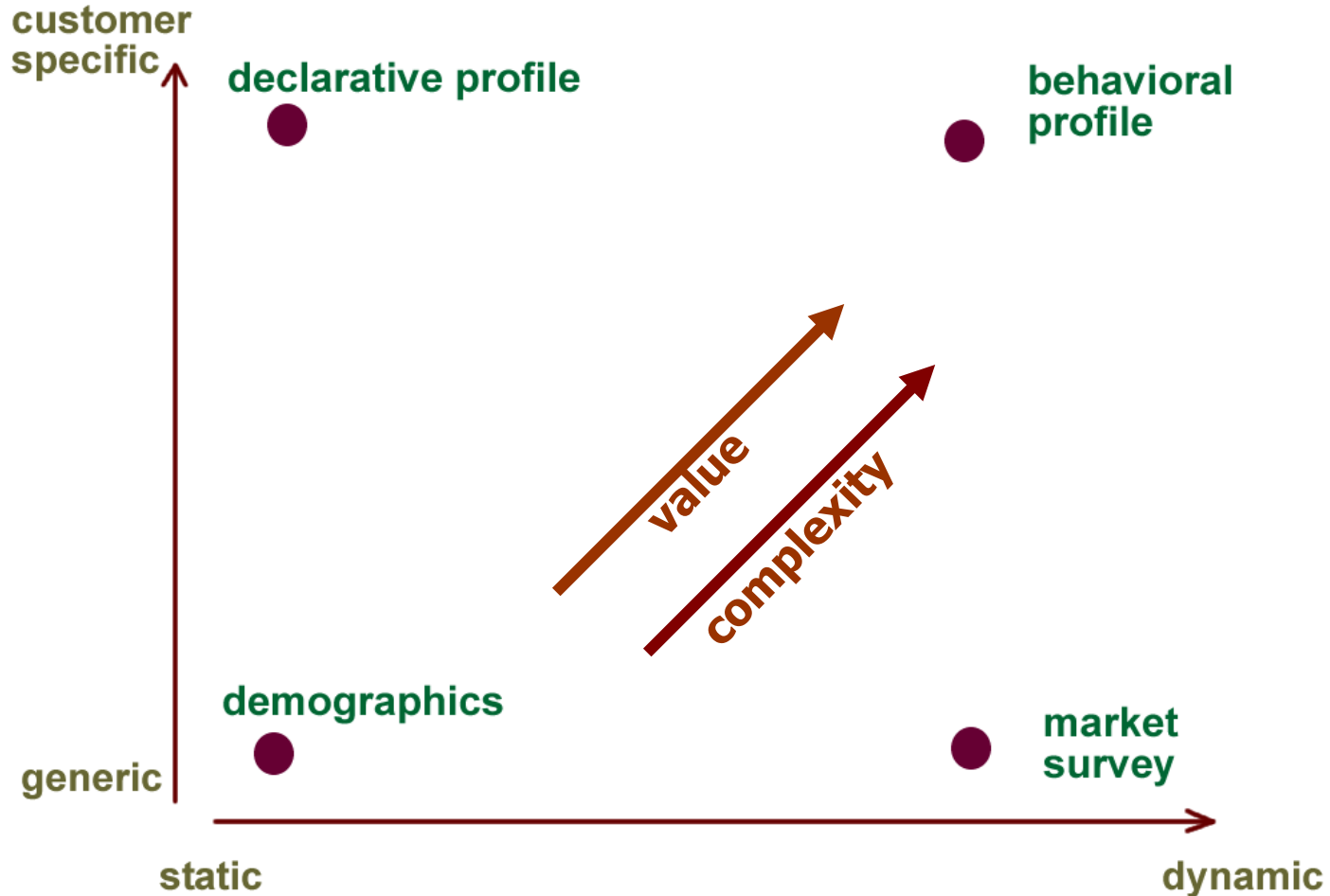
Environment

Education

Finance and Fraud

Health

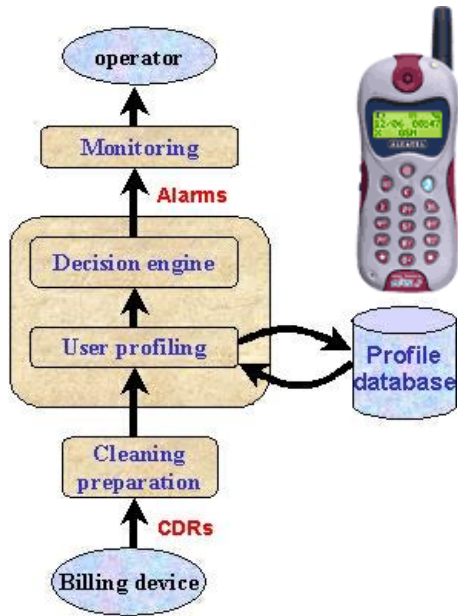
Customer Intelligence



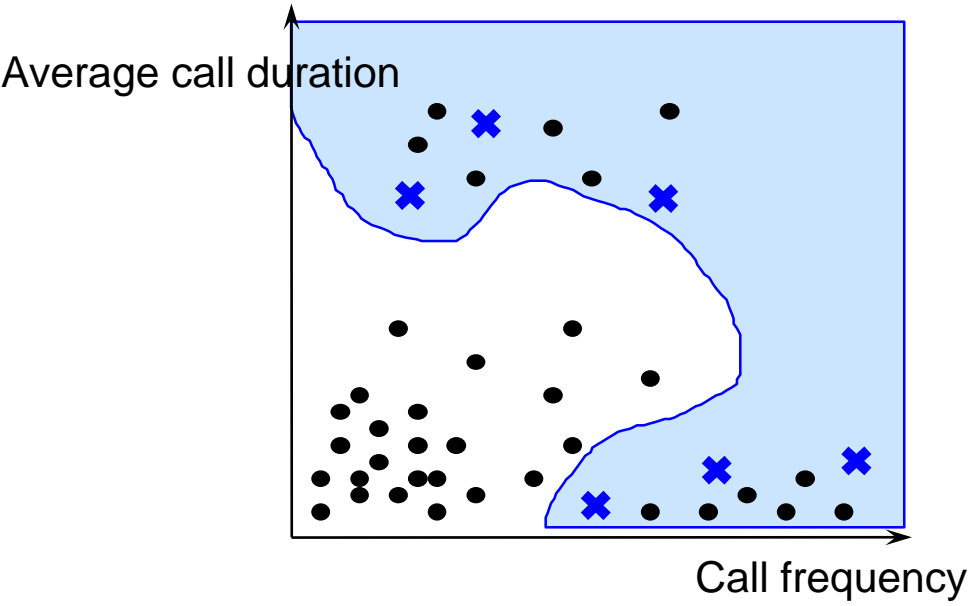
Customer Intelligence



Fraud detection on mobile phone network

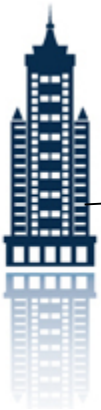


	Short Duration	Long Duration	High Frequency	International	Same Destination	Off Peak	Call Forwarding	Behaviour Change
Direct calling		X	X	X			X	
ABX fraud	X		X		X	X		X
Freephone fraud	X		X		X			X
Premium rate fraud		X	X		X			X
Subscription fraud			X					
Handset theft		X	X	X	X			X



