

DATA

MINING

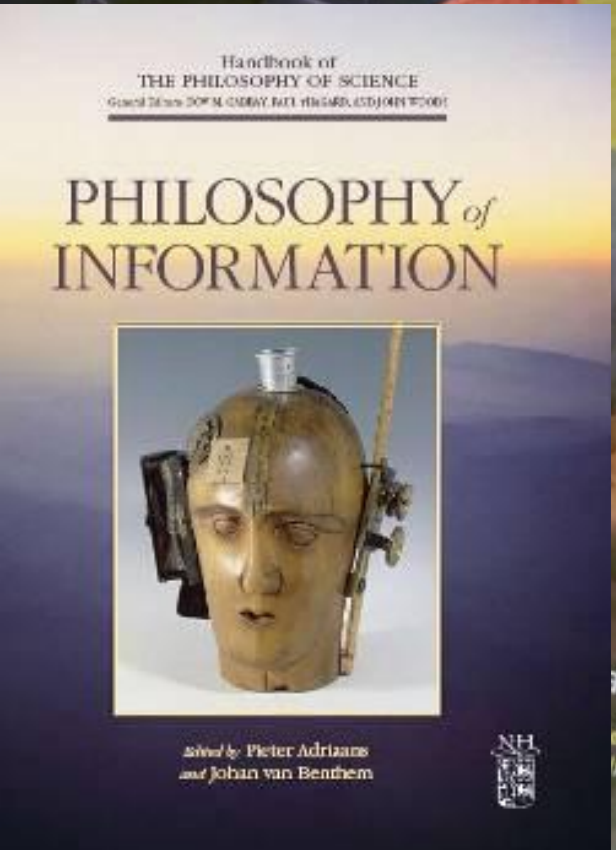
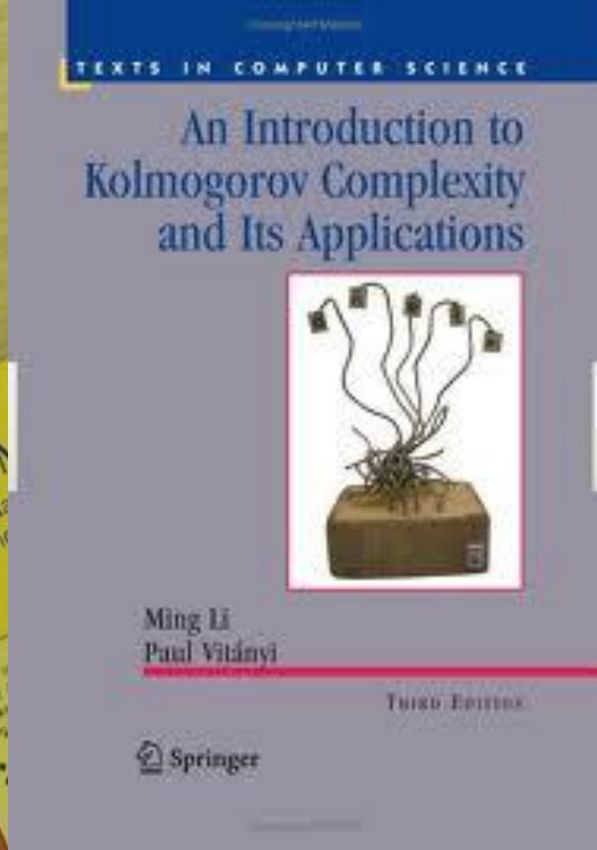
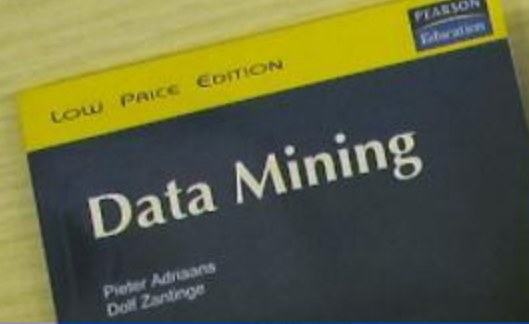
Pieter Adriaans
Dolf Zantinge


ADDISON-WESLEY



Looking Backward, Looking Forward

PTDM @ ICDM
Pieter Adriaans
Universiteit van
Amsterdam



Looking backward 1985-2012

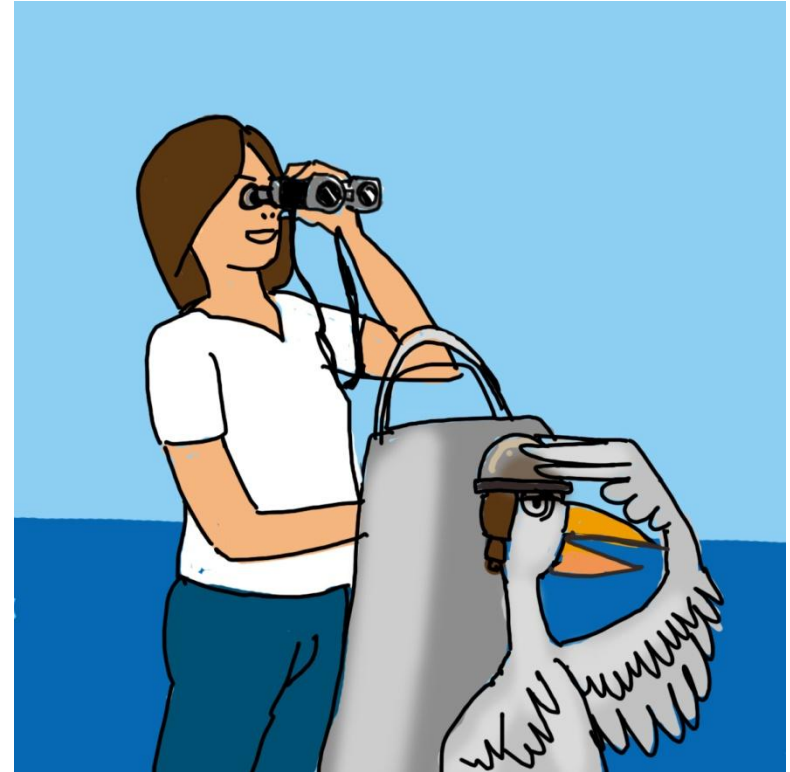
More than 200 AI,ML,DM projects

- AI, Expert Systems, Prolog Data Base technology (OBIS 1988, Captains, 1992)
- Data Mining: Workstation, Client Server, Mature ML research, large data bases
- Time dimension (ASM, ICT)
- Structured Data (GI, EMILE)
- Dynamic Systems (JSF, Robosail)
- E-science (VI-e)
- Big Data (Commit, Data2Semantics)



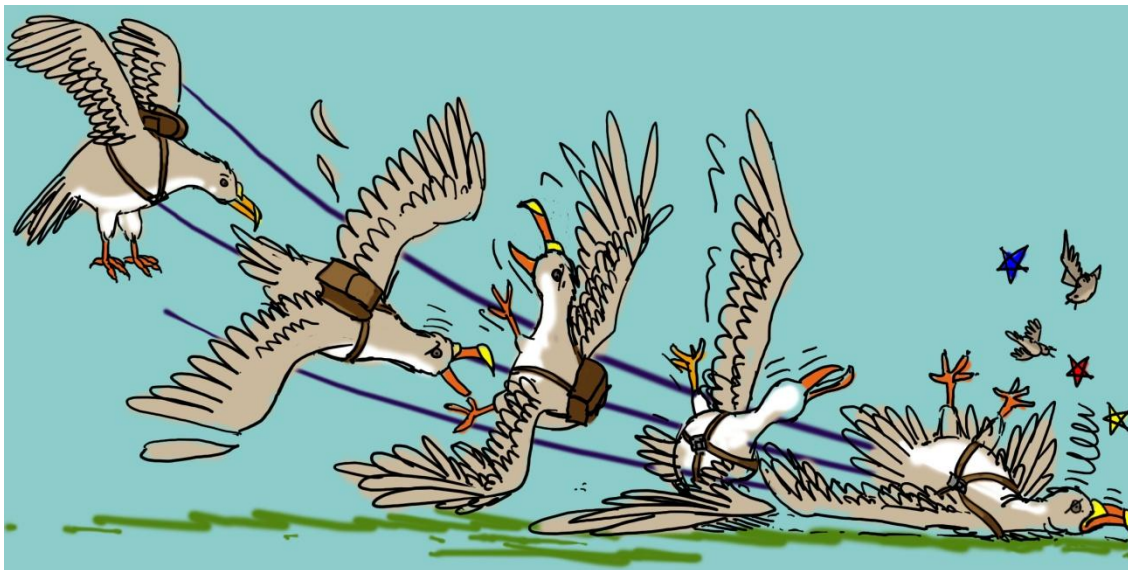
Looking forward

- Complexity measures:
facticity
- Empirical incompleteness
theories
- e-science methodology
for the 21st century



1992: Captains

- Crew Availability planning KLM
- Oracle Database
- Prolog planning algorithms
- Genetic Algorithms for Bid Prediction



Plan board

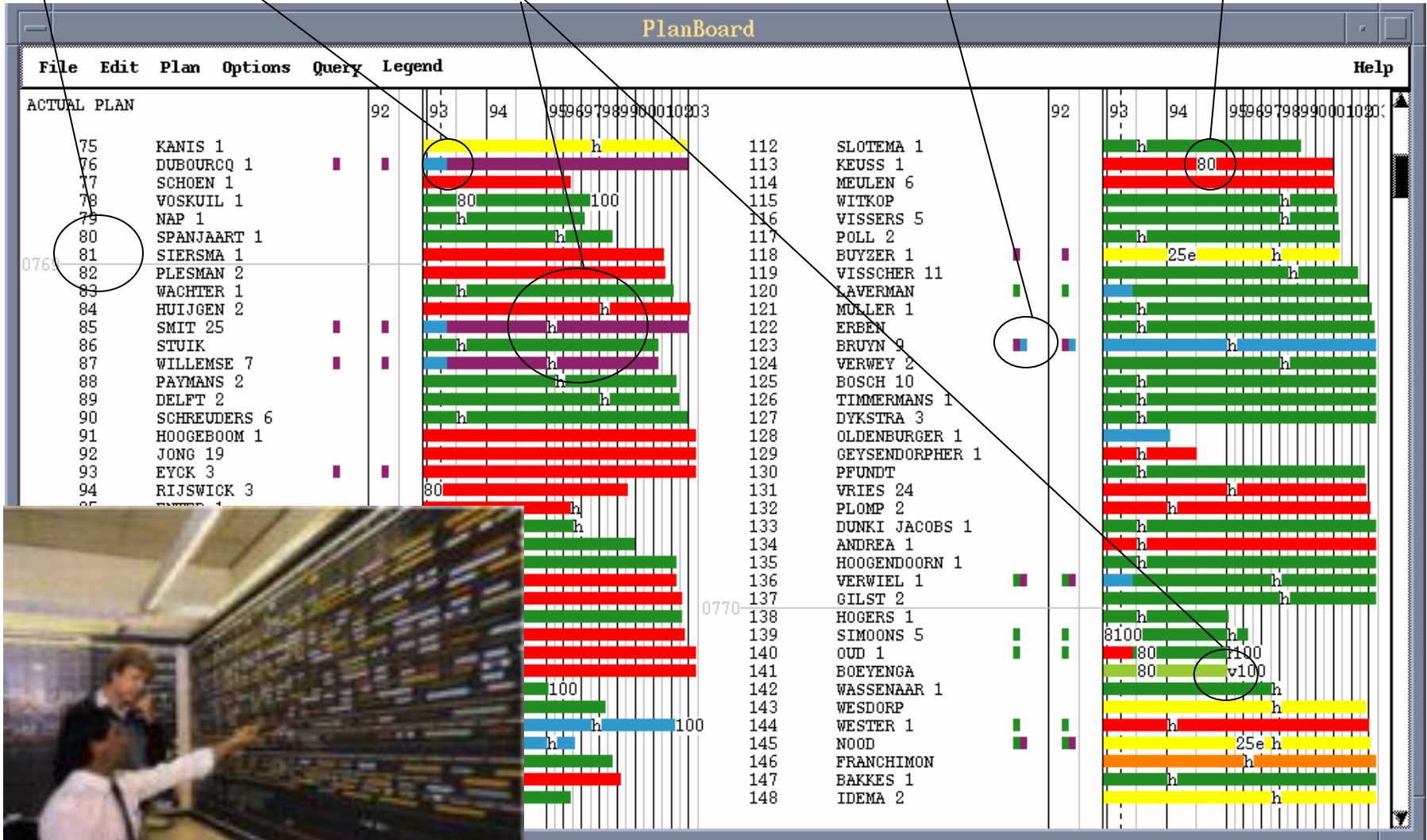
Seniority

Transitions

Horizontal/Vertical binding

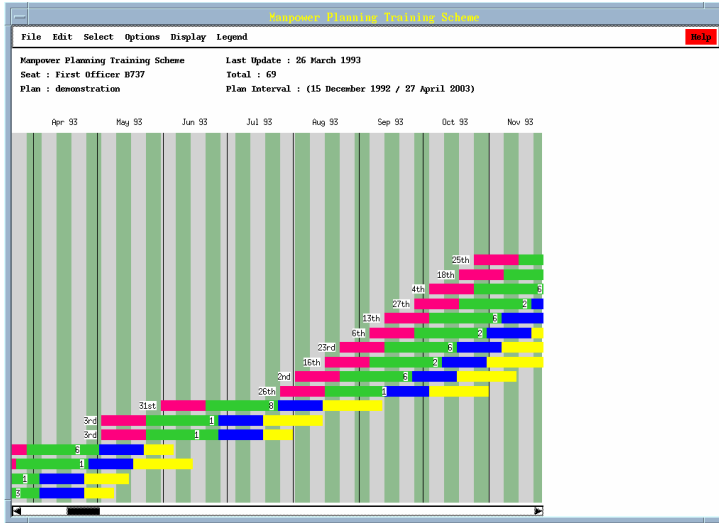
Pilot bids (preferences)

Part Time

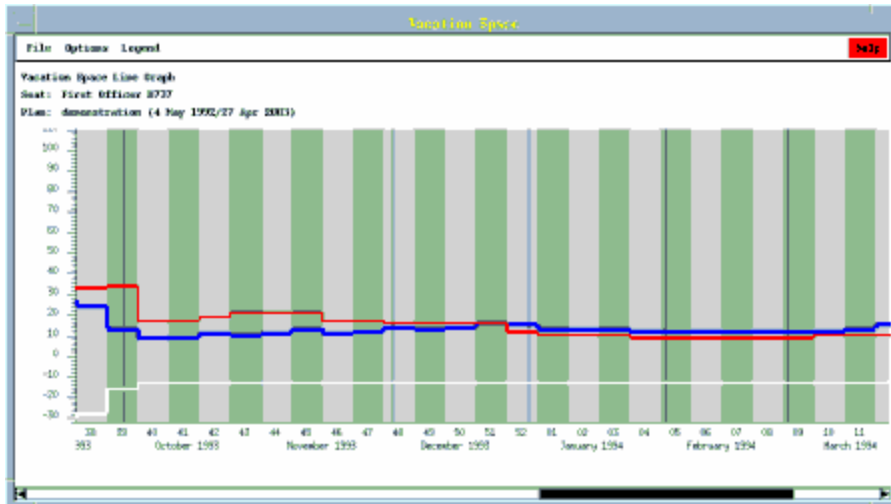


Manpower planning

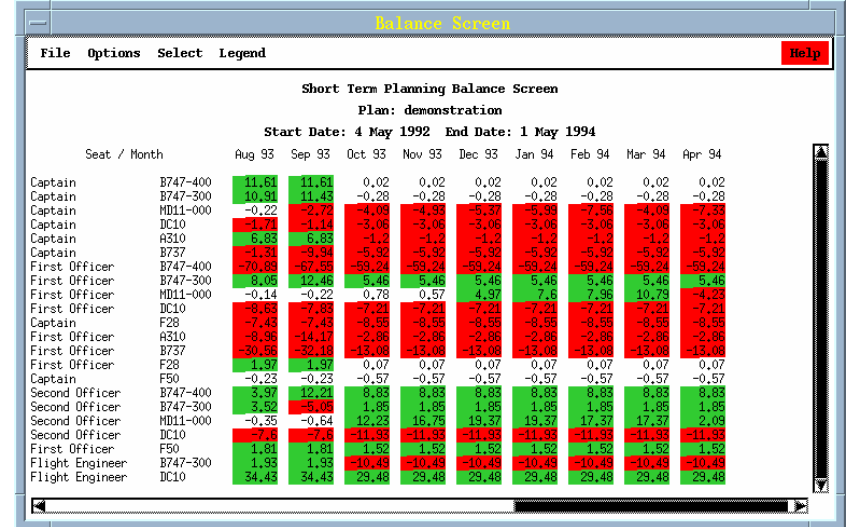
Manpower Planning



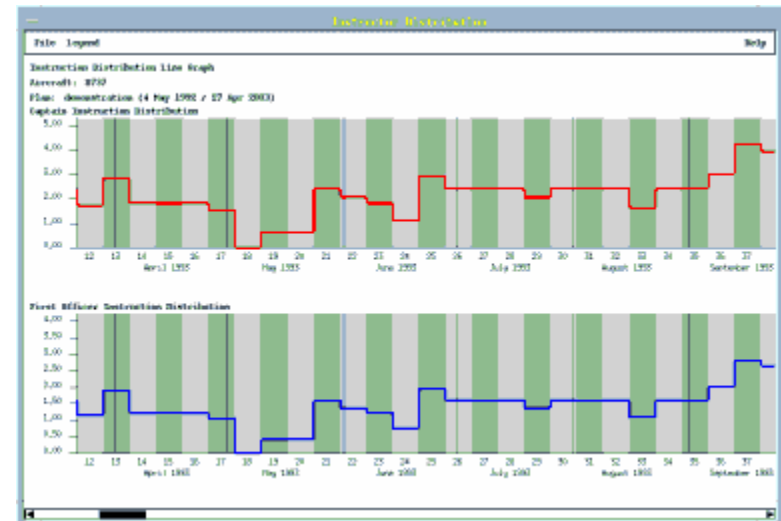
Vacation Distribution



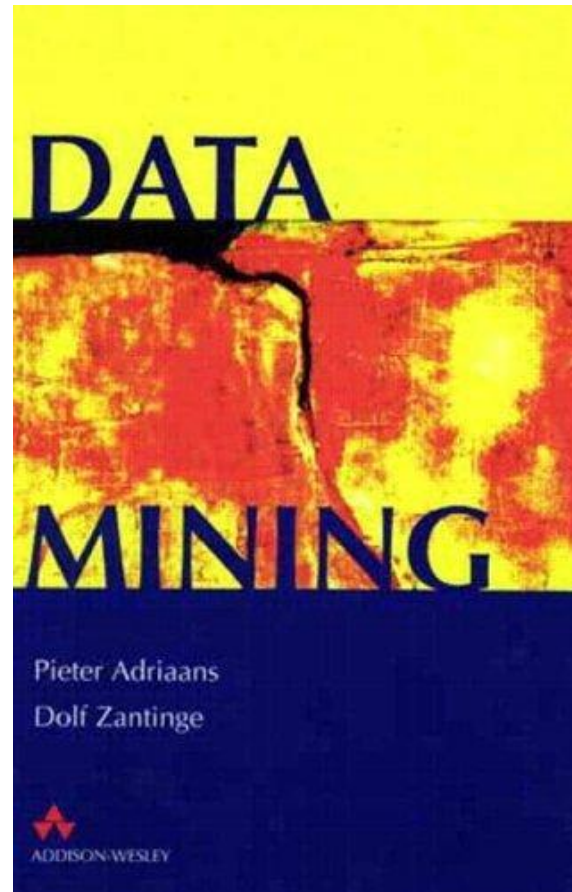
Seat Survey



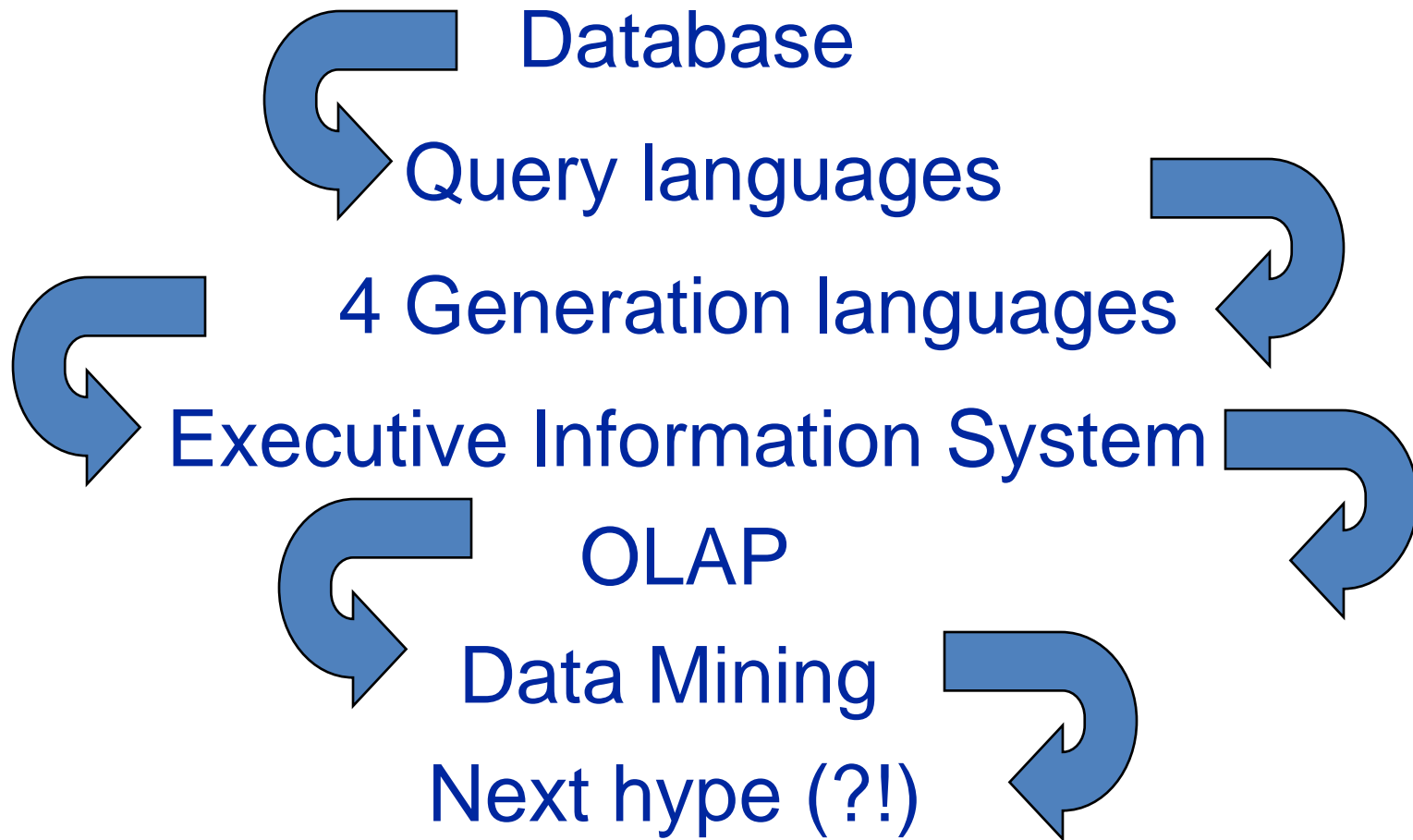
Instruction Distribution



1996



The ancestors of Data Mining

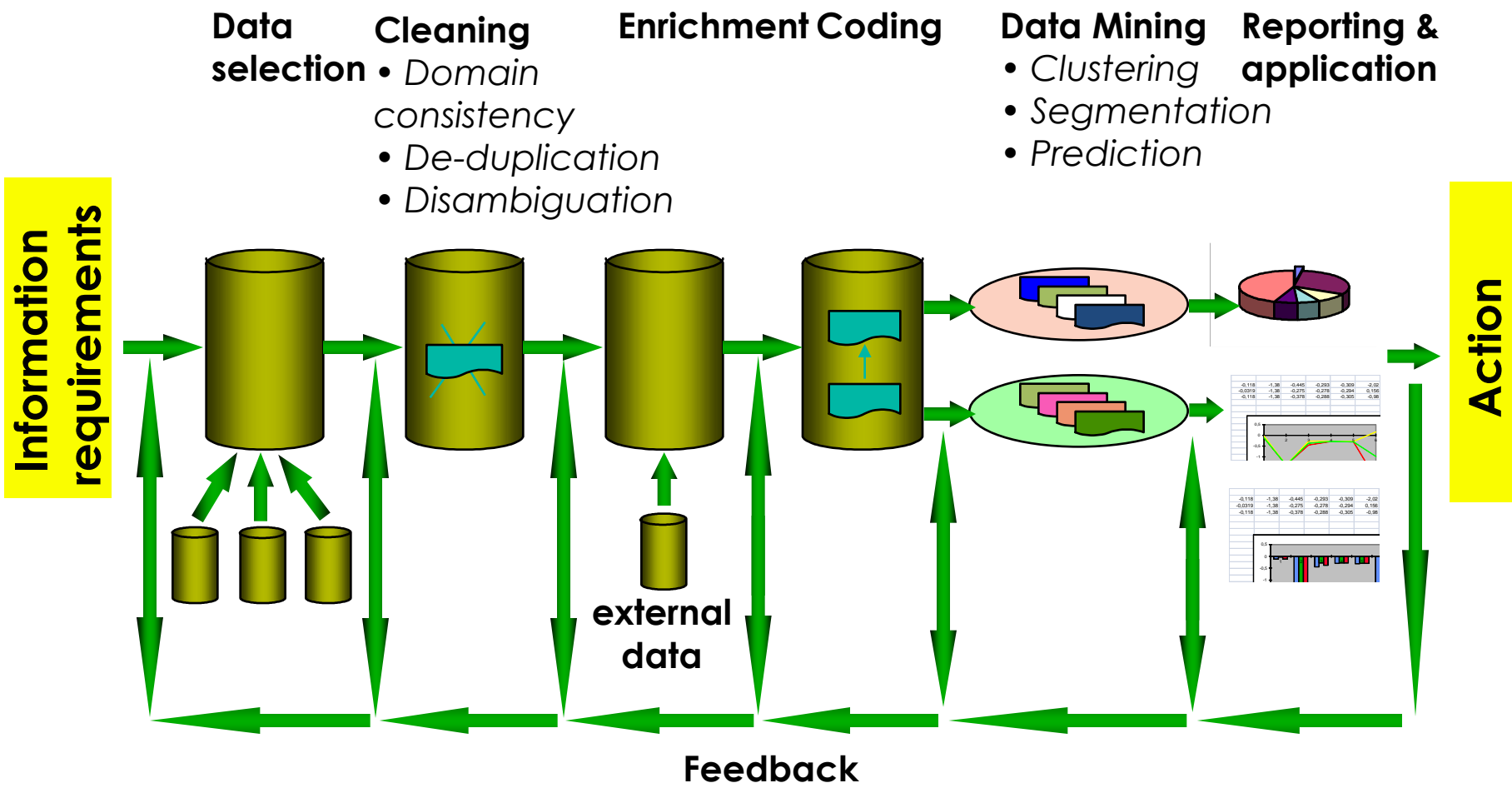


KDD Definition

Knowledge **D**iscovery in **D**atabases is the non trivial extraction of implicit, previously unknown and potentially useful knowledge from data

(after Frawley Et al. 1991)

The KDD process



1997: Adaptive System Management

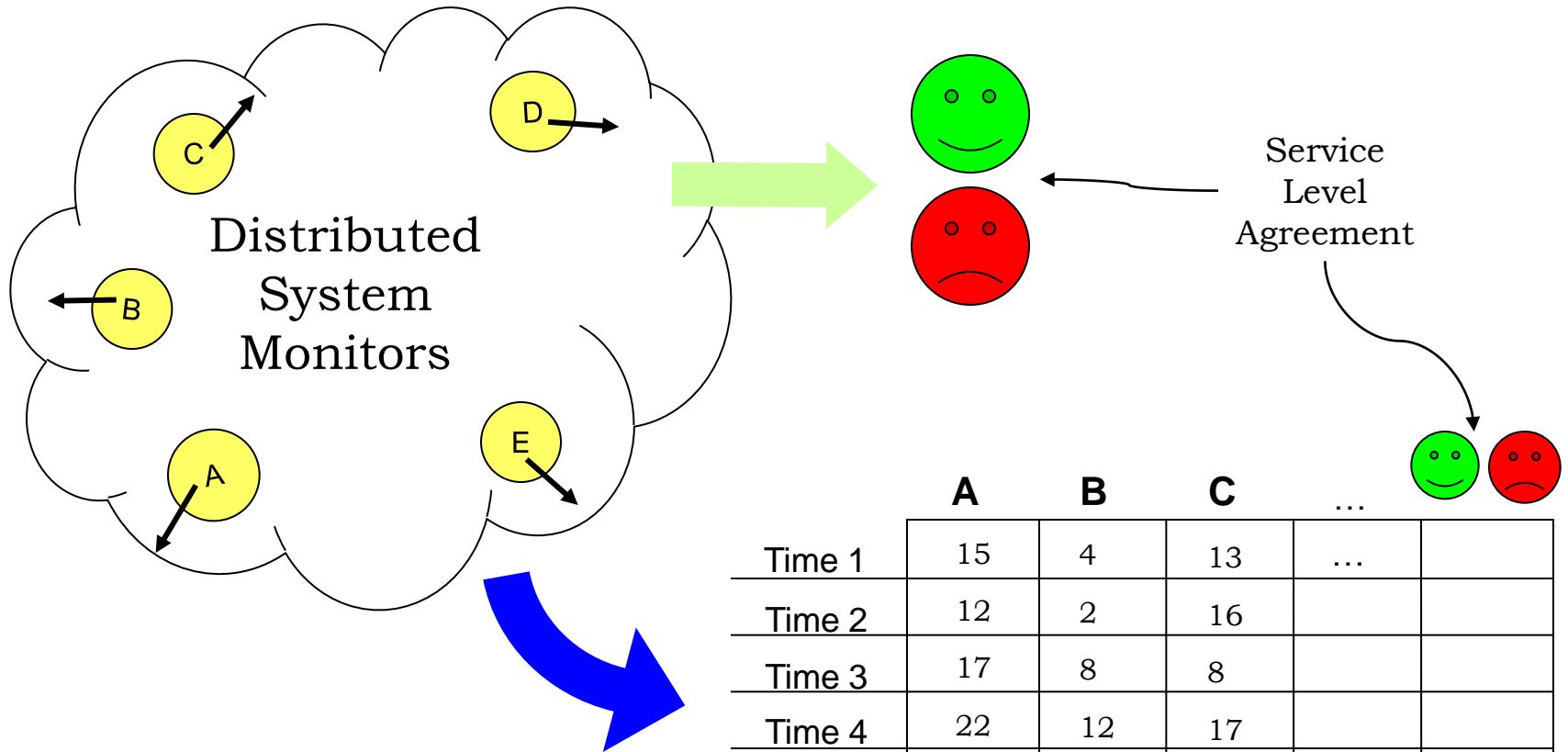
- the use of machine learning and data mining techniques to create self-learning models of the constantly changing IT-environment that allow us to predict the behavior of systems in the future and take timely automated remedial action to prevent system failure or decreased availability of the system.

Analysis of IT infrastructures using knowledge discovery

The image displays several windows from the Adaptive Systems Management (ASM) software suite, illustrating the analysis of IT infrastructures using knowledge discovery.



- Adaptive Systems Management (ASM):** The main interface window at the top, showing menu options (File, ASM, Help) and tool icons for Model Builder, Continuity Analysis, Monitor Mining, Application Mining, and Reconnect. It is currently displaying 'Strain Analysis: Flight 248'.
- ASM Regression Window:** A scatter plot titled 'Straingauge outerwing and Vertical acceleration in aircraft assenstelsel'. The y-axis is 'Vertical acceleration in aircraft assenstelsel' (ranging from 46.0 to 585.0) and the x-axis is 'Straingauge' (ranging from -356.0 to 0). A regression line is shown with the equation $Y = -0.40556124 * X + 124.453964$ and a correlation of -0.60042584 . Buttons for 'Calculate', 'Rank Y', 'Close', and 'Help...' are visible.
- Continuity Analysis - Strain Analysis: Flight 248:** A window showing a time-series plot of '25:Straingauge outerwing - Flight 248' with a target (SLA) of -200 . The plot shows values fluctuating around zero over time. A date stamp '4/14/70 11:41 PM' is present. A blue text overlay reads 'asing trend, slope = -4.8118277E-6' and 'section: 15-Mar-71 7:54:56 PM'. A clock shows '8:45 AM'.
- Straingauge outerwing Decision Tree:** A decision tree diagram for 'Straingauge outerwing'. The root node is 'Straingauge outerwing <= -200.0' (Aircraft for 248, Propensity: 2.0, Support: 6948). It branches into two main paths:
 - Vertical acceleration in aircraft assenstelsel > 299.0** (Aircraft for 248, Propensity: 78.7, Support: 230):
 - Speed relative to air > 4055.0 (Aircraft for 248, Propensity: 100.0, Support: 176)
 - Speed relative to air <= 4055.0 (Aircraft for 248, Propensity: 0.3, Support: 54)
 - Vertical acceleration in aircraft assenstelsel <= 299.0** (Aircraft for 248, Propensity: 0.2, Support: 9719):
 - Pitch Rate > 161.0 (Aircraft for 248, Propensity: 7.8, Support: 276)
 - Pitch Rate <= 161.0 (Aircraft for 248, Propensity: 0.0, Support: 9442)

Monitoring in the system

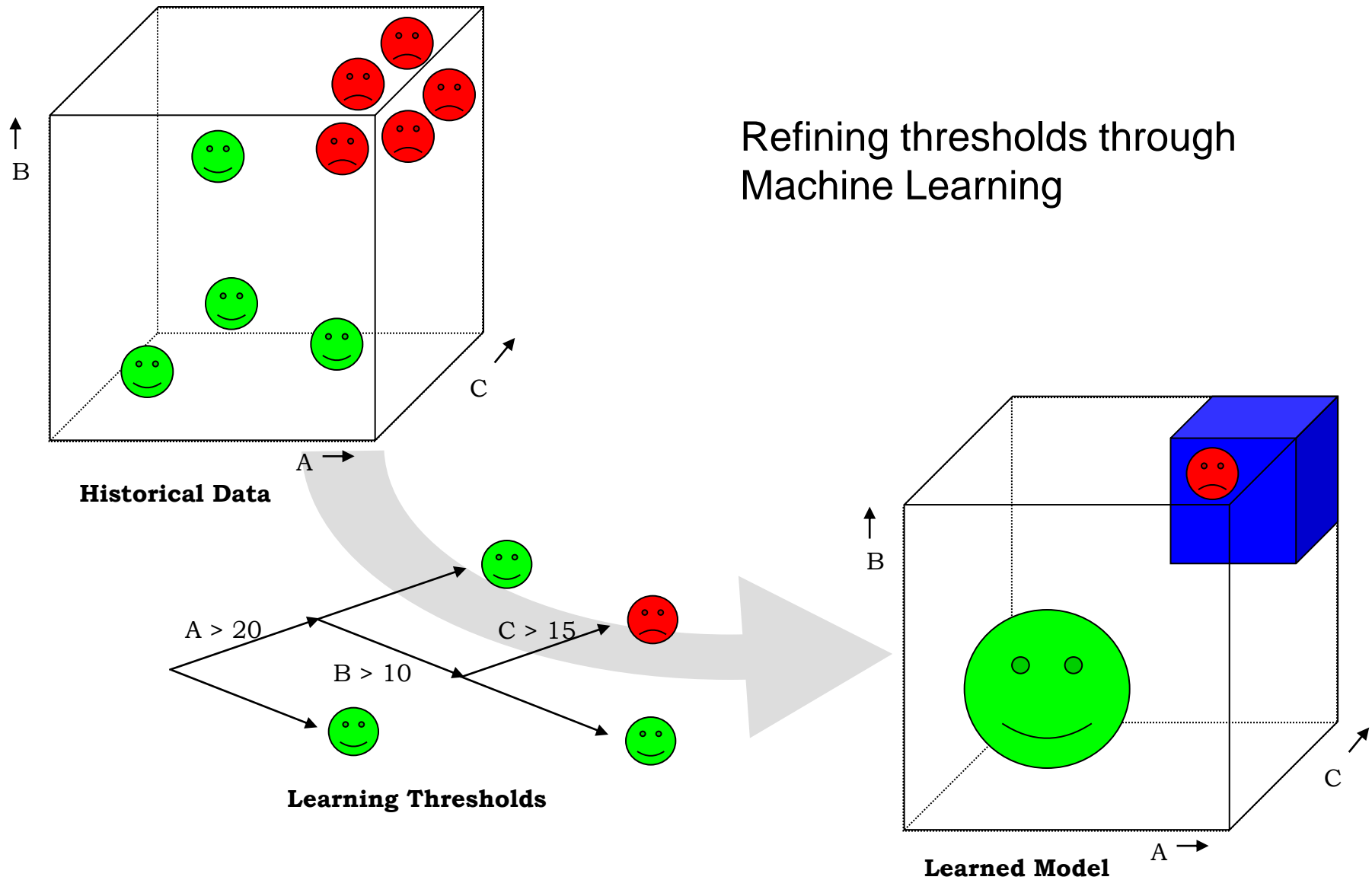


Monitor Legend

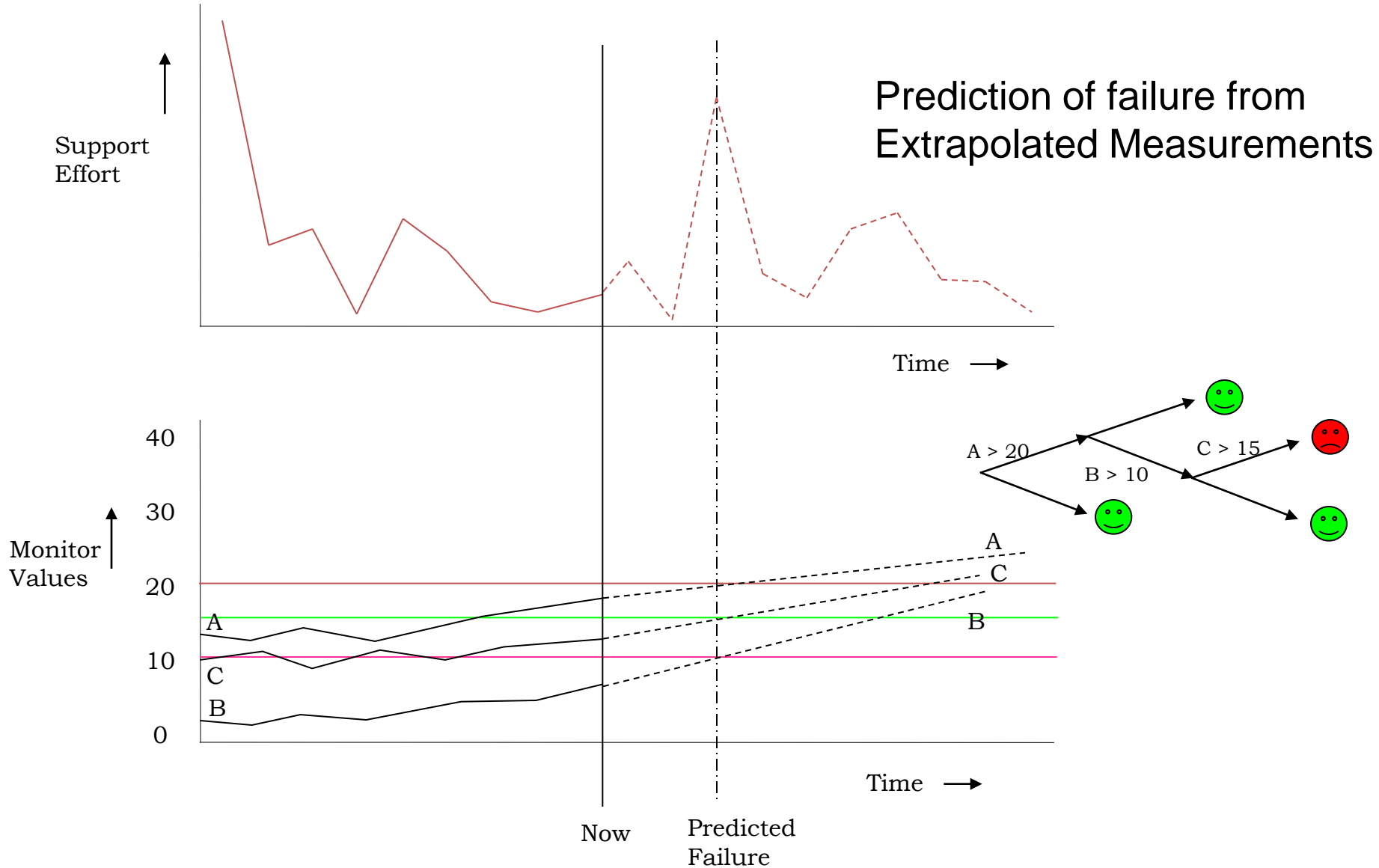
- A: Nr. of Processes
- B: Nr. of Failed Jobs
- C: Login time Application
- D: System Load
- E: CPU Idle time
- etc.