Electric Market Segmentation

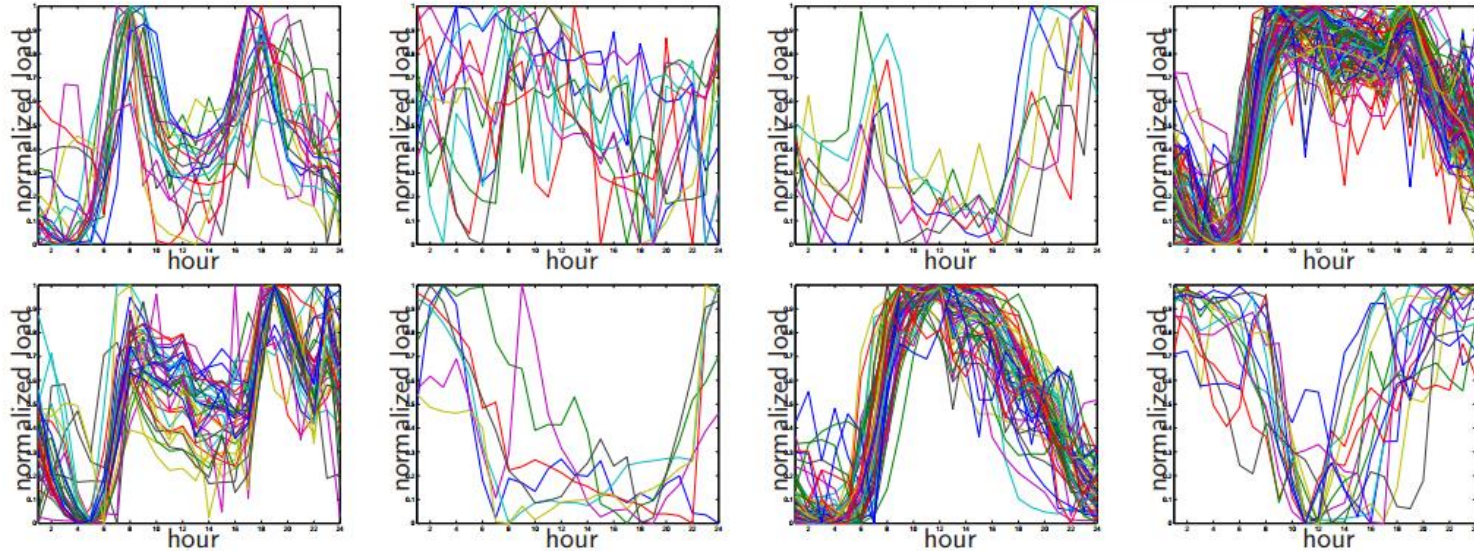
Problem & Objectives



Electric Market Segmentation

Data

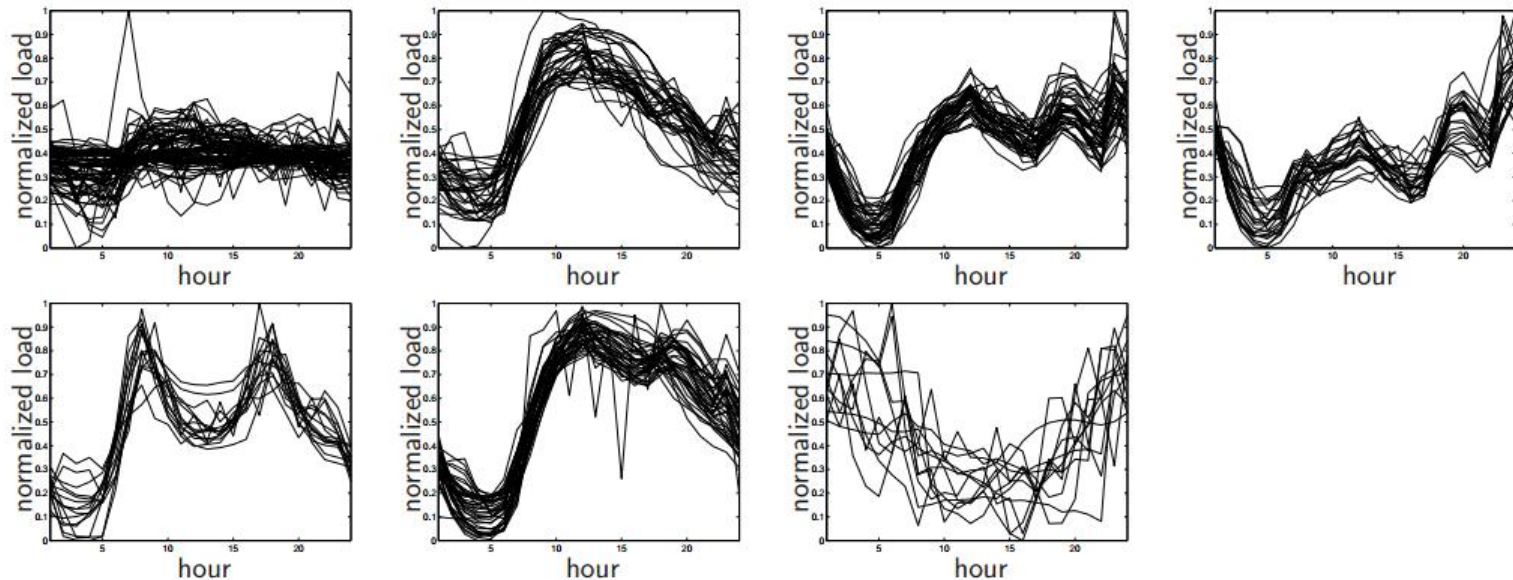
Power load: 245 substations, hourly (5 years)
Periodic AR modelling: dim reduction 43.824 \rightarrow 24
k-means applied after dimensionality reduction