	A	B	C	...		
Time 1	15	4	13	...		
Time 2	12	2	16			
Time 3	17	8	8			
Time 4	22	12	17			
...			
Time k						

Refining thresholds



Prediction of failure

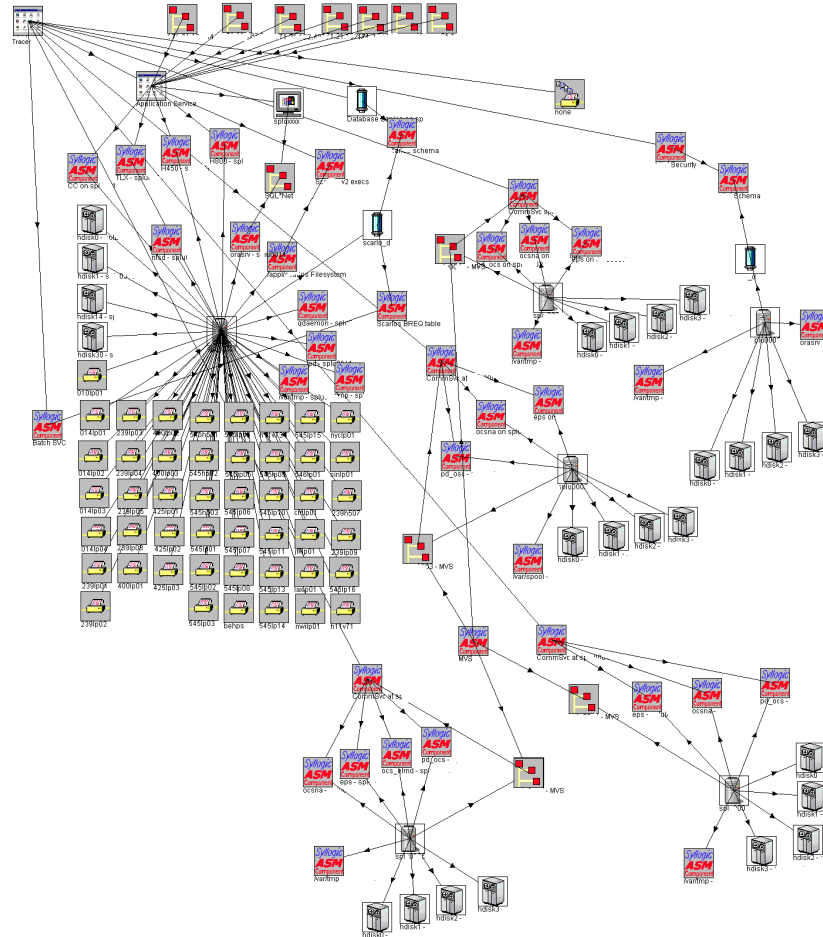


Patents

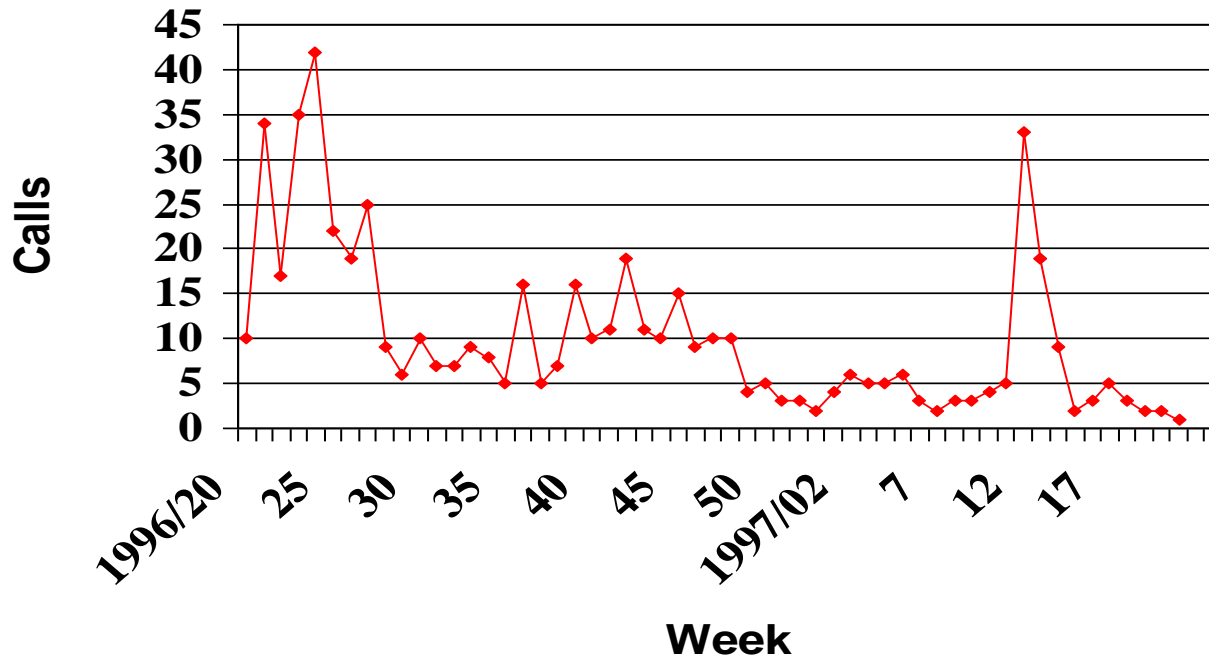
- Adriaans et al. United States Patent US 6,311,175 B1
System and method for generating performance models of complex information technology systems.
- Adriaans et al. United States Patent US 6,313,390 B1
A method for automatically controlling electronic musical devices by means of real-time construction and search of a multi-level data structure.
- Adriaans et al. United States Patent US 6,393,387 B1
System and method for model mining complex information management systems.

Experiments

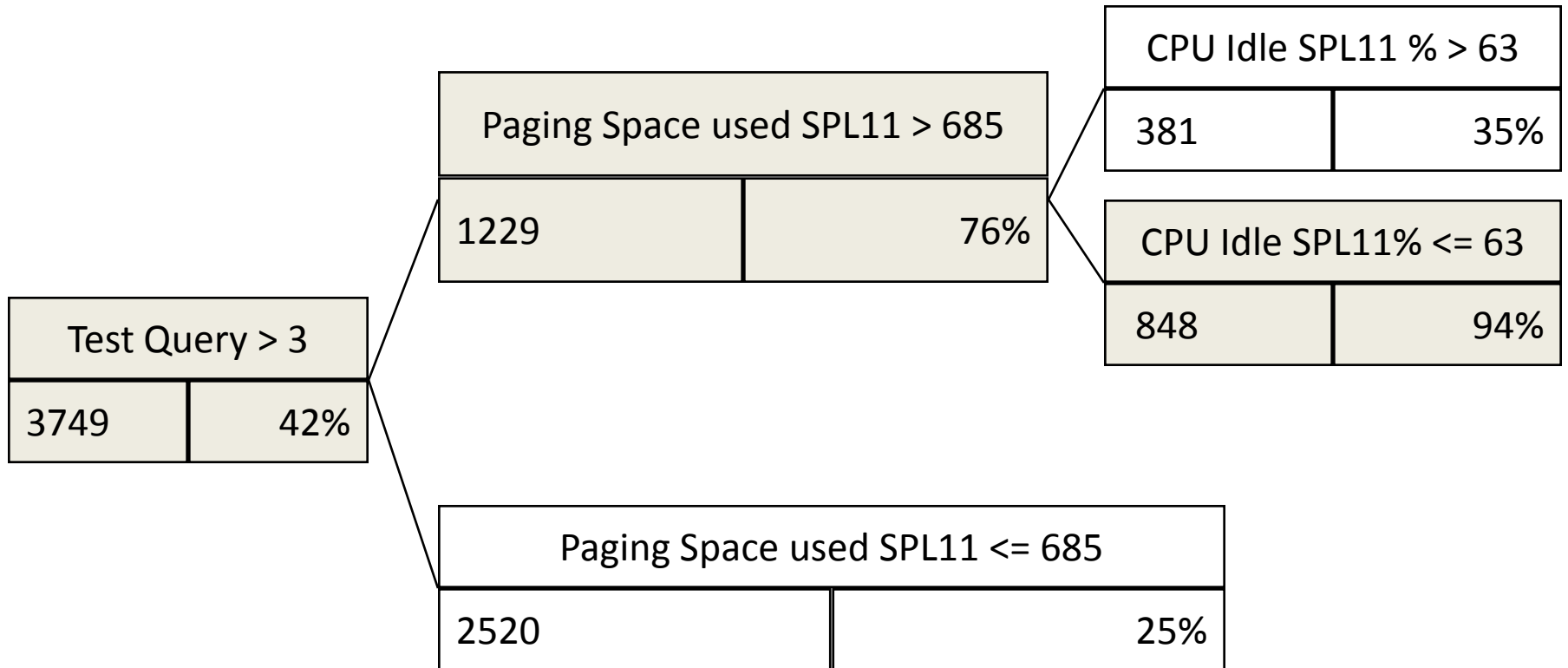
- 140 components
- 250 monitors
- 2 months
- 3500 snapshots



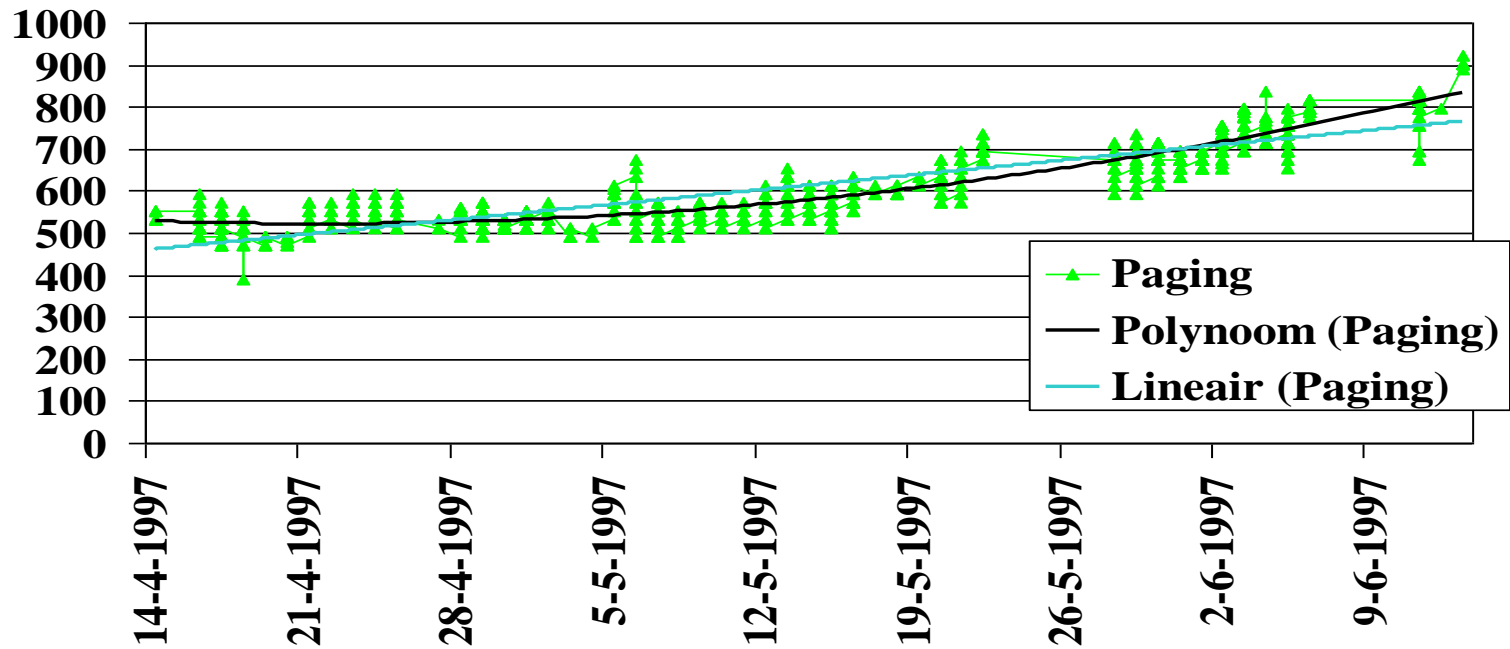
Helpdesk calls application



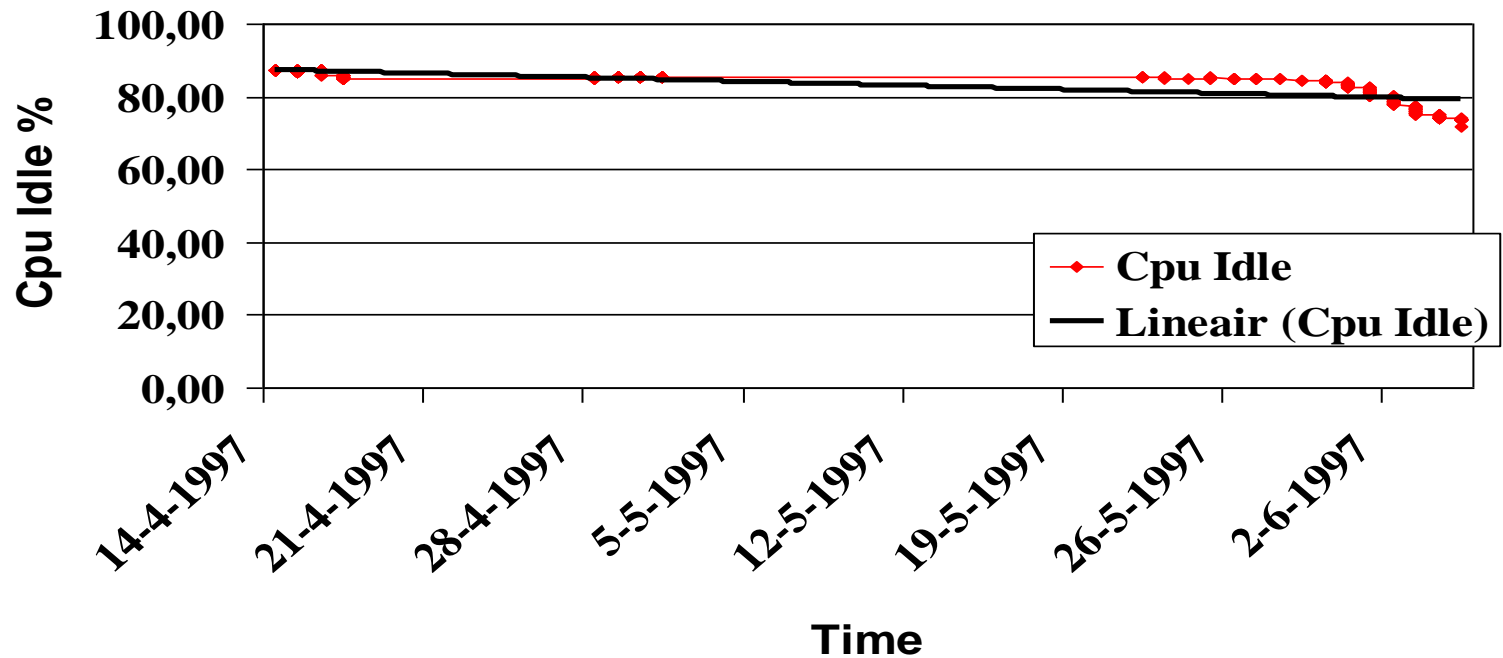
Decision tree Scarlos performance




Paging space develops polynomially



CPU Idle Trend



1998-2005

The logo features the word "ROBOsail" in a stylized font. "ROBO" is in a dark blue, rounded, sans-serif typeface, while "sail" is in a lighter blue, more fluid, rounded font. The text is set against a background of a light blue globe with a white grid of latitude and longitude lines. A thick, wavy line in two shades of blue (light and dark) runs horizontally across the bottom of the text, resembling a wave or a stylized underline.

ROBOsail



Pole Balancing

State:

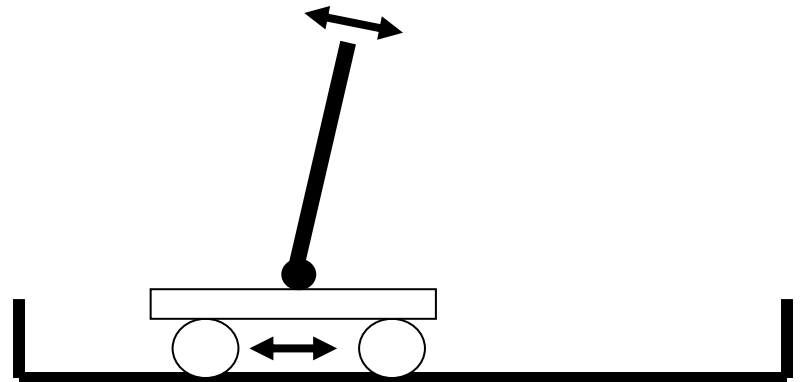
- Velocity of the cart
- Position of the cart
- Position of the stick
- Velocity of the stick

Action:

- Force on cart

Reward:

Stability of the pole



Ship Balancing

State:

- Velocity of the rudder
- Position of the rudder
- Course
- Log speed
- Apparent Wind speed
- Apparent wind angle
- Heel
- etc...

Action:

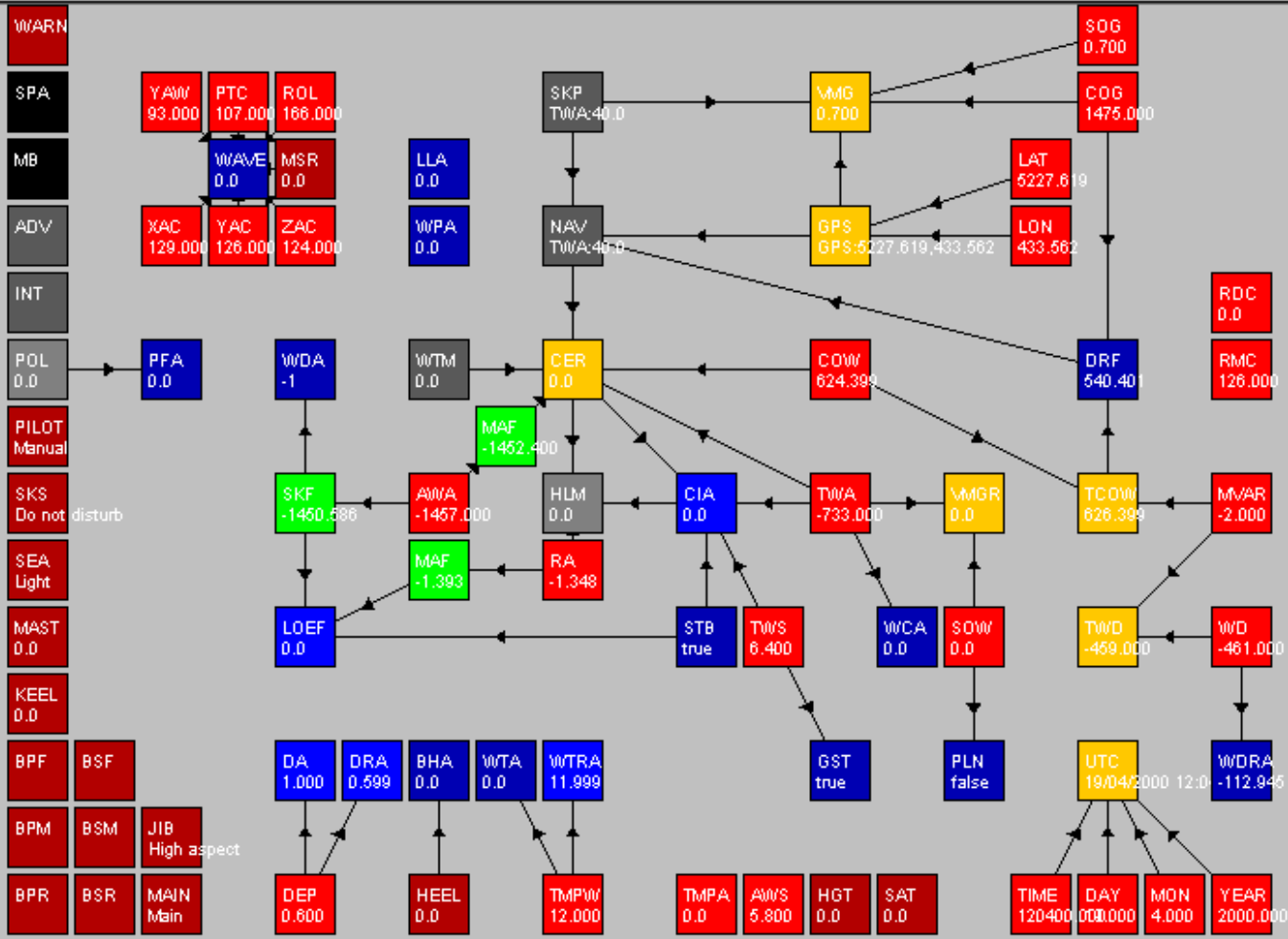
- Force Exercised on the rudder

Reward:

- Velocity of the yacht



12:04:00



Mon Mar 20 11:56:16 CET 2000, 229kb, Description

	View	Compress	Delete	Import	Export	Playback	Scim			
R	BSR	PTC	ROL	YAW	XAC	YAC	ZAC	RMC	RDC	PL
								125.0	0.0	
		129.0	114.0	127.0	126.0	131.0	121.0	125.0	0.0	
		129.0	114.0	127.0	130.0	128.0	117.0	125.0	0.0	
		129.0	115.0	127.0	129.0	127.0	117.0			
								125.0	0.0	
		131.0	115.0	127.0	128.0	126.0	116.0			
								125.0	0.0	
		128.0	114.0	127.0	127.0	127.0	118.0	124.0	0.0	
		128.0	114.0	127.0	128.0	124.0	119.0	124.0	0.0	
		129.0	115.0	127.0	126.0	125.0	118.0	125.0	0.0	
		127.0	115.0	127.0	124.0	122.0	116.0			
								125.0	0.0	
		129.0	115.0	127.0	125.0	125.0	124.0			
								125.0	0.0	

```

----- SYS: Program mode: offline
----- HELM: Helmsman timer on 250 ms
----- INIT: Initializing movingAverageFilter #3, MAF
----- INIT: Initializing simpleKalmanFilter #4, SKF
----- INIT: Initializing movingAverageFilter #2, MAF
----- INIT: Initializing ModelBuilder, MB
----- INIT: Initializing Statespace, SPA
    
```



GPS Waypoints

Lon 00' 00" 000 N

Lat 00' 00" 000 E

Set

Ping

Sail selection

[Rif 1 (R1)]

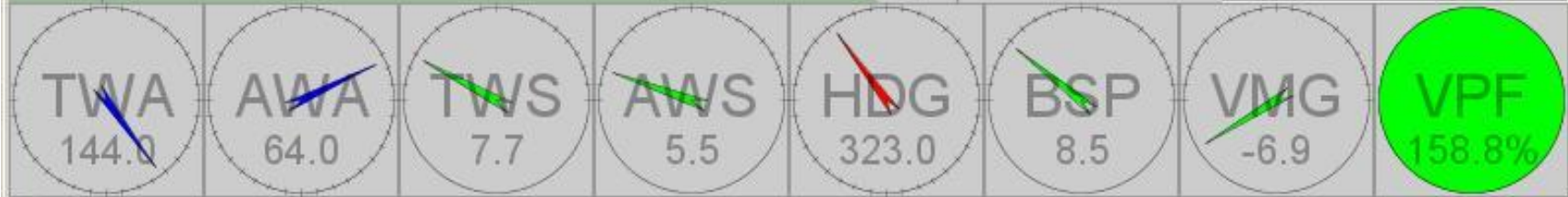
[KF;Kleine Fok]

[H;Halfwinder]

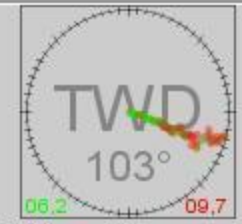
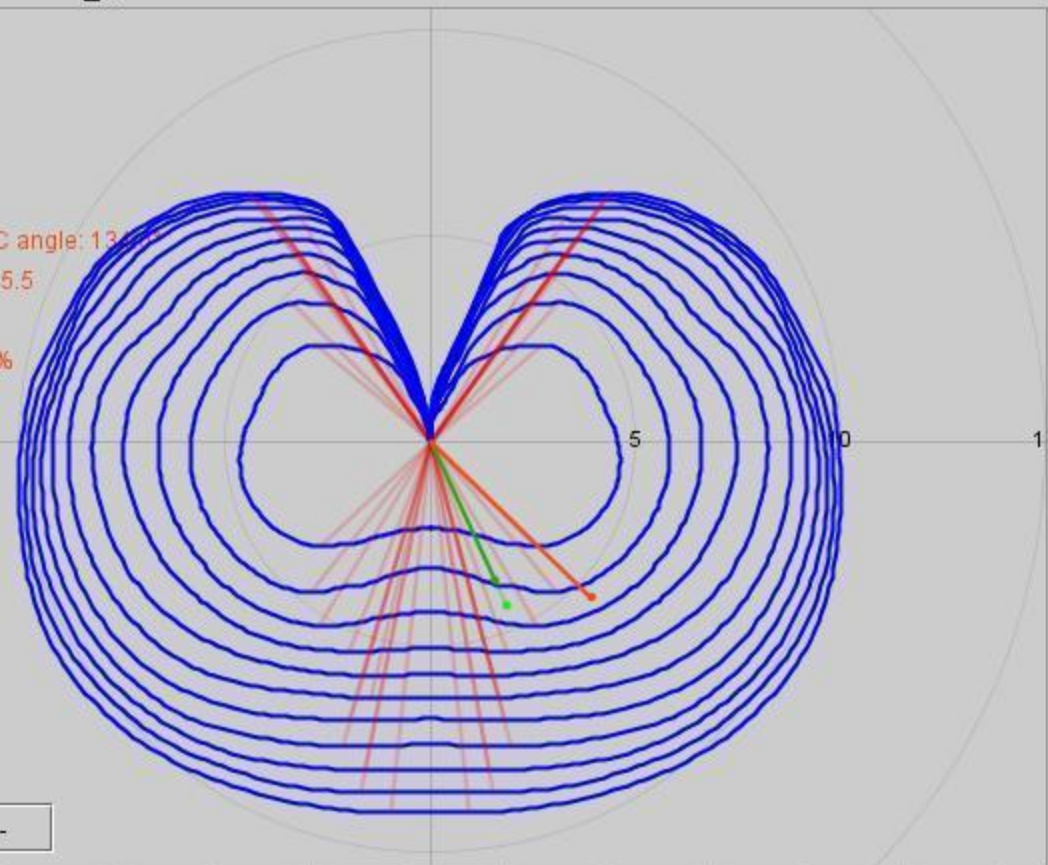
Timer

0:01:30

▶ || □



TWA: 155.0°
 TWS: 7.0
 BSP: 3.8
 BSPP: 4.4
 VPF: 85%
 Suggested VMC angle: 13°
 Expected BSP: 5.5
 VMC: 5.1
 VMC gain: +36%



GPS Waypoints

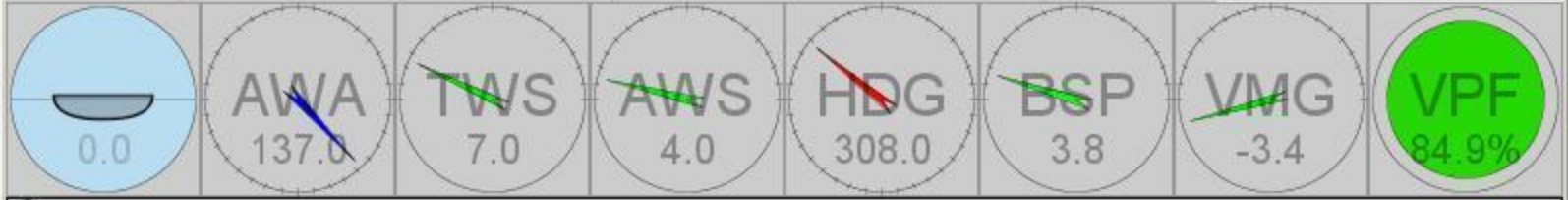
Lon N ▾
 Lat E ▾

Sail selection

[Rif 1 (R1)] ▾
 [GF; Grote Fok] ▾
 [K; Kluiver] ▾

Timer

0:00:25



File View Options Help

config

view

settings

auto manual

optimal

min/max

control

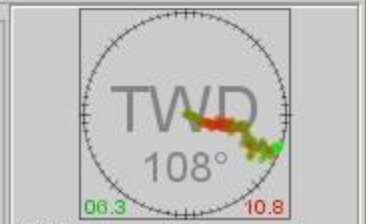
bsp	eff
4.1	0.0
tws	twa
9.1	58.8
aws	awa
12.2	38.2

advisor 5.8kn 44.0%

M

GF

K



GPS Waypoints

Lon 00' 00" 000 N

Lat 00' 00" 000 E

Set

Ping

Sail selection

Mainsail (M)

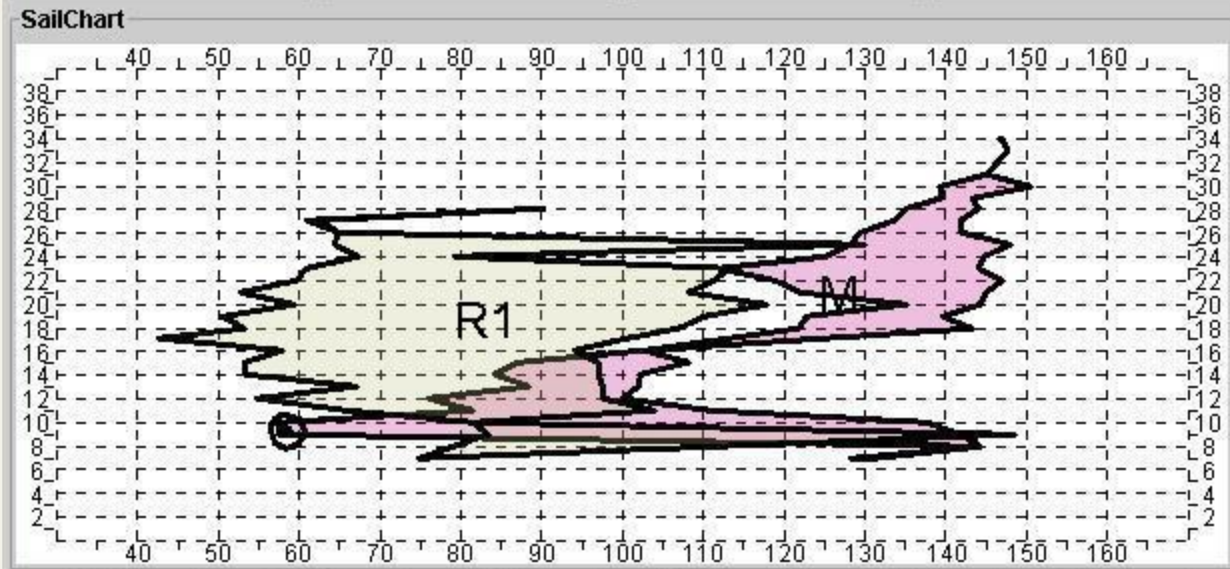
Kleine Fok (KF)

Kluiver (K)

Timer

0:00:00

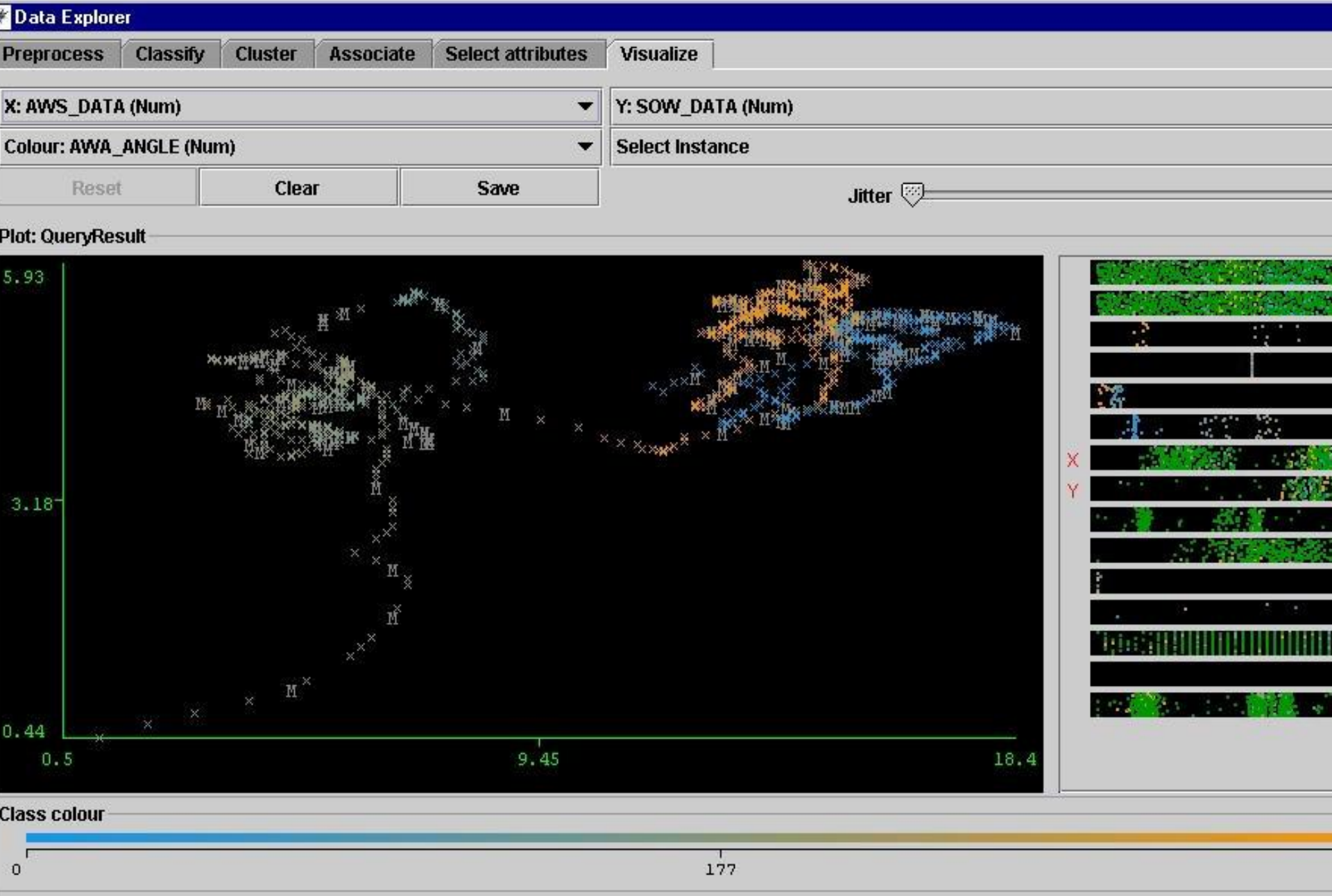
▶ || □



TWA 52.7	AWA 35.0	TWS 10.8	AWS 13.9	HDG 55.0	BSP 4.50	VMG 2.7	VPF 68.9%
-------------	-------------	-------------	-------------	-------------	-------------	------------	--------------