Electric Market Segmentation

Expertise In Action

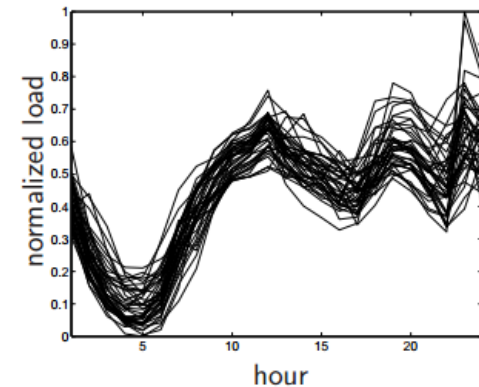
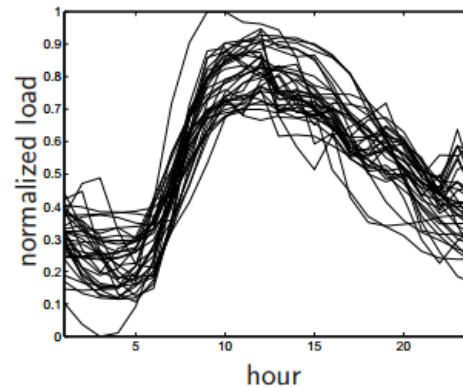
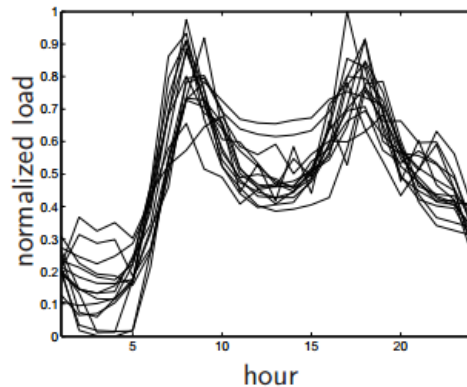
Application of kernel spectral clustering, directly in high dim $d = 43.824$
Model selection on kernel parameter and number of clusters



[Alzate, Espinoza, De Moor, Suykens, 2009]

Electric Market Segmentation

Problem Solved



Electricity load: 245 substations in Belgian grid (1/2 train, 1/2 validation)
 $x_i \in \mathbb{R}^{43.824}$: spectral clustering on **high dimensional data** (5 years)

3 of 7 detected clusters:




- 1: **Residential profile**: morning and evening peaks
- 2: **Business profile**: peaked around noon
- 3: **Industrial profile**: increasing morning, oscillating afternoon and evening

Bankruptcy Prediction

Problem



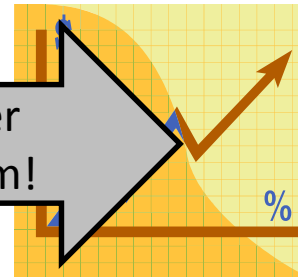
What we know!

-  : Capital
-  : Credit
-  : Total assets

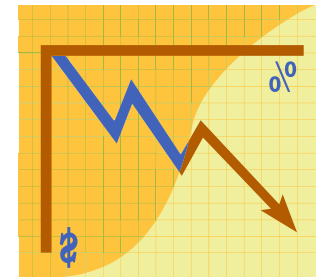
In total : 40 candidate inputs

What we don't!

Classifier
Algorithm!



or



?

Bankruptcy Prediction

Solution

Cross-validation

Only a subset of inputs is important
↓
Concise classifier!