Mining for Sensor Calibration

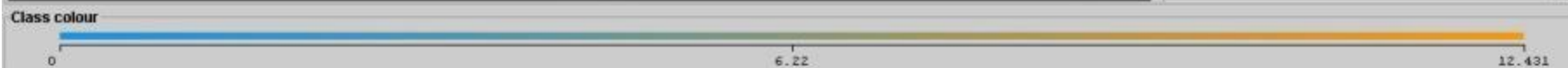
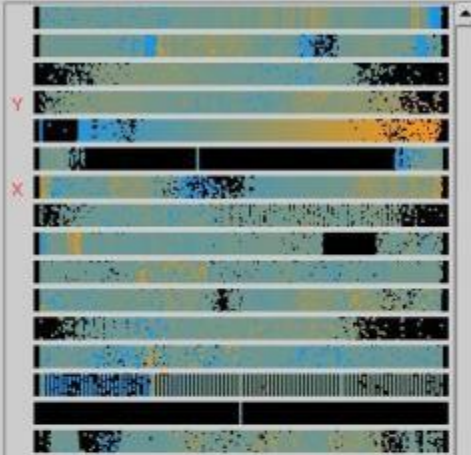
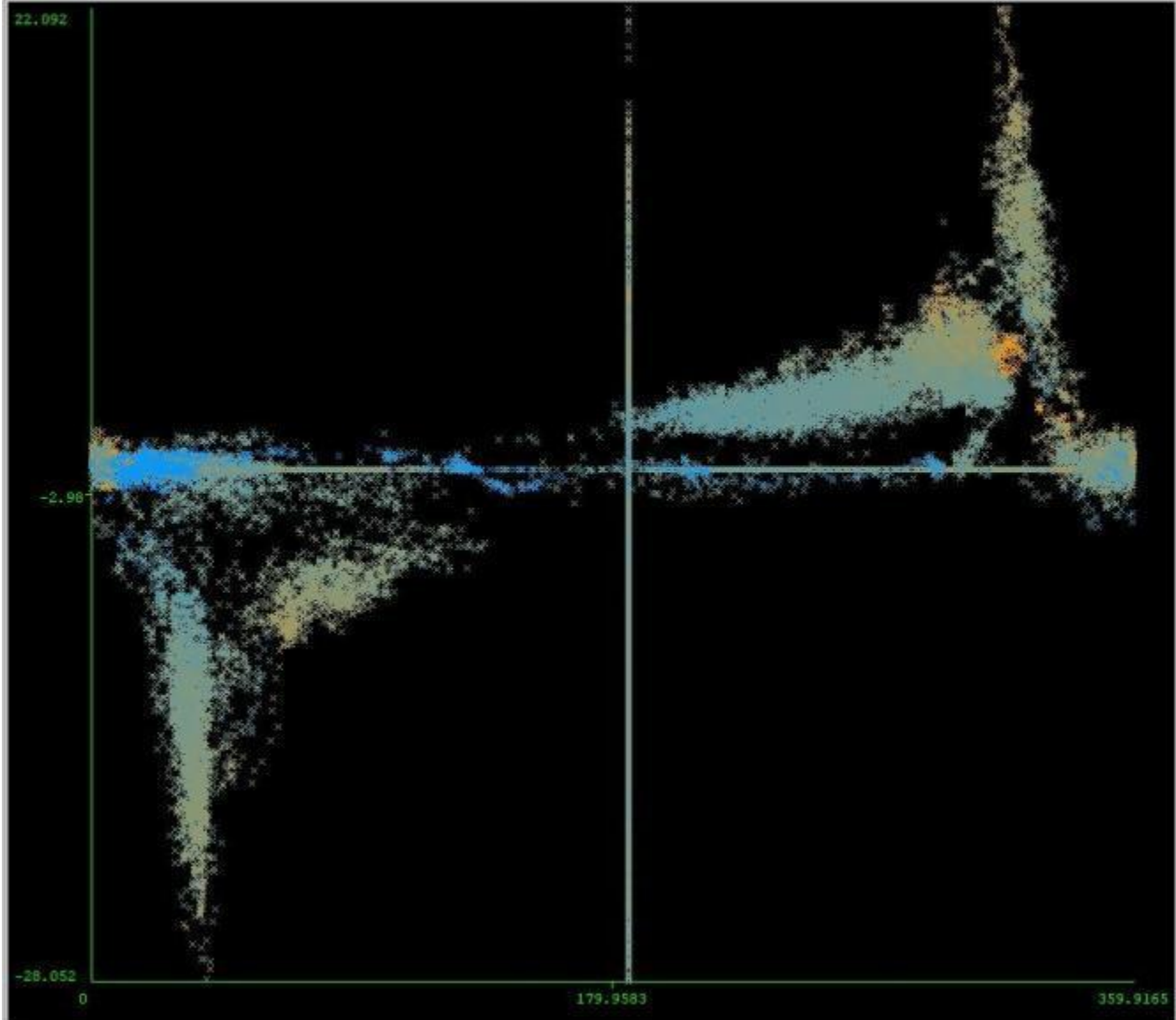


Mining for Sensor Calibration

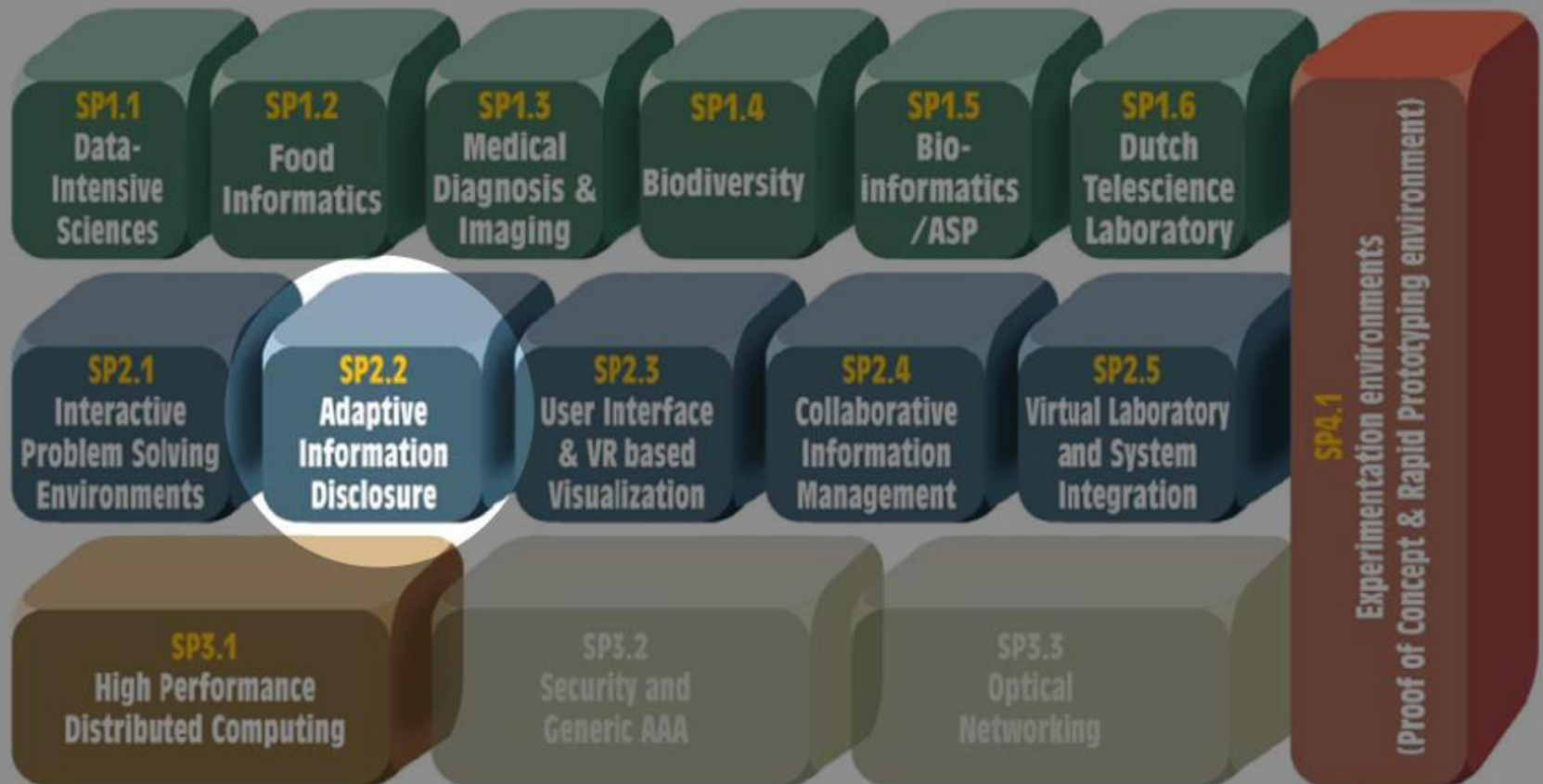
- Possible causes:
 - Log sensor gives asymmetric data
 - Wind sensor gives asymmetric data
 - Boat itself is asymmetric
 - Boats with masts > 25 m perform better over port-bow on Northern hemisphere

X: AWA_ANGLE (Num) Y: RLL_DATA (Num)
Colour: SOW_DATA (Num) Select Instance
Reset Clear Save Jitter

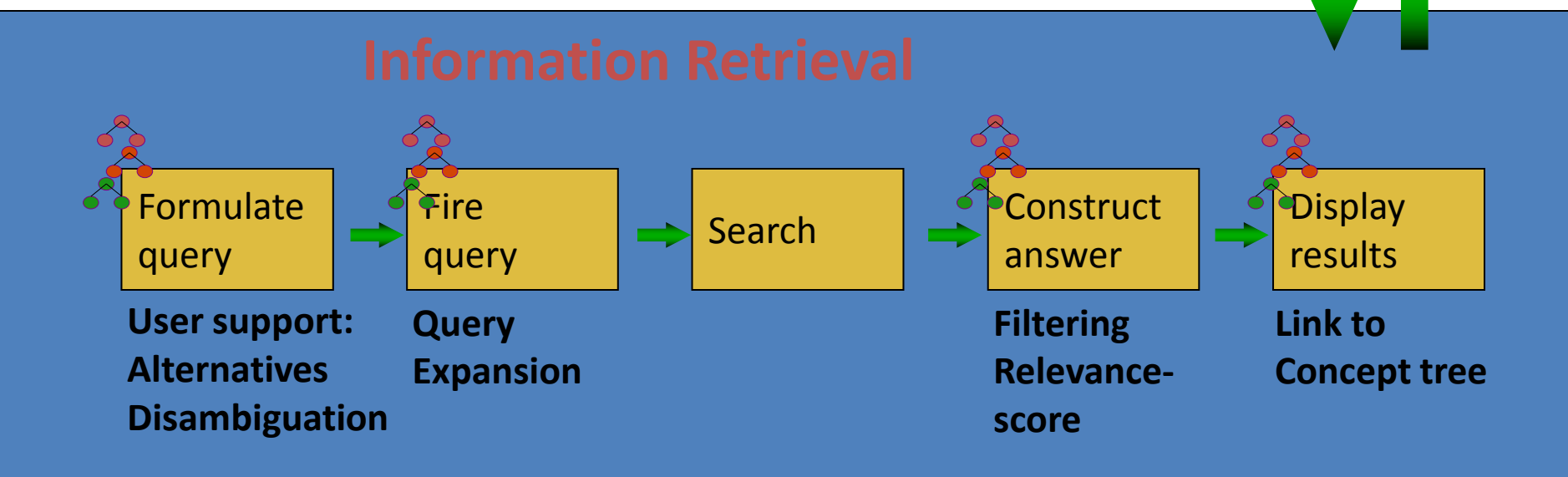
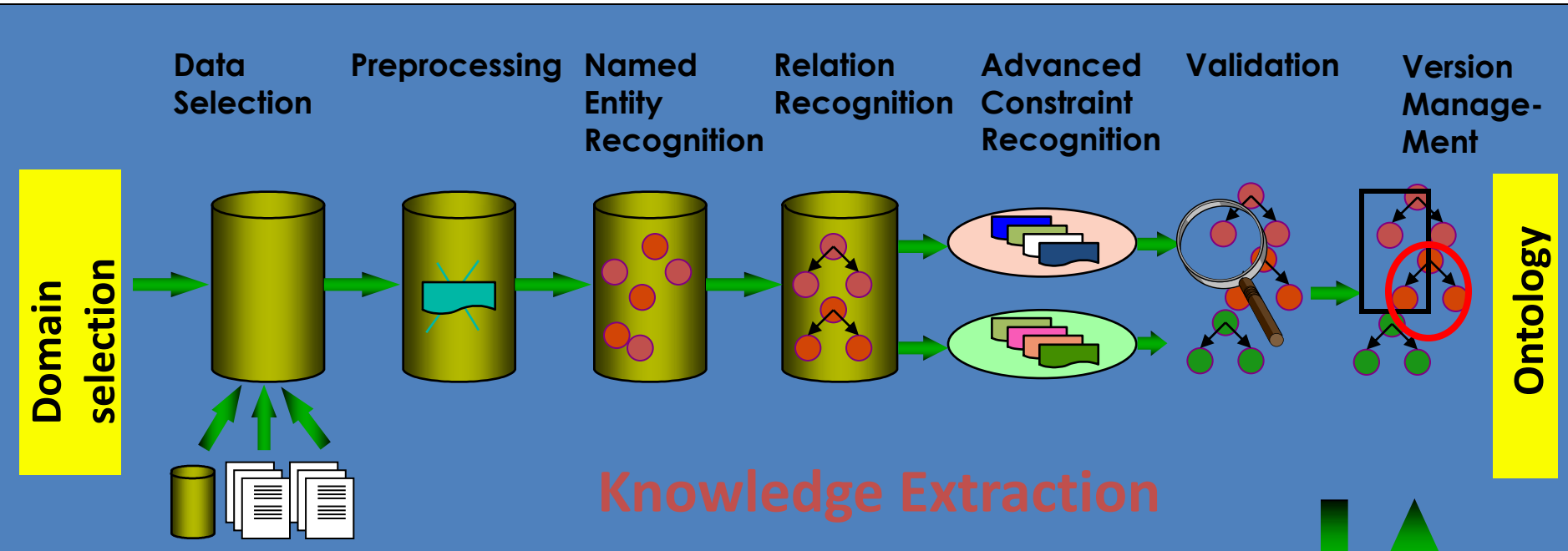
Plot: QueryResult-weka.filters.unsupervised.attribute.ReplaceMissingValues

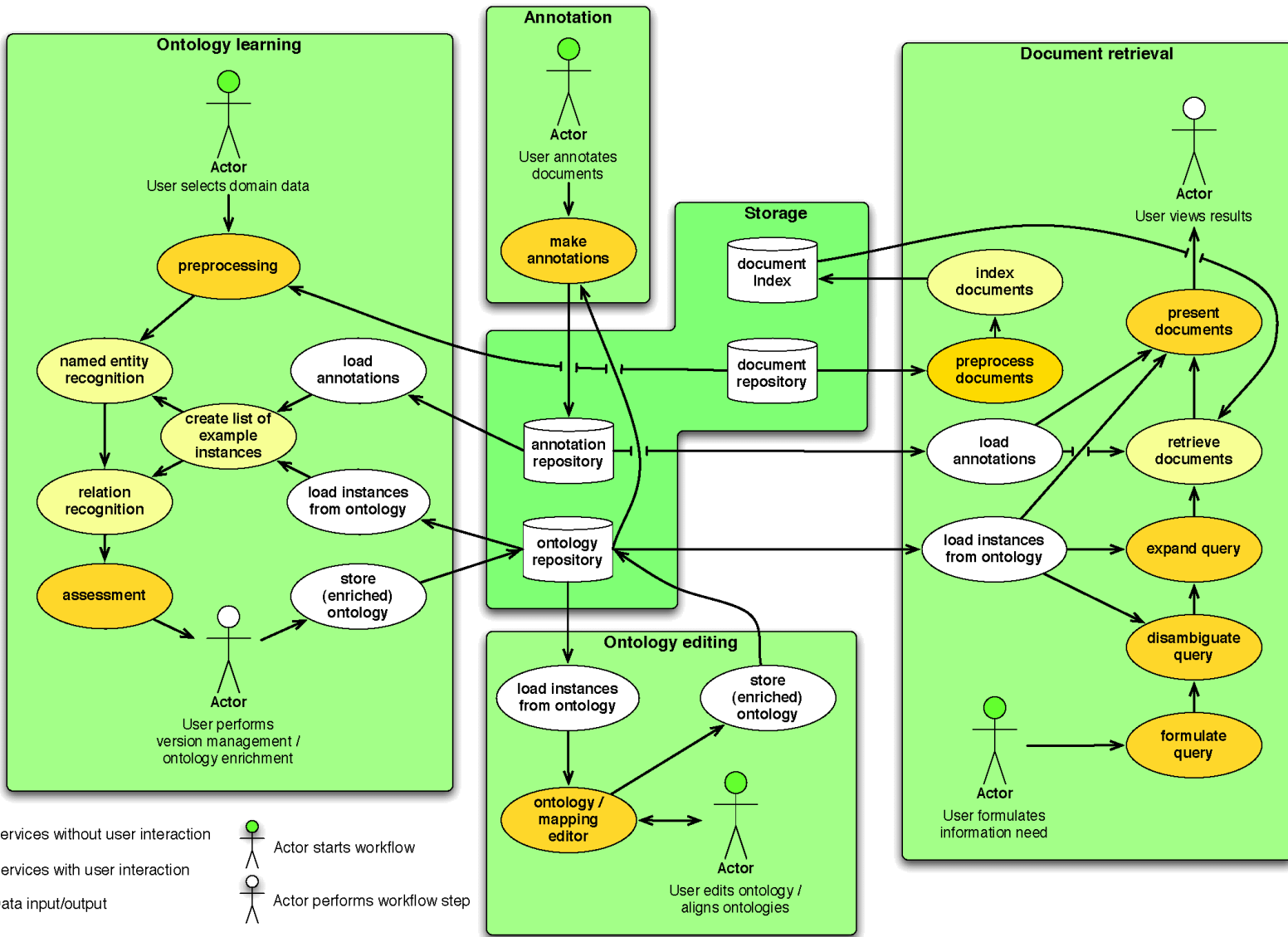


2003-2009 Adaptive Information Disclosure (AID) participating in the VL-e (Virtual lab for e-Science) project



Adaptive Information Disclosure





The AIDA toolbox
for knowledge extraction and knowledge management
in a Virtual Laboratory for e-Science

2010-2015: D2S: From Data to Semantics

for data publishers



Challenges

How to share scientific data?
How to access, analyse and interpret the data?
How to communicate and publish results?

Domains

- Life-Sciences
- Humanities



PHILIPS

End Users

- Elsevier
- Philips
- DANS



DANS

How do we speed up the transfer of scientific data, information, knowledge from a research paper into actionable form?

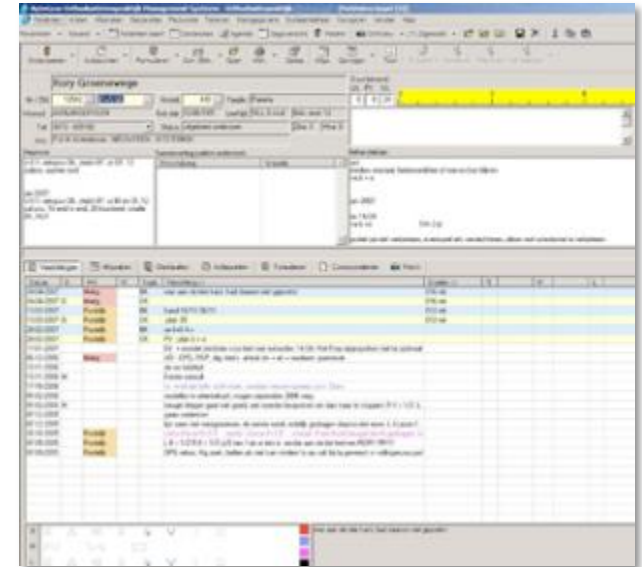
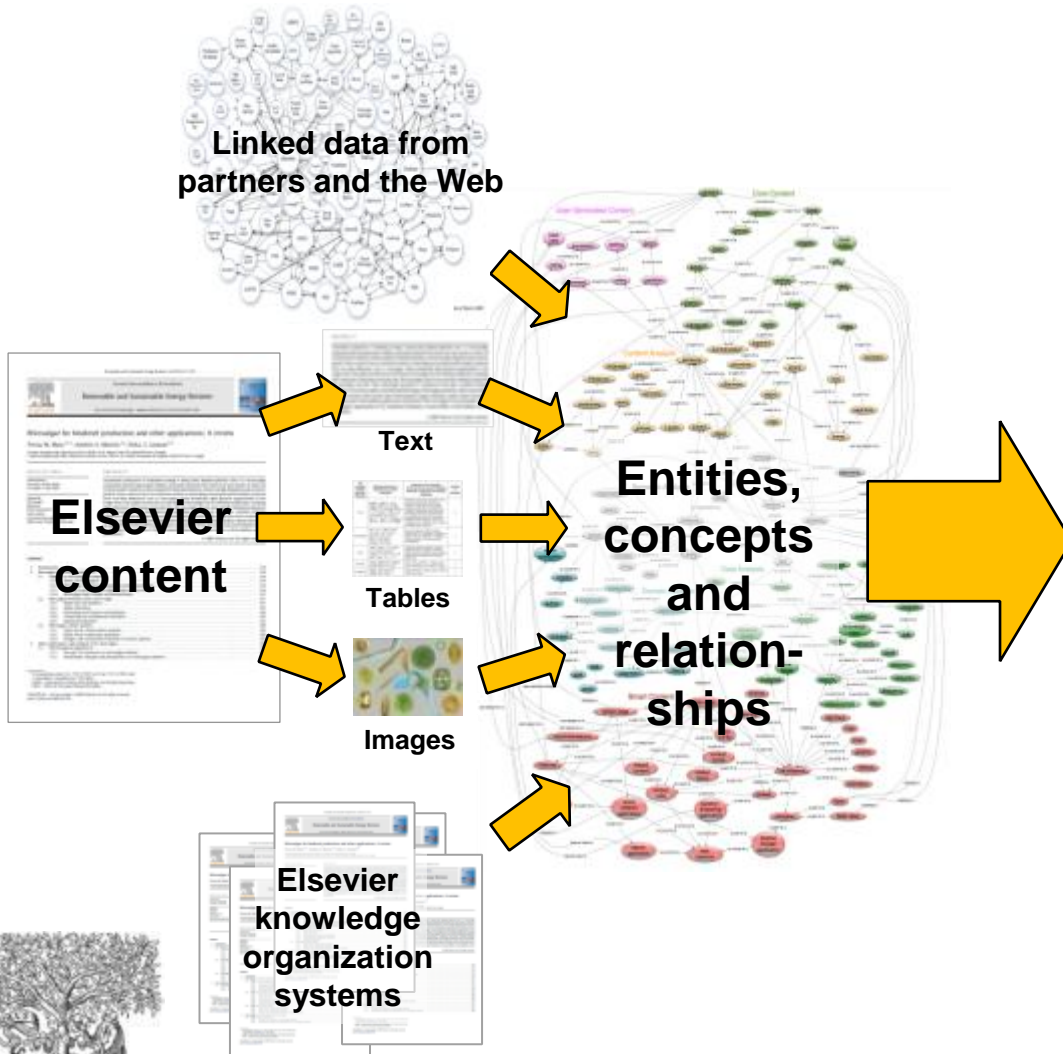


- Looking to become a provider of clinical rules
- Map clinical guidelines against recent papers
- Link to data sets of research paper



- looking to quickly consume research findings into Clinical Decision Support systems
- current guideline update cycle is 5 years

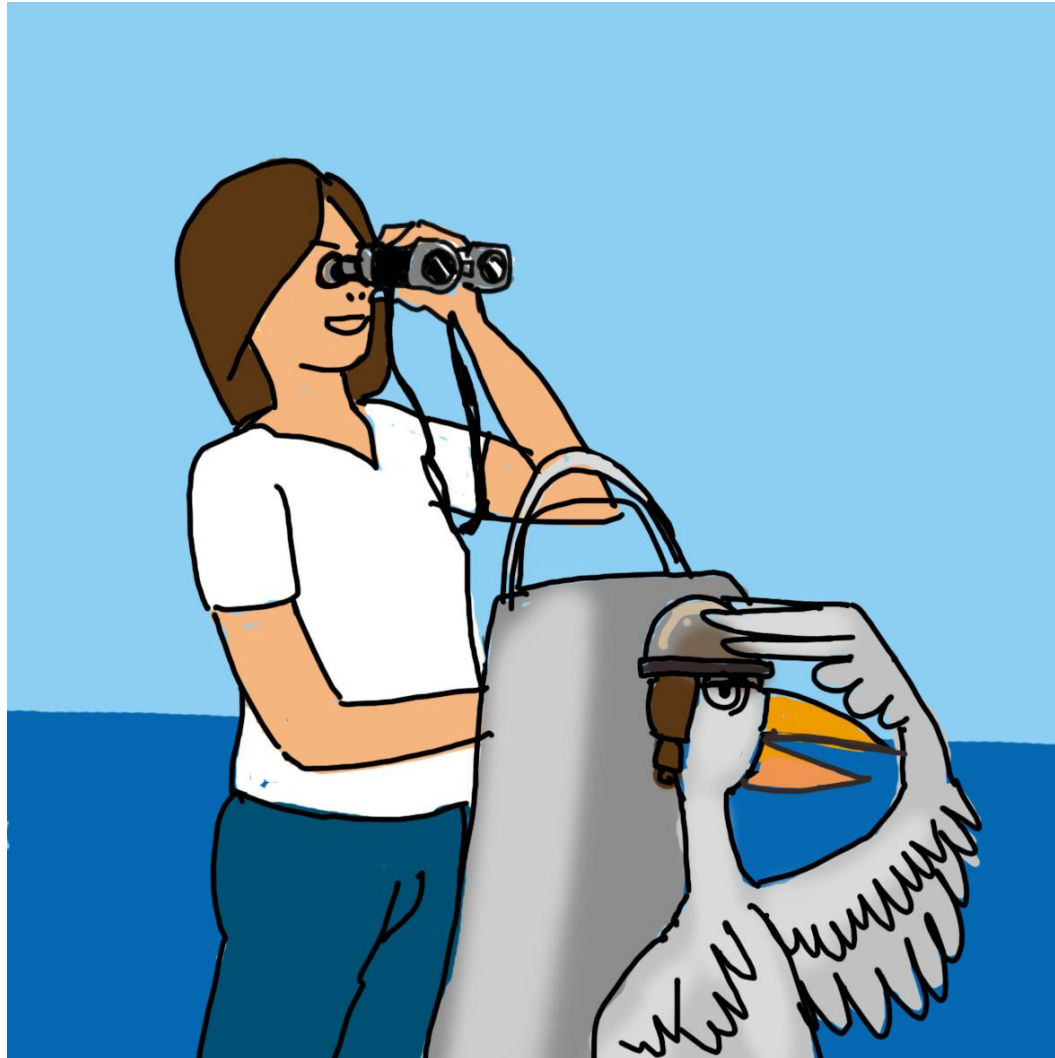
Smart content at Elsevier & Philips



Philips Hospital Information System



Looking Forward



A hand holding a magnifying glass over a target with a bullseye. The target has a red bullseye in the center, surrounded by yellow and black concentric rings. The hand is positioned on the left side of the target, holding the handle of the magnifying glass.

EXPERT OPINION

Contact Editor: **Brian Brannon**, bbrannon@computer.org

The Unreasonable Effectiveness of Data

Alon Halevy, Peter Norvig, and Fernando Pereira, Google

Eugene Wigner's article "The Unreasonable Effectiveness of Mathematics in the Natural Sciences"¹ examines why so much of physics can be neatly explained with simple mathematical formulas

such as $f = ma$ or $e = mc^2$. Meanwhile, sciences that deal with things rather than elementary particles often require elegant math-

behavior. So, this corpus could serve as the basis of a complete model for certain tasks—if only we knew how to extract the model from the data.

Learning from Text at Web Scale

The biggest successes in natural-language-related machine learning have been statistical speech recognition and statistical machine translation. The reason for these successes is not that these tasks are easier than other tasks; they are in fact much harder. Document classification that

Optimizing Discovery in the Big Data Era

The Netherlands eScience Center (NLeSC) supports and reinforces multidisciplinary and data-intensive research through creative and innovative use of ICT in all its manifestations. To stimulate this enhanced Science (eScience) NLeSC works as a network organization focused on collaboration, with the aim to change scientific practice by making large-scale data analysis possible across multiple disciplines.



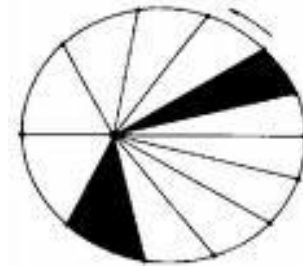
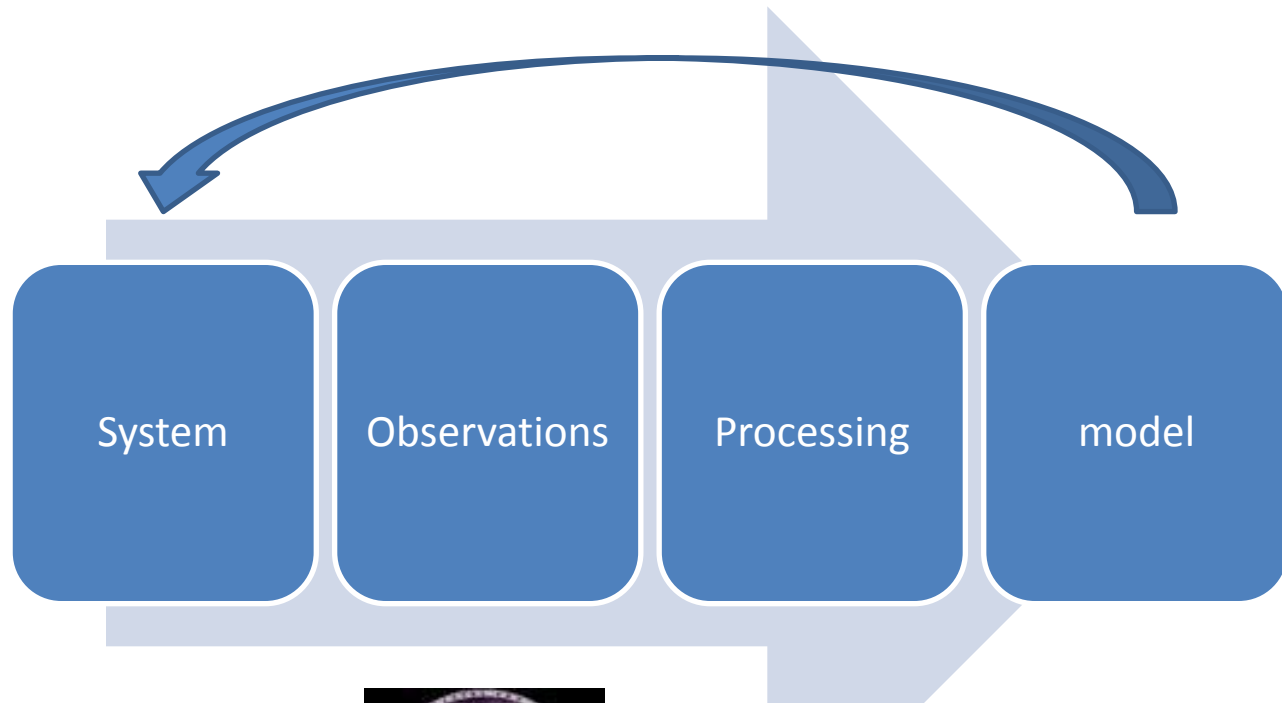
News

- [ASCI course: A Programmer's Guide for Modern High-Performance Computing](#)
 - [Dutch Research Consortium at SC12](#)
 - [eScience climatology project to reap benefits of access to US supercomputers](#)
 - [Carl-Christian Buhr of Neelie Kroes' cabinet visits NLeSC](#)
 - [EU blijft streven naar universeel ID](#)
- [▶ Read more](#)

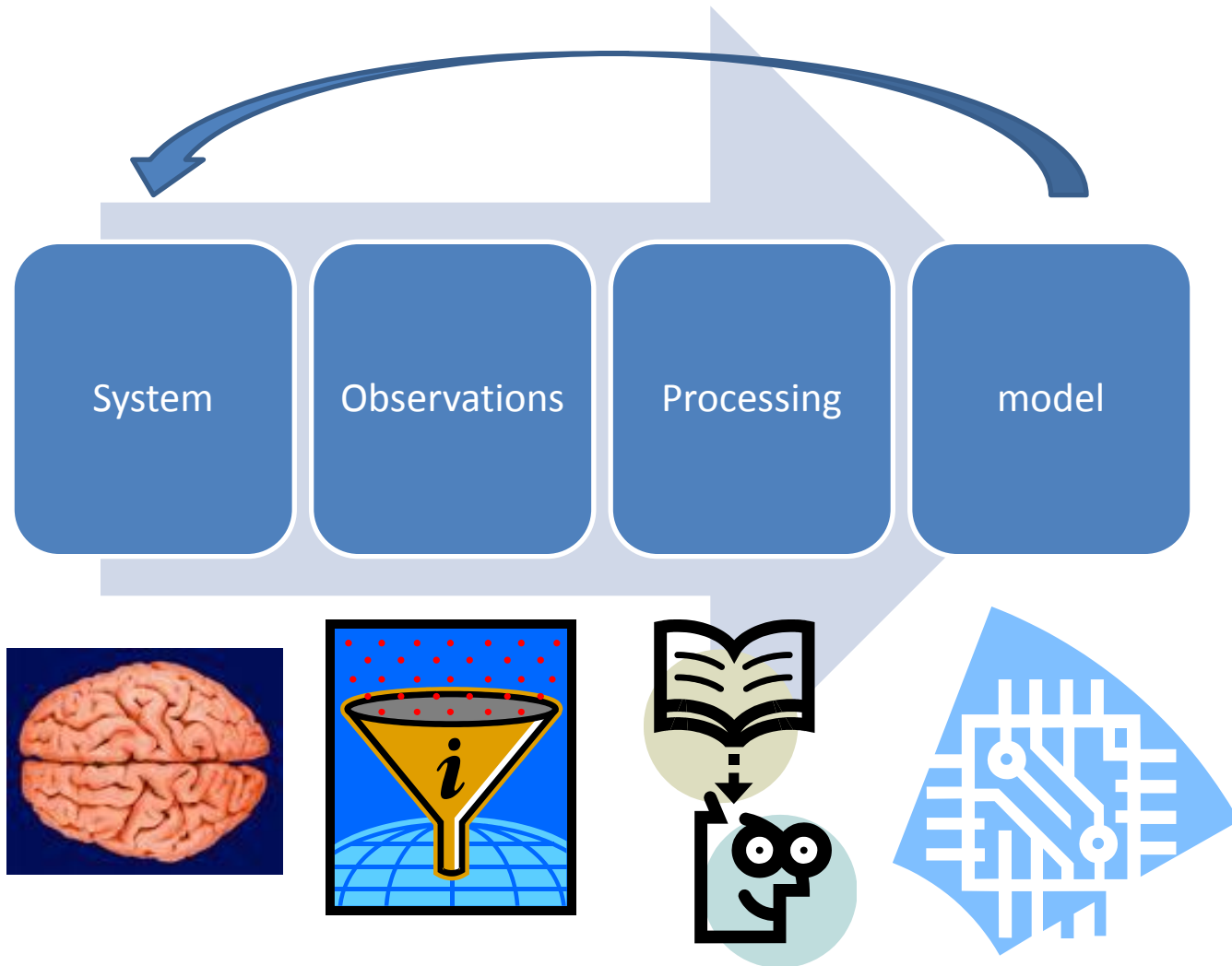
Agenda

- [Lustrum opleiding Bio-informatica 30 november 2012](#)
 - [Data sharing workshop 14 December 2012](#)
 - [e-BioGrid momentum 17 December 2012](#)
 - [27th SARA Superdag 19 December 2012](#)
 - [Free webinar Text Mining: a new way to discover food knowledge 17 January 2013](#)
- [▶ Read more](#)

Research cycle



Research cycle



Science in the 21st century

- Size and complexity of problems
 - Climate studies,
 - Human cognition,
 - The Cell
 - Elementary structure of matter,
 - Language
 - Social networks

Research objective facticity project

- Formulate a mathematical theory about an objective measurement of the amount of model information in a data set