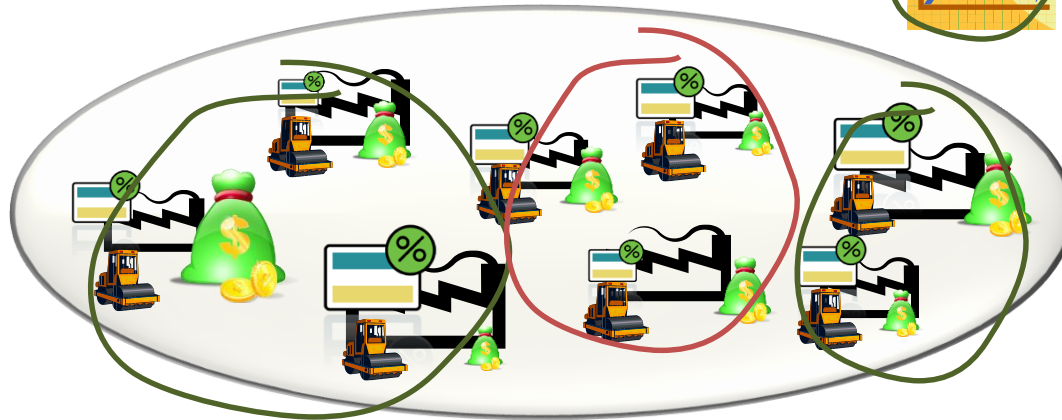
⇒ > 90% accurate

Classifier

Decides on basis of inputs



or



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

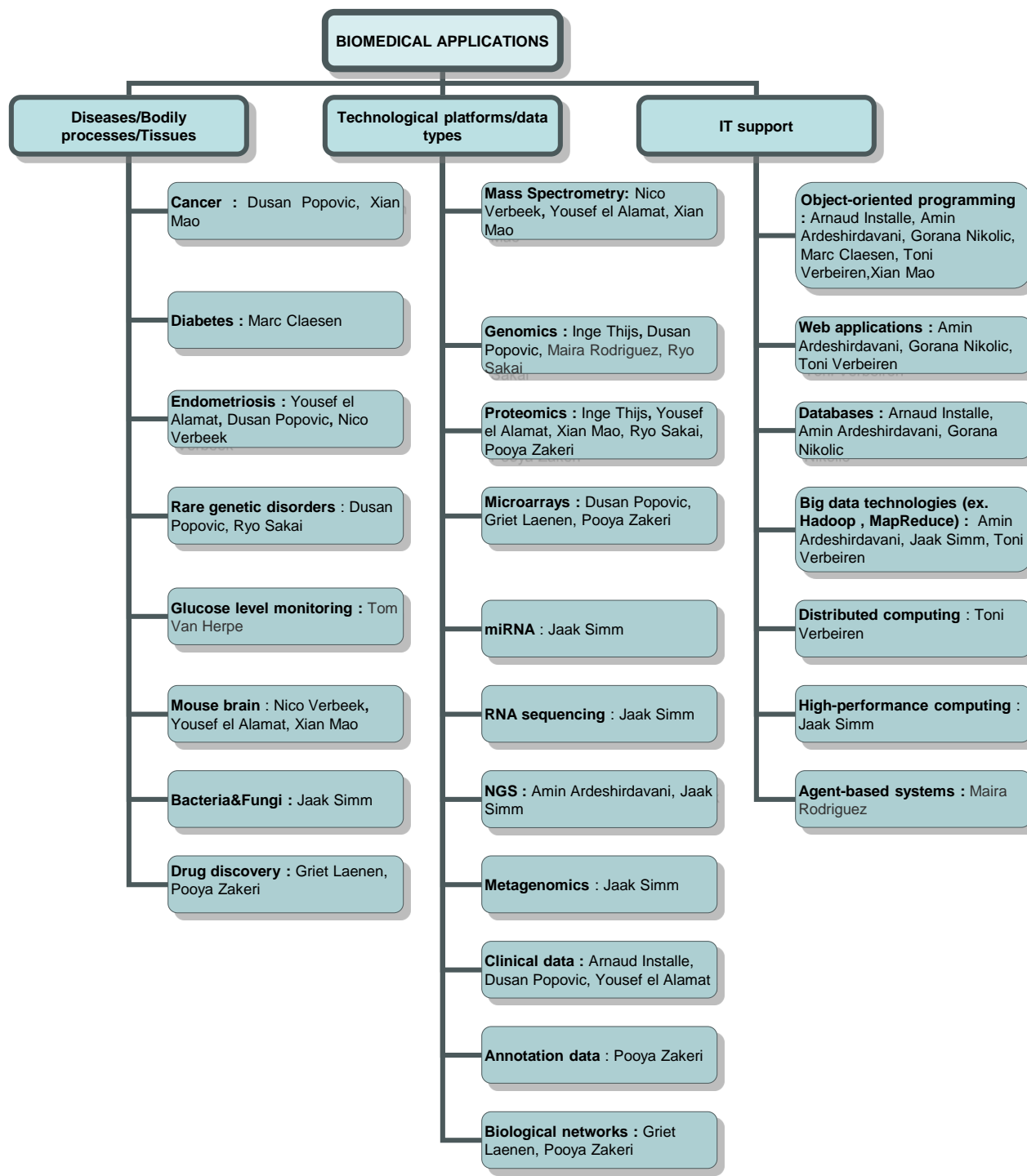
Industry

Environment

Social networks

Finance and Fraud

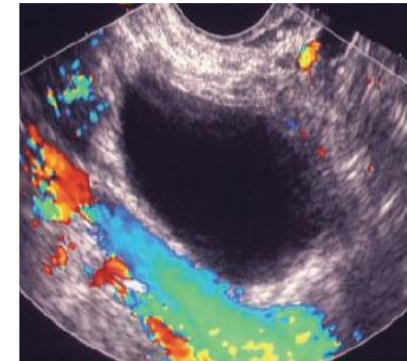
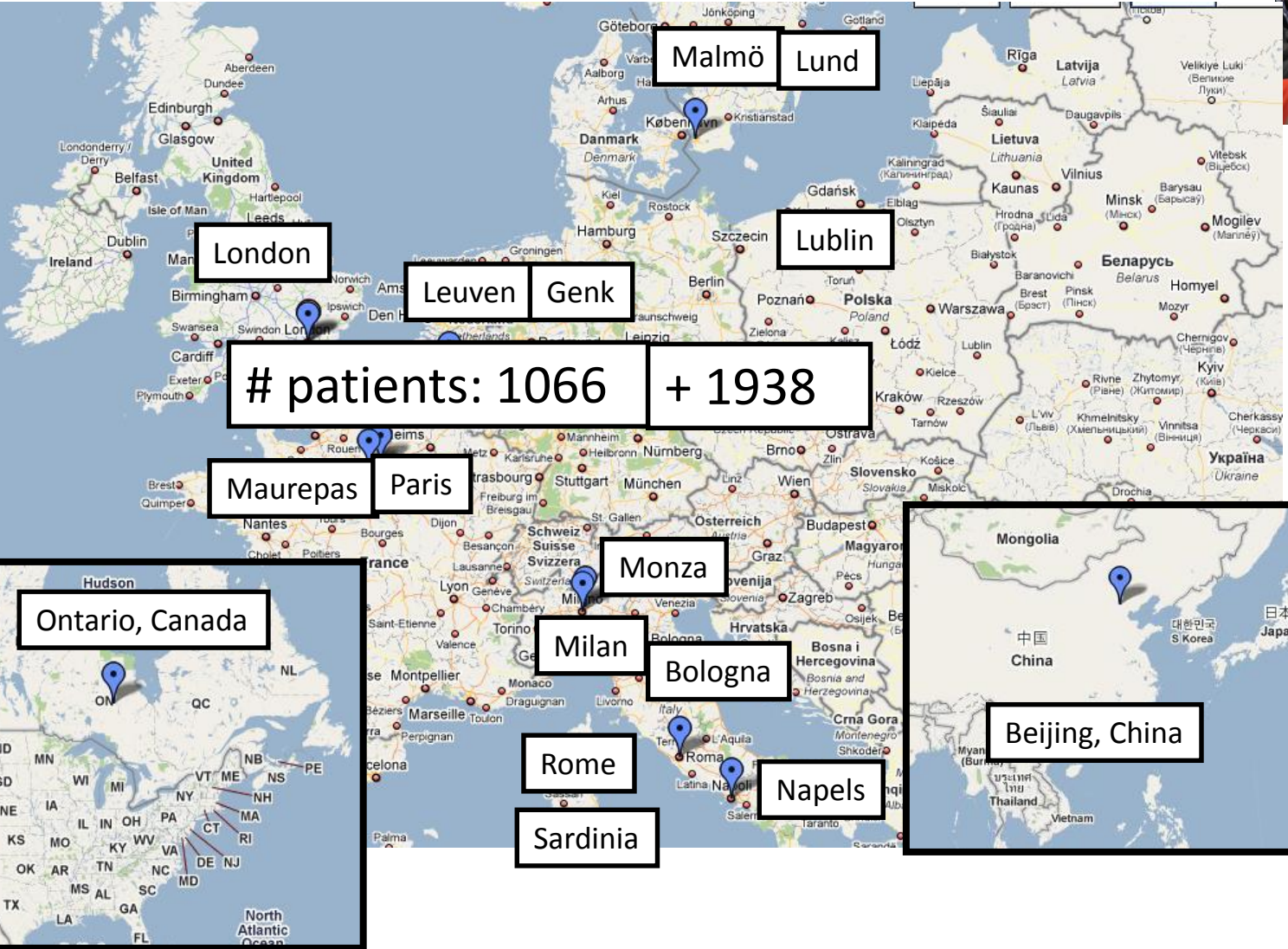
Health