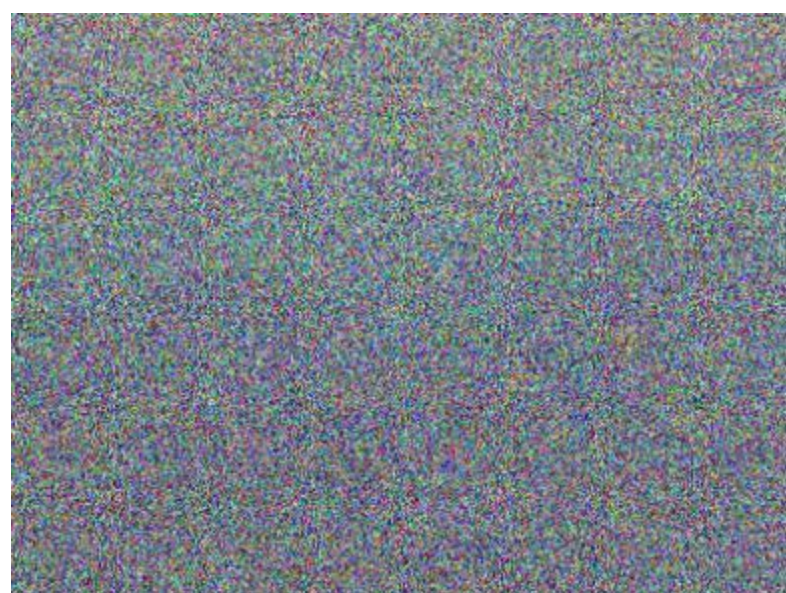
25 % noise



50 % noise

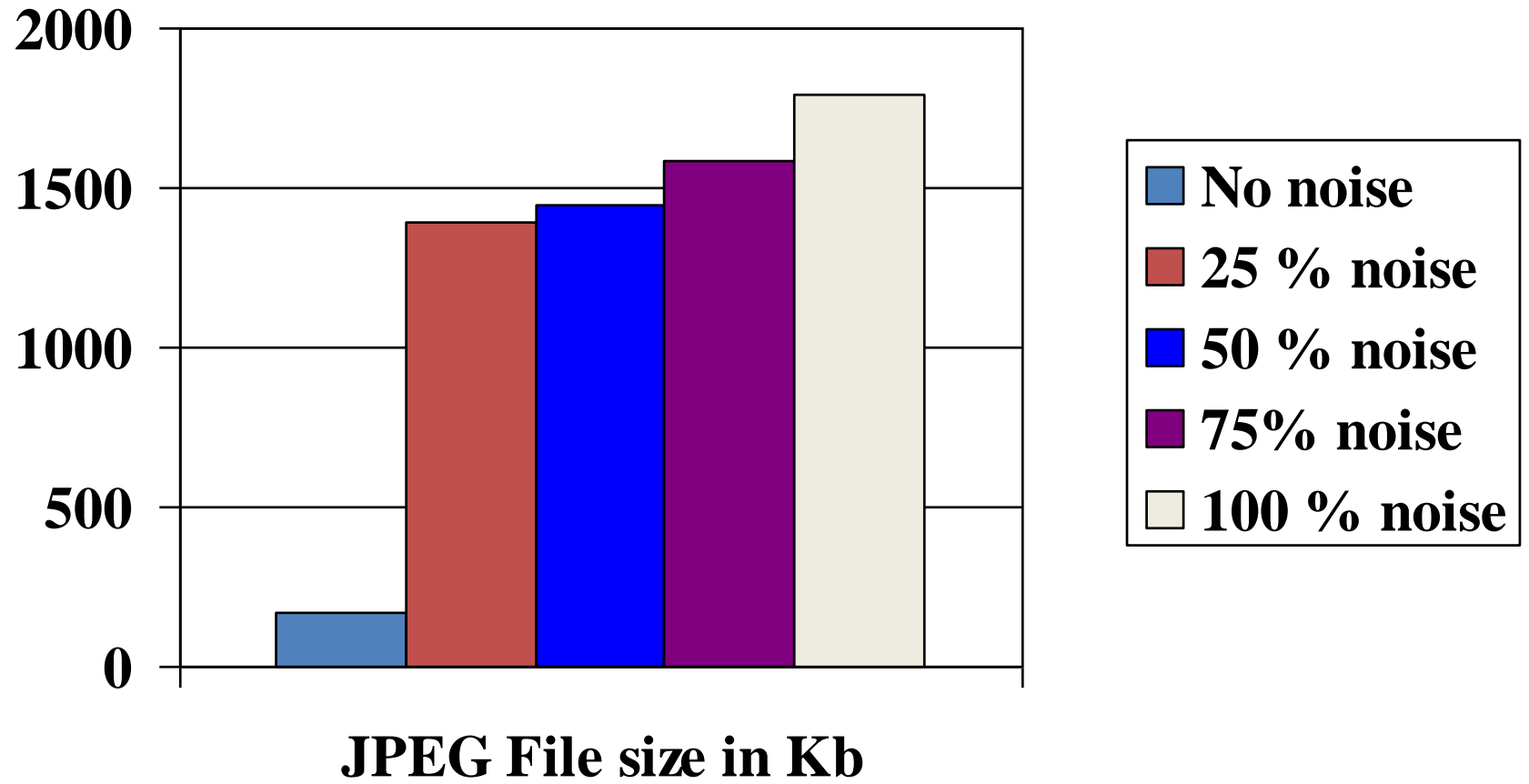


75 % noise



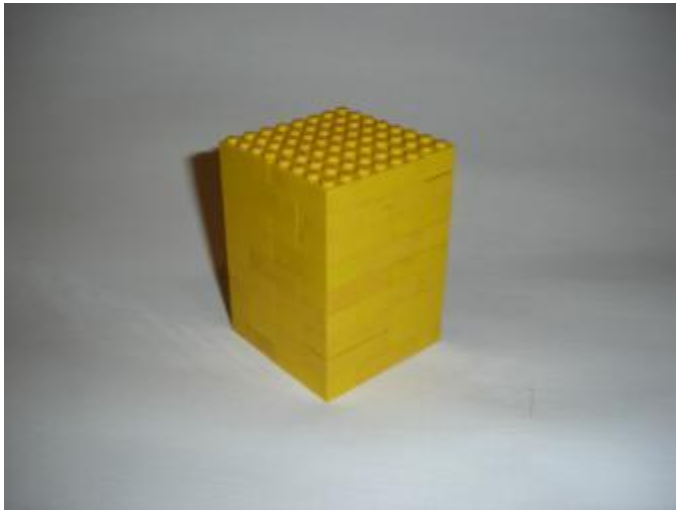
100 % noise

JPEG File size with noise added



Between order and chaos: facticity

Order



Facticity



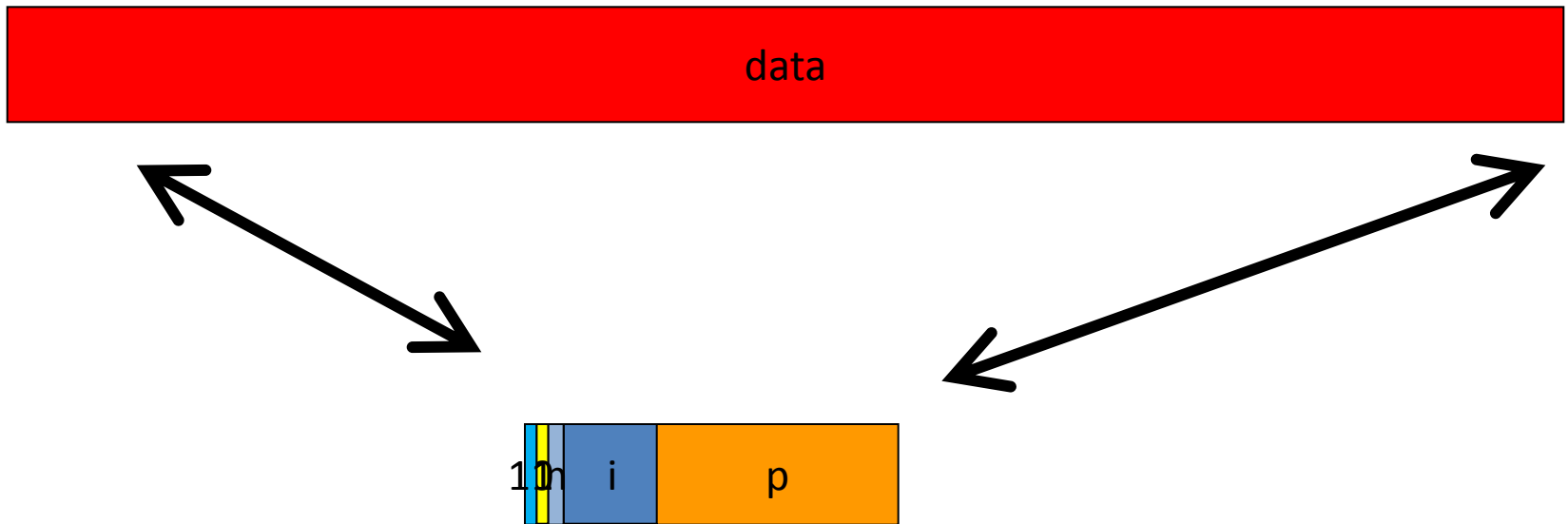
Chaos



Facticity (Questions)

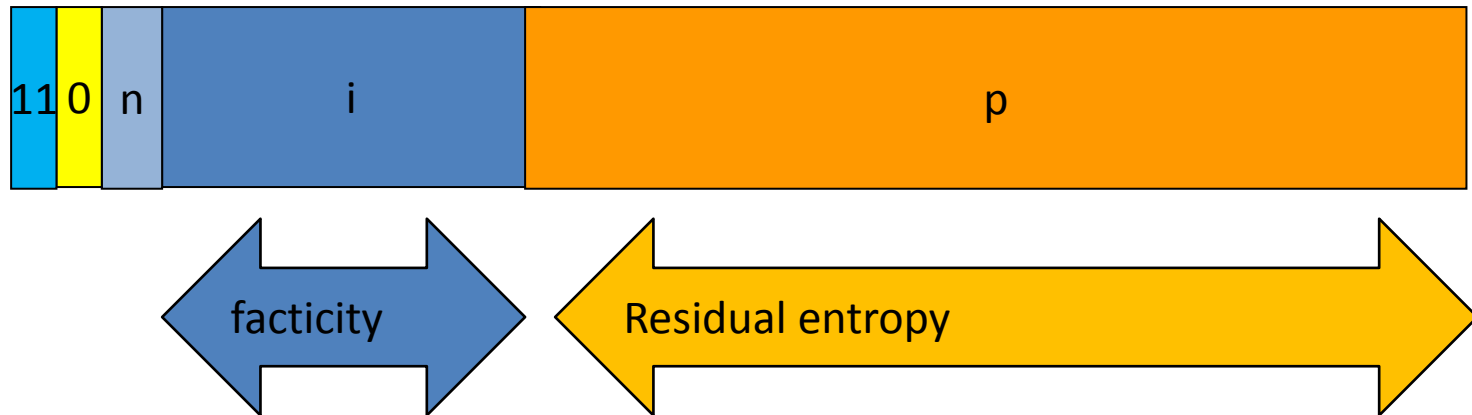
- Is there a universal model extraction method? (yes)
- Are there possibly more optimal models? (yes)
- Can we be sure that it is “definite” i.e. all relevant model information is extracted and nothing more (tricky, but yes)
- Is facticity stable? (No, but more a feature than a bug)

Turing two-part code compression



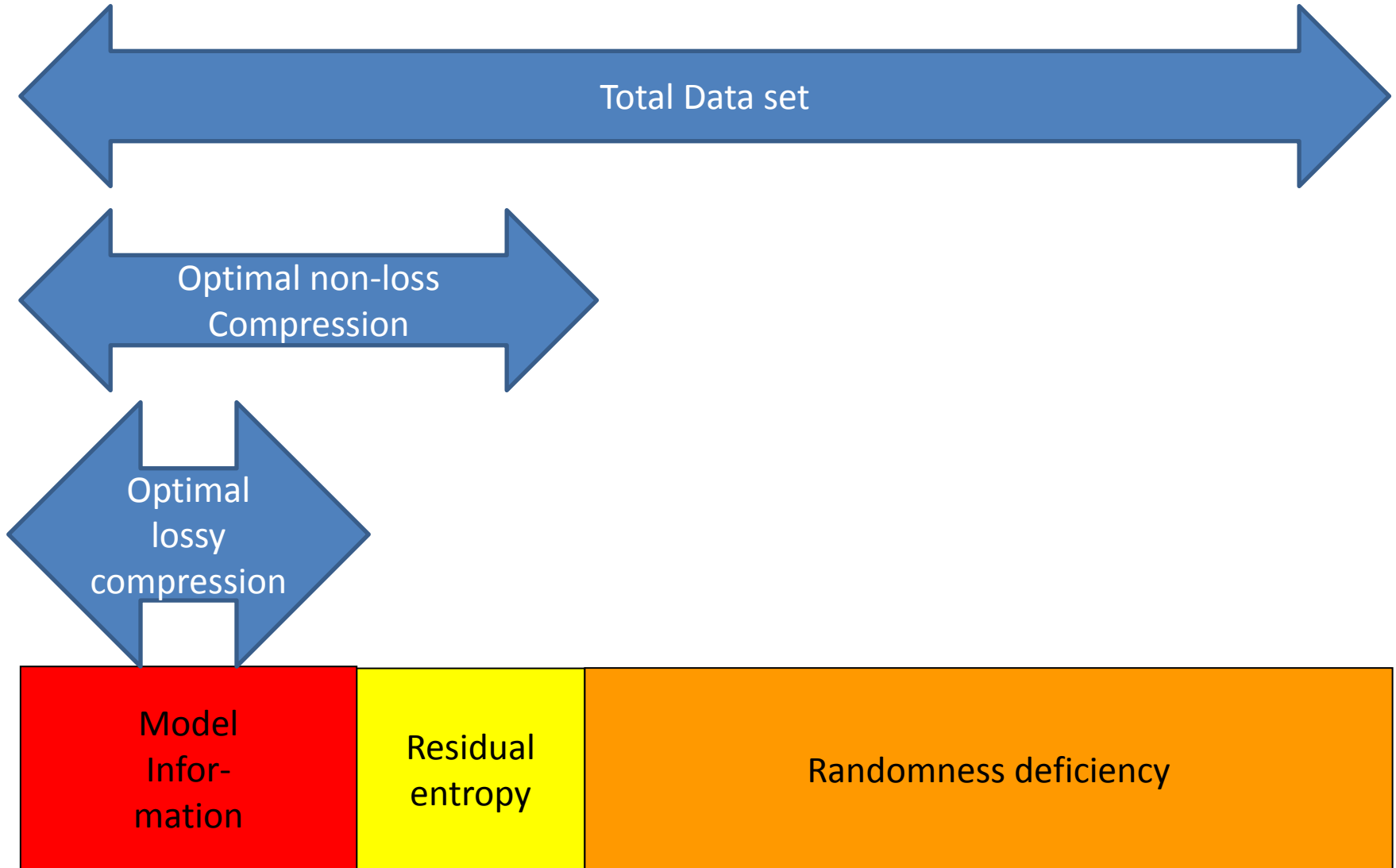
Definition of facticity: the amount of self-descriptive information in a dataset

$$\phi^U(x) = \arg \min_i \{ |i| : \{ \min_{i,p} \{ |\bar{i}| + |p| : U(\bar{i}p) = x \} \} \}$$



- The facticity $\phi(x)$ of a string x is the amount of self descriptive information in x .
- Facticity is definite!




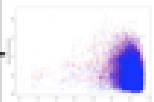









Crash Course Complexity Theory



existing complexity measures

- entropy estimators: (Shannon, Kolmogorov, VC-dimension)
- model complexity estimators (sophistication, computational depth, self-dissimilarity, facticity),
- network-analysis measures (structural graph-theoretic properties, centrality, betweenness, degree distributions)
- complex systems analysis (multi-scale analyses, robustness, dynamics)
- and their theoretical connections.

Atlas of Complexity: Comparative Analysis of Metrics & Datasets

Datasets Metrics		1-D CA	2-D CA	Bit Strings	RNA sequence data	Protein sequence data	Human text	Microarray & GWAS	Networks: Concept Web		RANDOM-ORG
Shannon Entropy											
Mutual Information											
Crutchfield Complexity											
Local Entropy Variance (per Herman)											
Critical Slowing Down (per Scheffer)											
PCA (per Lude)											
Facticity (per Pieter)											
Zip compression											
MIC (per Yahya)											



CVZ: dure medicijnen voor drie zeldzame ziektes uit basispakket

CVZ

'De verhouding tussen kosten en baten is onacceptabel'

'Het doorslaggevende argument is de buitengewoon ongunstige kosteneffectiviteit'

80
 € 40.000,-
 € 100.000,-

ziekte van Pompe

ziekte van Fabry

dure medicijnen voor ziekte van Pompe en ziekte van Fabry uit basispakket
 naar het schrappen van meer medicijnen voor zeldzame aandoeningen.

Het College voor Zorgverzekeringen
 drie medicijnen voor twee zeldzame
 en niet langer vergoeden uit het
 et. Dat staat in twee conceptadviezen, waar de NOG

door Niels

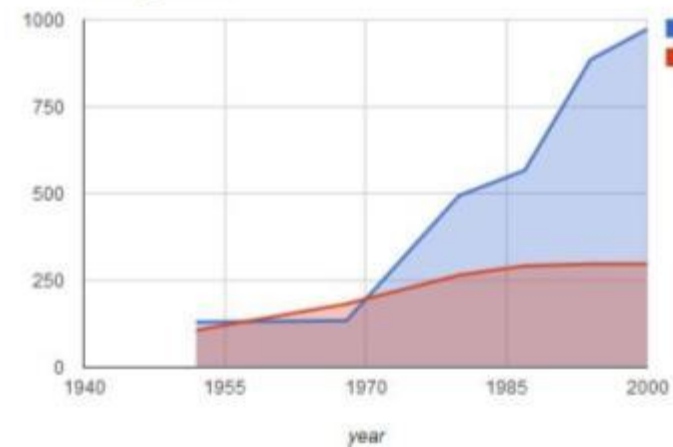
February 26, 2010
 30 Million Americans, 7,000 Rare Diseases

Posted: 07:01 PM ET
 By Peter Saltonstall
 Pres. & CEO of NORD - National Organization of
 Rare Disorders

While the eyes of the nation are directed toward
 Washington and the ongoing Health Care
 debate, the opportunity to raise
 awareness of rare diseases
 often flies below the radar of the
 national consciousness - rare diseases.

states, "rare" refers to conditions
 that affect fewer than 200,000 Americans.
 The National Institutes of Health (NIH)
 estimates that there are nearly 7,000 such diseases. These
 diseases collectively affect nearly 30 million
 Americans, or one out of every 10 people.
 Many of these uncommon—and in
 some cases, unknown—conditions. That's one out of every
 10 people in the nation, patients are plagued by unexplainable
 conditions, forcing them to endure difficult
 diagnosis.

DSM growth



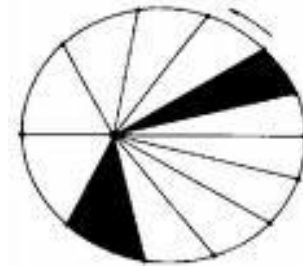
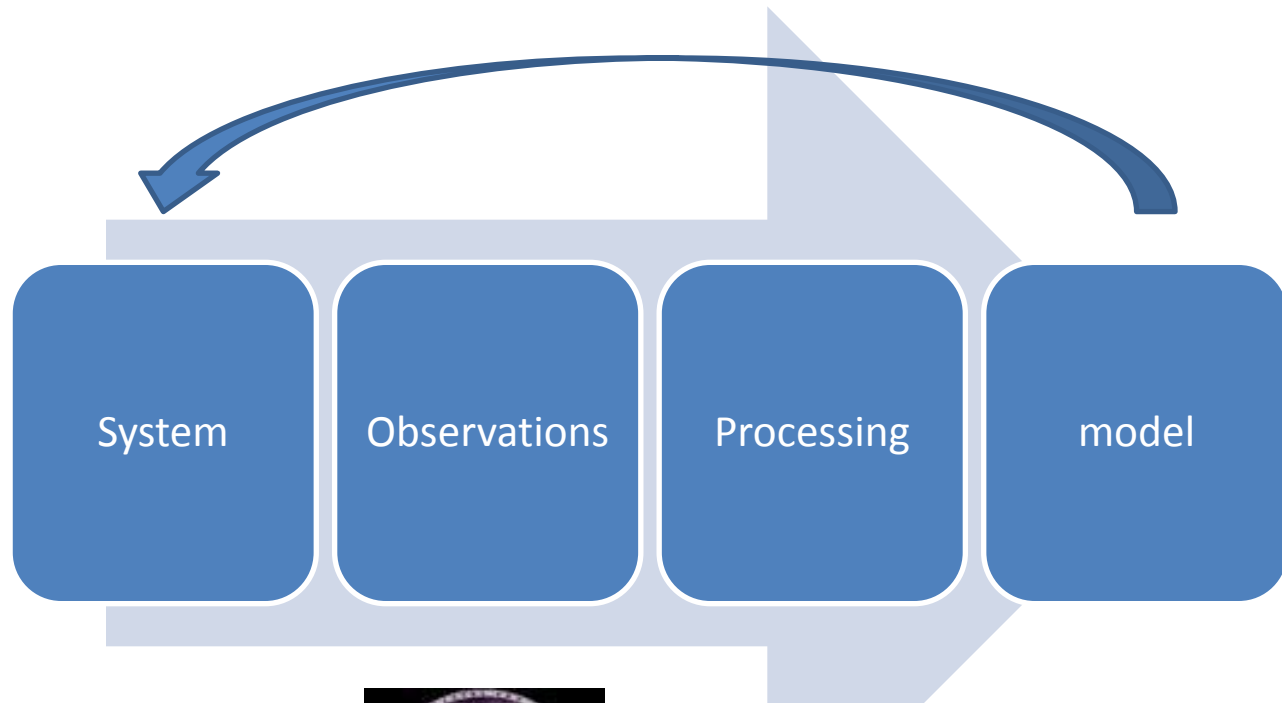
Research Questions

- What is a medically relevant model of an individual human being?
- How complex is such a model?
- How individual is such a model?
- How difficult is it to learn or make such a model?

medically relevant model of an individual human being

- A computational model (or program) that allows us to predict, explain (and if possible, help cure) any disease an individual human being X could have or get.

Research cycle



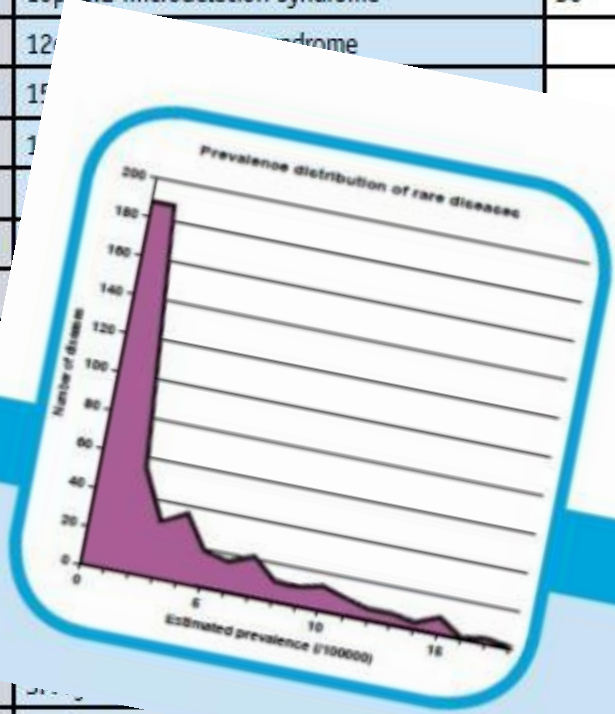
Research Methodology Proposal I

- Make a list of all possible diseases.
- For every disease study a relevant group of patients empirically and make a model that allows us to make deterministic predictions.

- How much diseases are there?
 - 2000?
 - 7000?