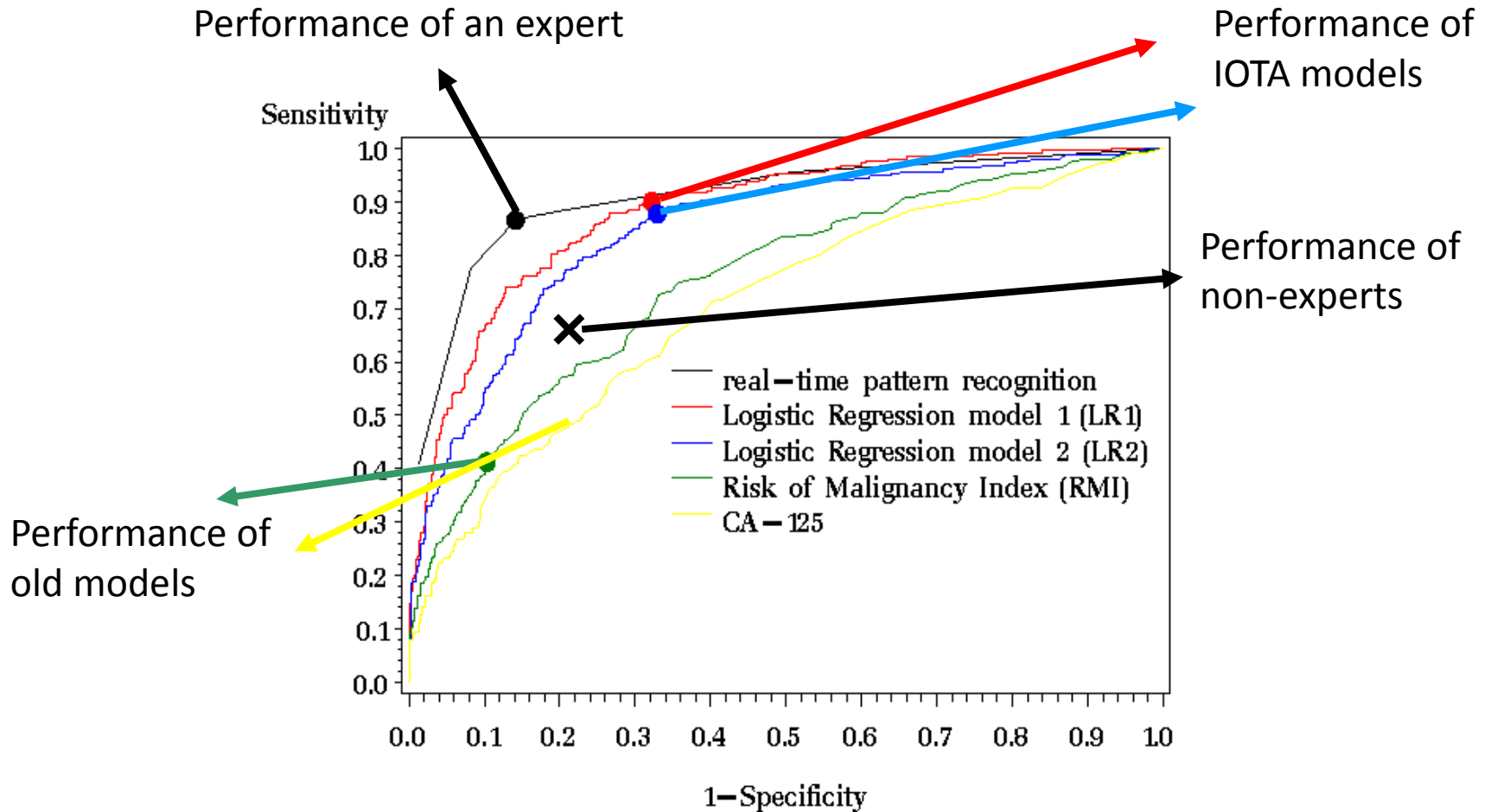
Participatory



IOTA app:
population
based
assessment
of ovarian
tumor malignancy:

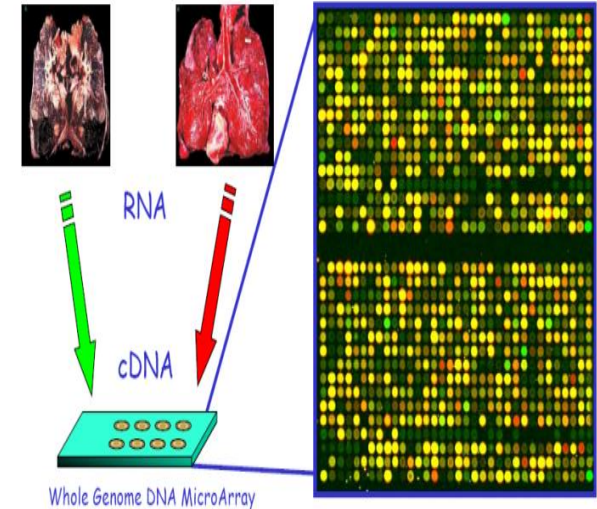
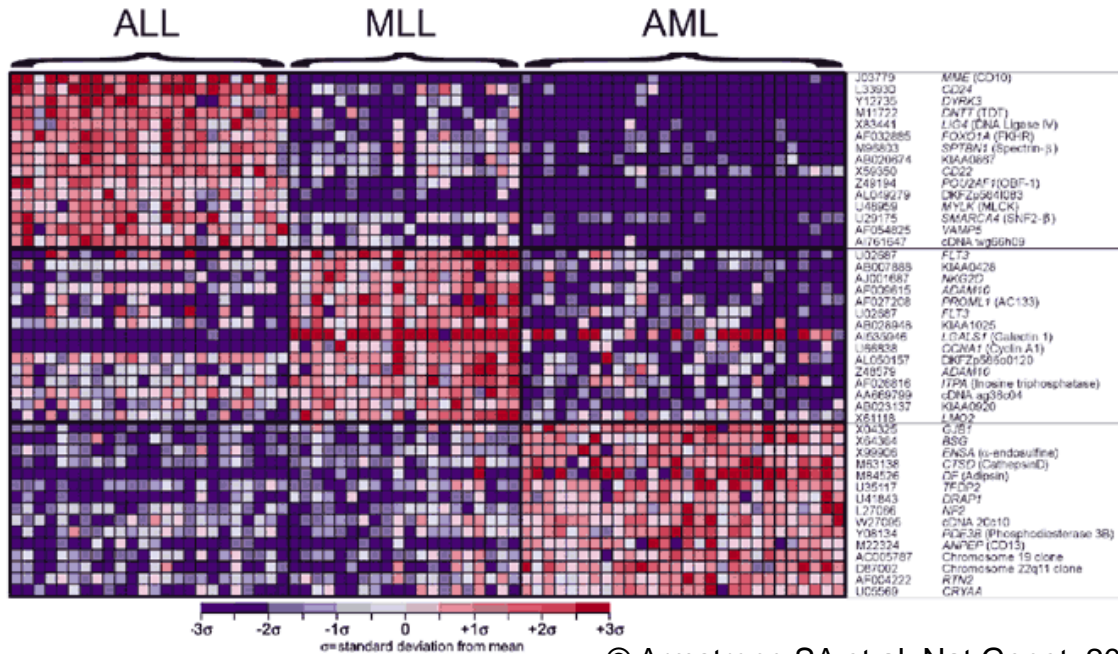


Performance



You share, we care !

Genomic markers for Leukemia



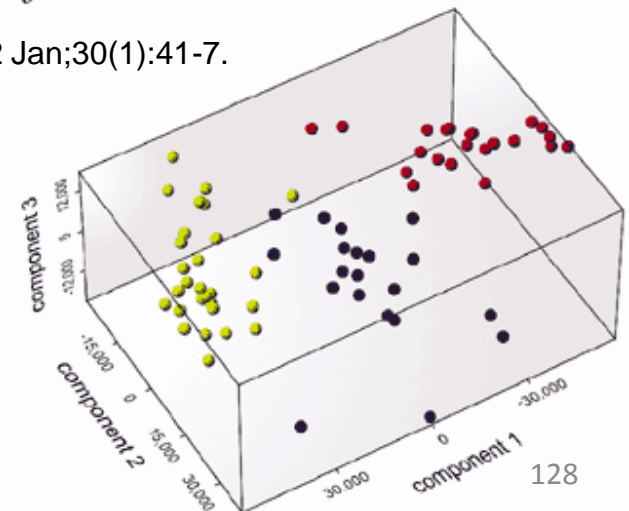
b

12 600 genes

72 patients

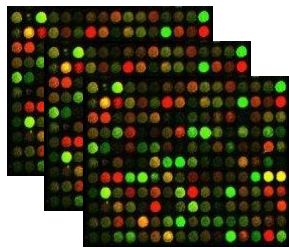
- 28 Acute Lymphoblastic Leukemia (ALL)
- 24 Acute Myeloid Leukemia (AML)
- 20 Mixed Linkage Leukemia (MLL)

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

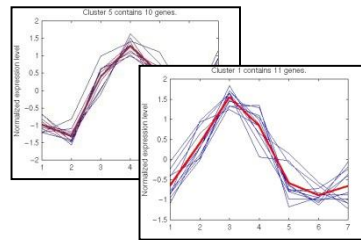


Genomic Data Fusion

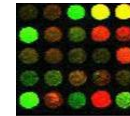
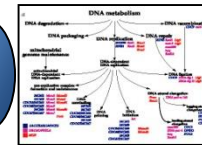
High-throughput genomics



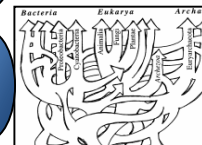
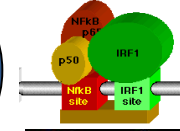
Data analysis



Information sources



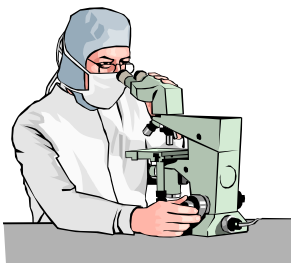
Her chaparral, at 11, she saw Michael Jackson performing on television and told Angelil that she wanted to be that big. Fine, said Angelil, who advised her to take 18 months off, during which she underwent a massive makeover that included plastic surgery, shorter hair and caps for the long incisors that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



Candidate genes

Name	Ensembl
TTR	ENSG00000118271
PAH	ENSG00000171759
G6PC	ENSG00000131482
IGF1	ENSG0000017427
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
BPA61	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

Validation



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

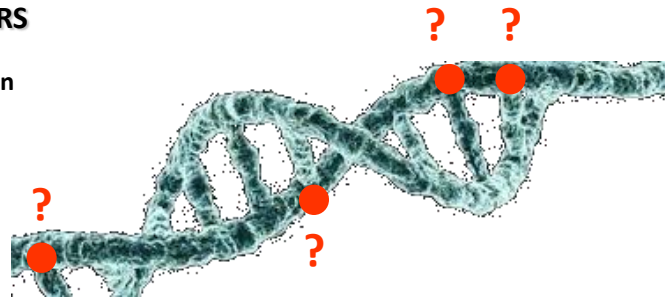
Mutation Prioritization

Problem & Objectives

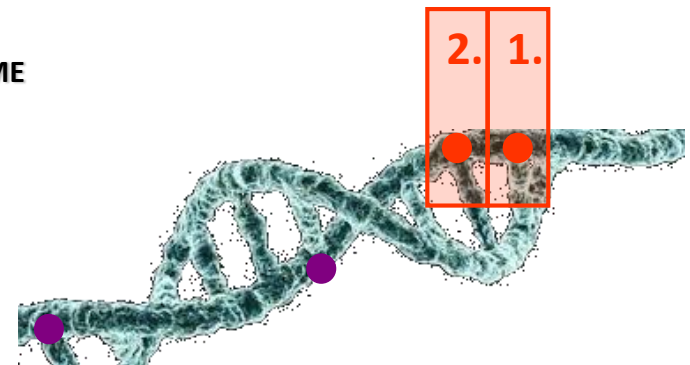


PREVALANCE OF GENETIC DISORDERS

Biomarker discovery as a classification problem



THOUSANDS MUTATIONS IN A GENOME



NEED TO PRIORITIZE MUTATIONS

Mutation prioritization

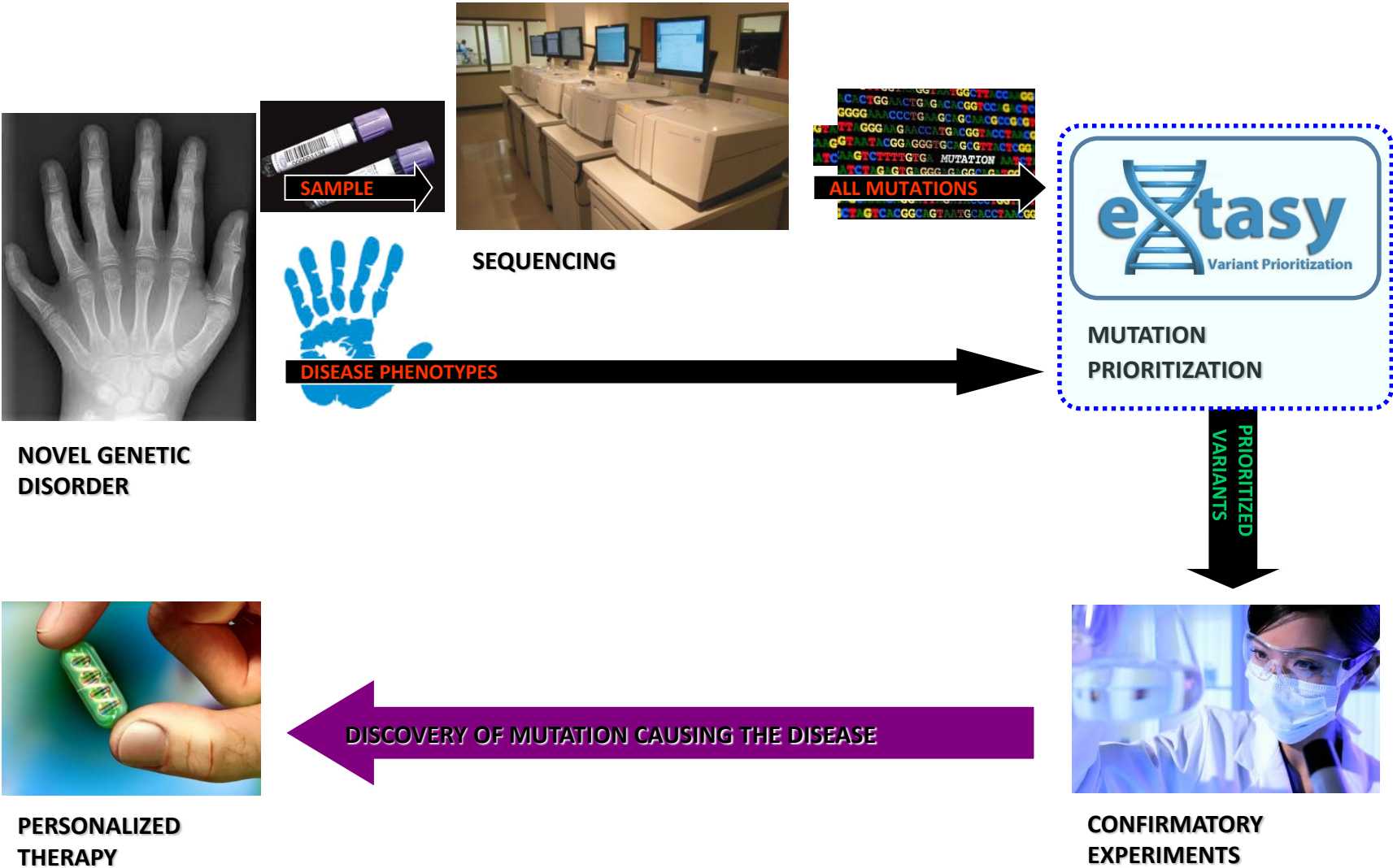
Data



Very sensitive
and confidential

Mutation prioritization

Expertise In Action



Mutation prioritization

Expertise In Action

ALL MUTATIONS

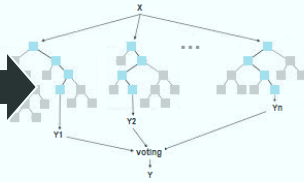
PolyPhen-2

Mutation Taster

SIFT

ANNOTATION

Other scores



DISEASE PHENOTYPES

Endeavour

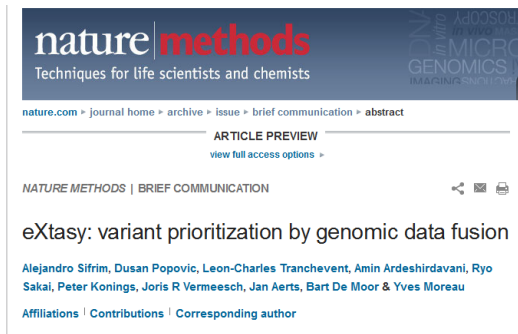
Phenotype scores

RANDOM FOREST CLASSIFIER

PRIORITIZED VARIANTS

Mutation prioritization

Problem Solved



nature methods
Techniques for life scientists and chemists

nature.com > journal home > archive > issue > brief communication > abstract

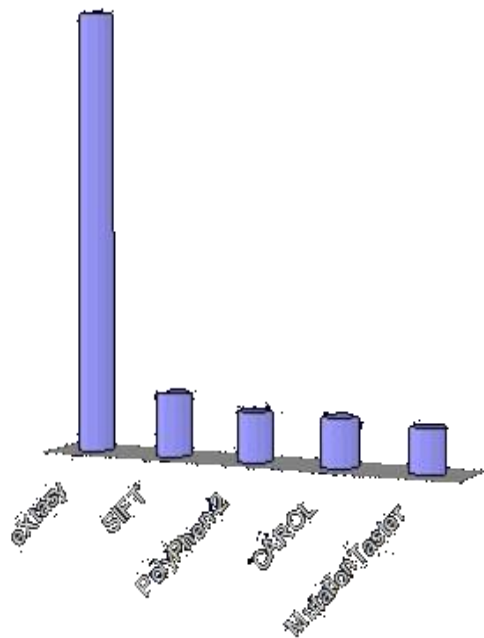
ARTICLE PREVIEW
[view full access options >](#)

NATURE METHODS | BRIEF COMMUNICATION

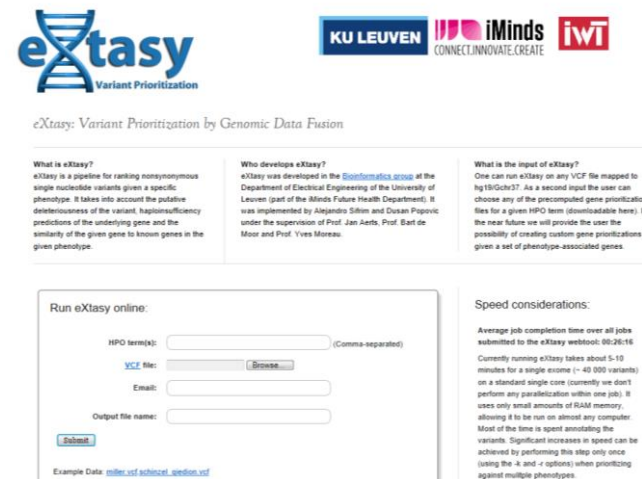
eXtasy: variant prioritization by genomic data fusion

Alejandro Sifrim, Dusan Popovic, Leon-Charles Tranchevent, Amin Ardeshirdavani, Ryo Sakai, Peter Konings, Joris R Vermeesch, Jan Aerts, Bart De Moor & Yves Moreau

[Affiliations](#) | [Contributions](#) | [Corresponding author](#)



10-FOLD INCREASE IN PRECISION



eXtasy
Variant Prioritization

KU LEUVEN | iMinds | IWT
CONNECT. INNOVATE. CREATE.

eXtasy: Variant Prioritization by Genomic Data Fusion

What is eXtasy?
eXtasy is a pipeline for ranking nonsynonymous single nucleotide variants given a specific phenotype. It takes into account the putative deleteriousness of the variant, haploinsufficiency predictions of the underlying gene and the similarity of the given gene to known genes in the given phenotype.

Who develops eXtasy?
eXtasy was developed in the [Bioinformatics group](#) at the Department of Electrical Engineering of the University of Leuven (part of the [Minds Future Health Department](#)). It was implemented by Alejandro Sifrim and Dusan Popovic under the supervision of Prof. Jan Aerts, Prof. Bart de Moor and Prof. Yves Moreau.

What is the input of eXtasy?
One can run eXtasy on any VCF file mapped to hg18/Gchv37. As a second input the user can choose any of the precomputed gene prioritization files for a given HPO term ([downloadable here](#)). In the near future we will provide the user the possibility of creating custom gene prioritizations given a set of phenotype-associated genes.

Run eXtasy online:

HPO term(s): (Comma-separated)

VCF file:

Email:

Output file name:

Example Data: [refiler.vcf](#) [schizacp_apsion.vcf](#)

Speed considerations:

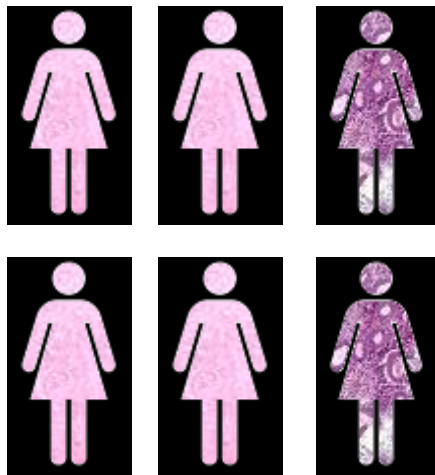
Average job completion time over all jobs submitted to the eXtasy webtool: 00:26:16

Currently running eXtasy takes about 5-10 minutes for a single enzyme (~40 000 variants) on a standard single core (currently we don't perform any parallelization within one job). It uses only small amounts of RAM memory, allowing it to be run on almost any computer. Most of the time is spent annotating the variants. Significant increases in speed can be achieved by performing this step only once (using the `-s` and `-r` options) when prioritizing against multiple phenotypes.

FREE AND EASY-TO-USE WEB TOOL

Non-invasive diagnosis of endometriosis

Problem & Objectives



PREVELANCE OF
ENDOMETRIOSIS



INVASIVE DIAGNOSTICS



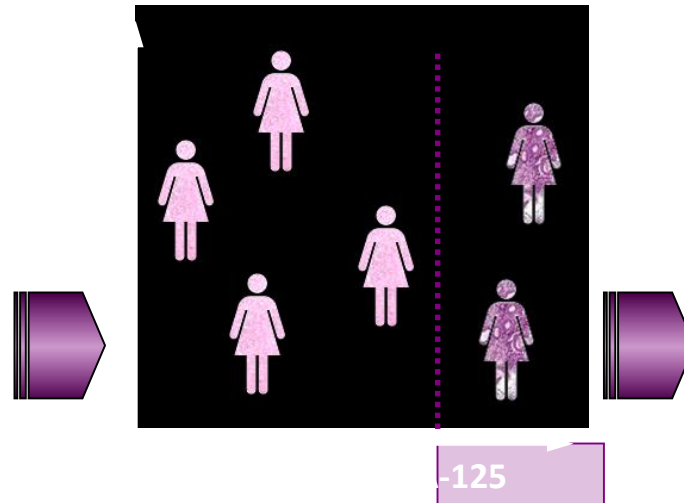
LONG DELAYS IN DIAGNOSIS

Non-invasive diagnosis of endometriosis

Expertise In Action



**PLASMA PROTEINS
LEVELS IN PATIENTS**



FEATURE SELECTION

Biomarker	Cycle phase	Cut-off
CA-125	All	> 12.5 U/ml
Glycodelin	All	> 18 ng/ml
VEGF	All	> 1.5 pg/ml
IGFBP-3	All	> 210 ng/ml
sICAM-1	All	< 243 ng/ml
CA 19-9	All	> 9.5 IU/ml
sICAM-1	Menstrual	< 254.6 ng/ml
IL-1 beta	Follicular	< 0.9 pg/ml
IL-6	Follicular	< 10 pg/ml
IFN-γ	Follicular	< 76 pg/ml
TNF-α	Follicular	< 45.6 pg/ml
IGFBP-3	Follicular	> 200 ng/ml
Glycodelin	Follicular	> 9.0 ng/ml
CA-125	Follicular	> 11.5 U/ml
CA-125	Luteal	> 13.5 U/ml
CA 19-9	Luteal	> 7.5 IU/ml

BIOMARKERS



Big Data

What

Who

Six dimensions

Data

Compute Infrastructure

Storage Infrastructure

Analytics

Visualization

Security & Privacy

Machine learning as a commodity

Expertise

Books & Spin-offs

Algorithms

Applications