ORPHA number	Disease or group of diseases	Estimated prevalence (/100,000)	Number of published cases or families
61197	16p11.2 microdeletion syndrome	30	
4063	12p19.1 deletion syndrome		4 cases
99318	15q22.1 deletion syndrome		150 cases
4065	17p11.2 deletion syndrome		
76	17p13.1 deletion syndrome		
68261	17p13.1 deletion syndrome		
9157	17p13.1 deletion syndrome		
63693	17p13.1 deletion syndrome		
617	17p13.1 deletion syndrome		
510	17p13.1 deletion syndrome		
007	17p13.1 deletion syndrome		
5	17p13.1 deletion syndrome		
616	17p13.1 deletion syndrome		
	3-methylcrotonylglyoxaluria		
7046	3-methylglutaconic aciduria type 1		
7047	3-methylglutaconic aciduria type 3	10	
975	46,XX disorder of sex development - skeletal anomalies		2 cases
43	46,XX gonadal dysgenesis	< 10	
68558	46,XY disorder of sex development - adrenal insufficiency		9 cases

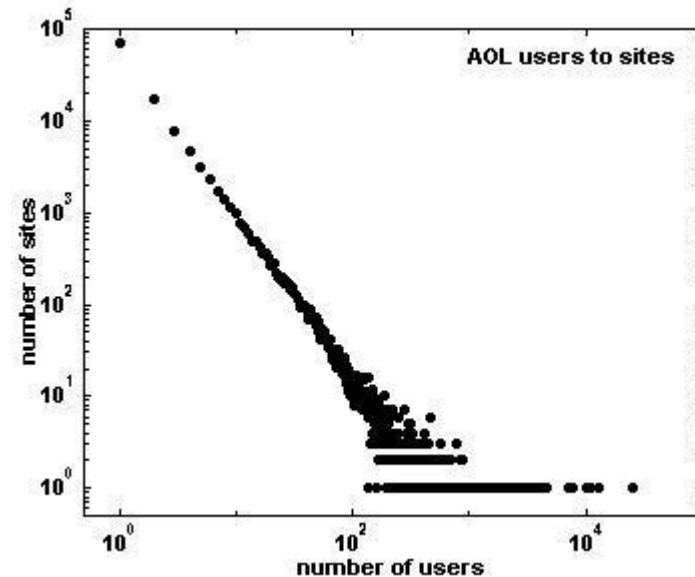
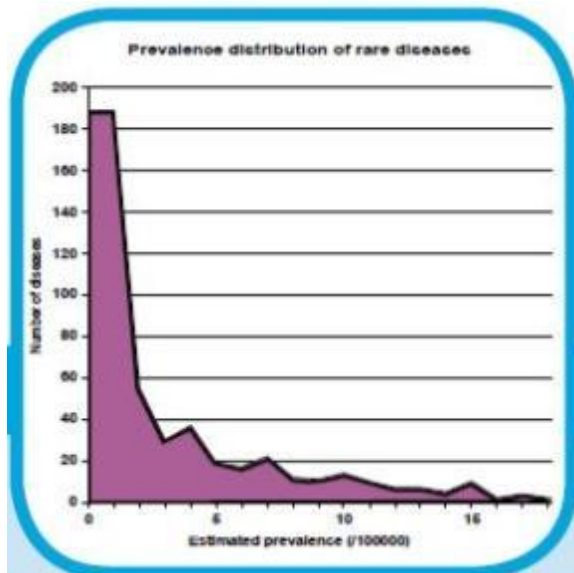
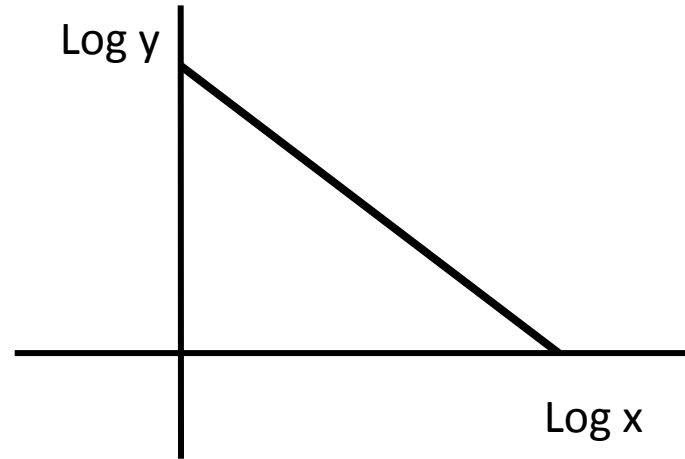
ORPHA Number	Disease or group of diseases	Estimated prevalence (/100,000)	Number of published cases or families
945	Acalvaria	< 1**	
67043	Acanthamoeba keratitis	1	
90301	Acanthosis nigricans - Insulin resistance - muscle cramps - acral enlargement		5 cases
	Acatlasemia	3.1	
		0.1	
			5 cases
			10
			cases
			cases
			100 cases
			100 cases
			100 cases
73214			
2221	Acquired hypertrophic cardiomyopathy		60 cases
99147	Acquired Von Willebrand syndrome		300 cases
36	Acrocollosal syndrome		34 cases



Orphanet Report Series
 Rare Diseases collection
 May 2012 | Number 1

Powerlaws

- $\log_c y = -a \log_c x + b$
- $y = c^b x^{-a}$



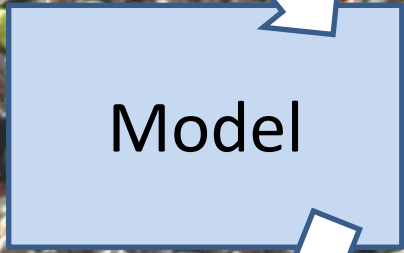
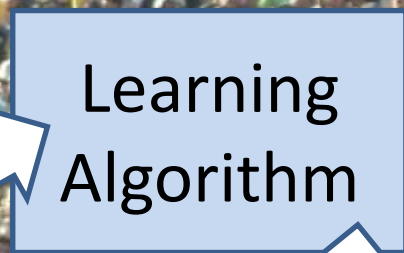
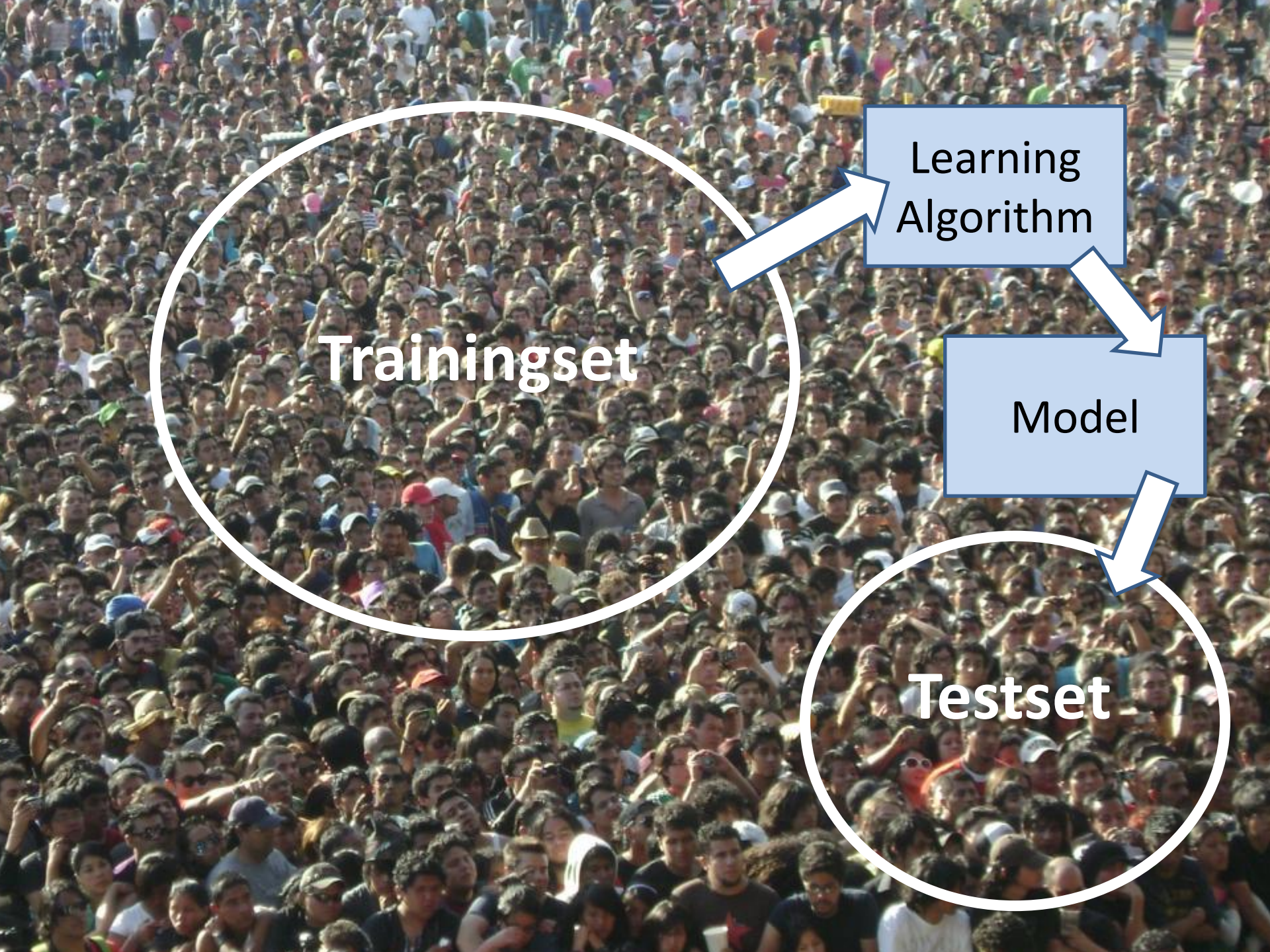
Problems!

- Power laws have no mean (i.e. it is infinite)
- If you double the population, you not only double the number of patients, but also the number of diseases.
- There is a (seemingly) unlimited supply of rare diseases.
- Most possible diseases are so rare that we never see them!

Research Methodology Proposal II

- Forget the list of all possible diseases.
- Forget deterministic prediction.
- Just study a relevant group of patients that is large enough to make statistically relevant risk assessments.

- How big a group of patients?



VC dimension

The bound on the test error of a classification model:

$$\text{training error} + \sqrt{\frac{h(\log(2N/h) + 1) - \log(\eta/4)}{N}}$$

Probability $1 - \eta$

Size of training set N
VC dimension: h
Condition: $h \ll N$

Problems again!!

Size of training set N

- Maximal:
current world population
= 0.7×10^{10}

VC dimension: h

- Maximal size of the
medically relevant model
of an individual human
being

Condition: $h \ll N$

- 10^8 bits \approx 12,5 Mb

12,5Mb!!

Maximal size of the
medically relevant
statistical model of an
individual human being.

Information: Some numbers

- 10^{10} bits Human genetic code
- 10^{14} bits Human brain
- 10^{32} bits Human body at quantum level
(10^{28} electrons)
- 10^{92} bits Total Universe

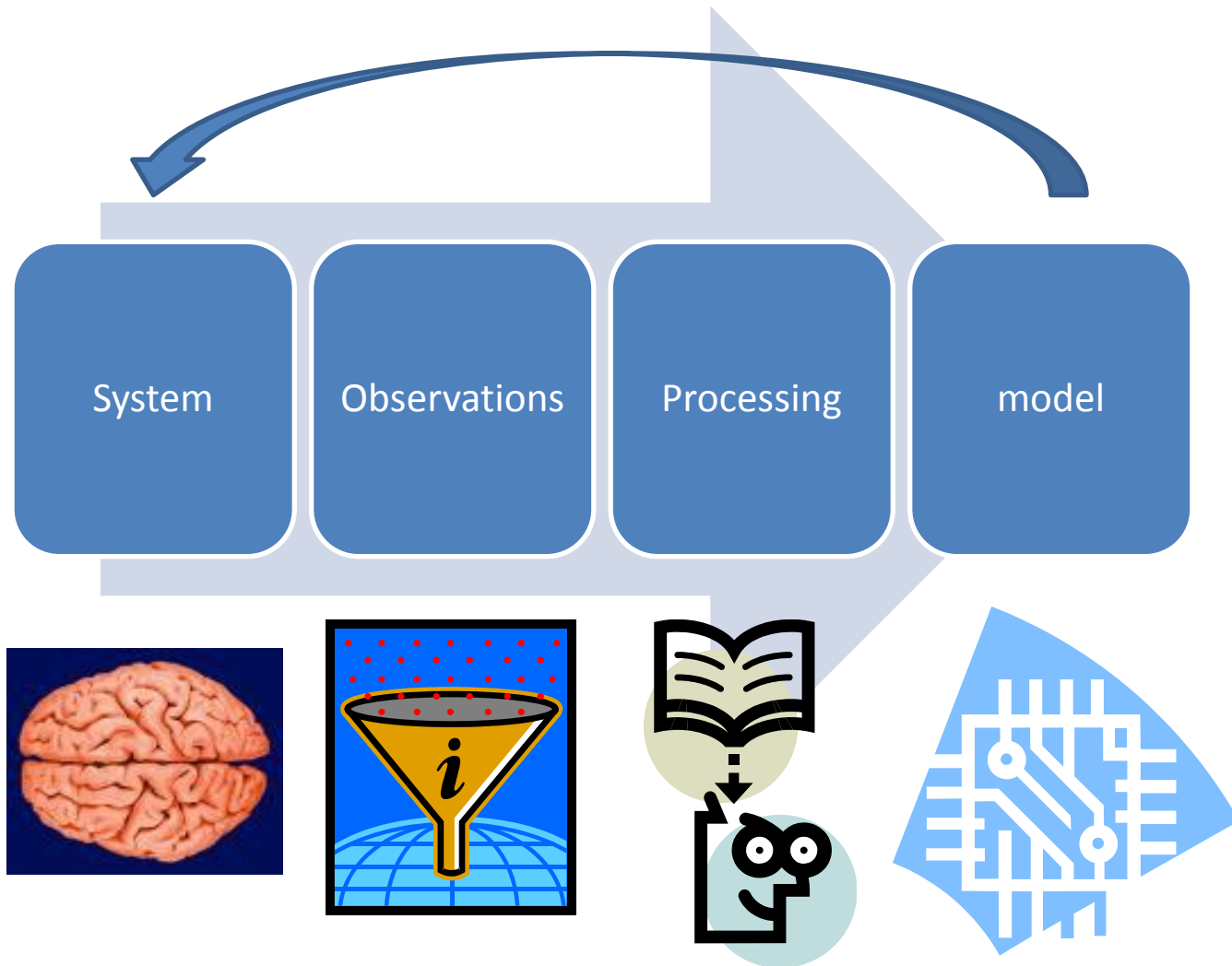
- 10^{123} Total number of
computational steps
since Big Bang (Seth Lloyd)



Empirical Incompleteness theory

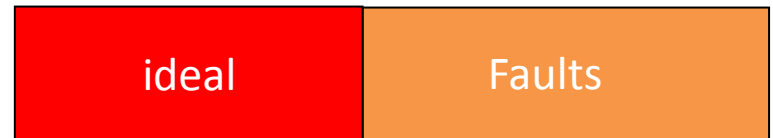
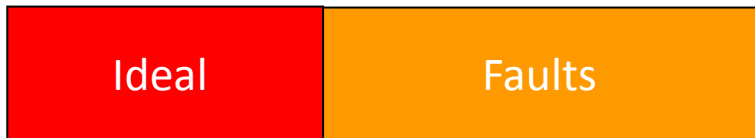
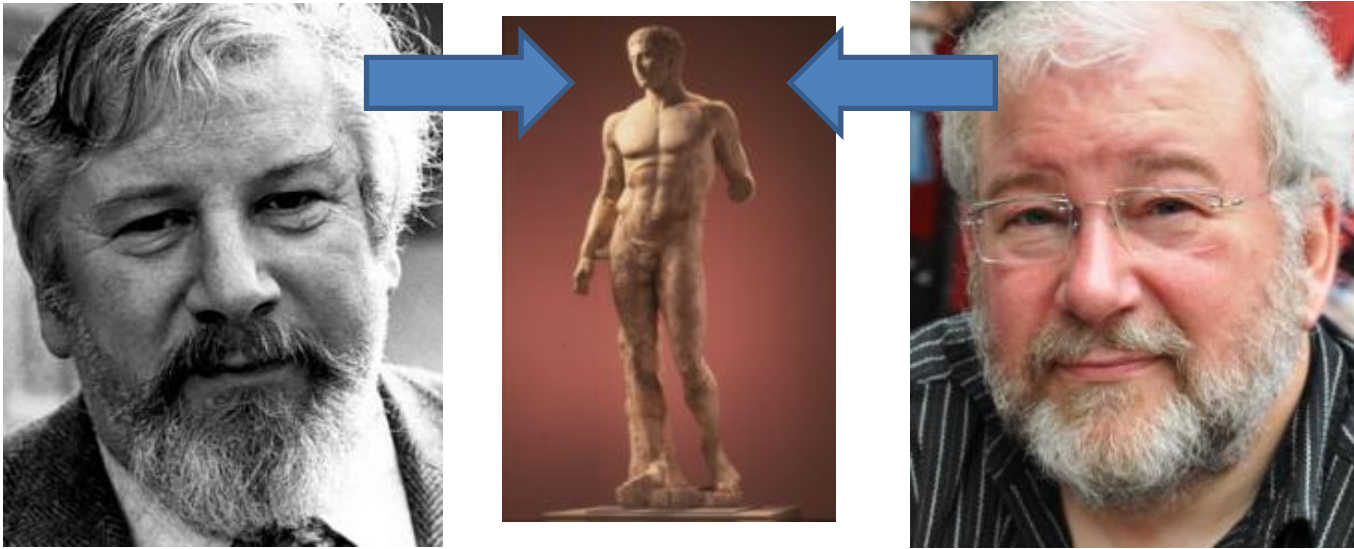
- Our universe contains classes of entities that have a population size that is much smaller than their relevant model complexity measured in bits.
- No adequate statistical models of such classes of entities based on frequency observations can be constructed.

Research cycle



Mutual model information

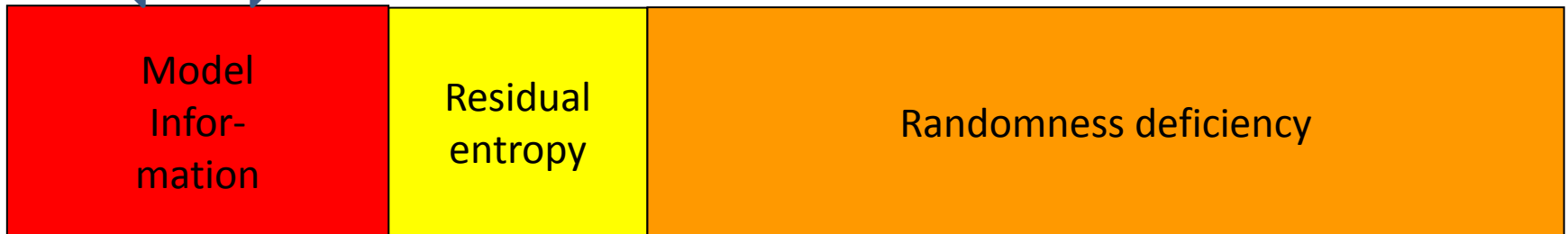
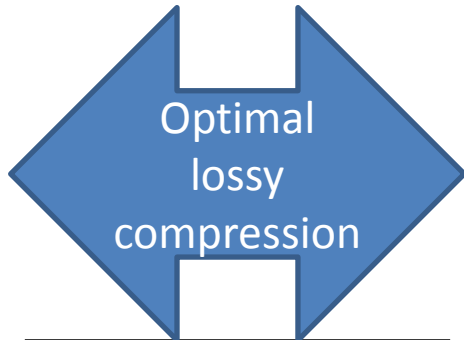
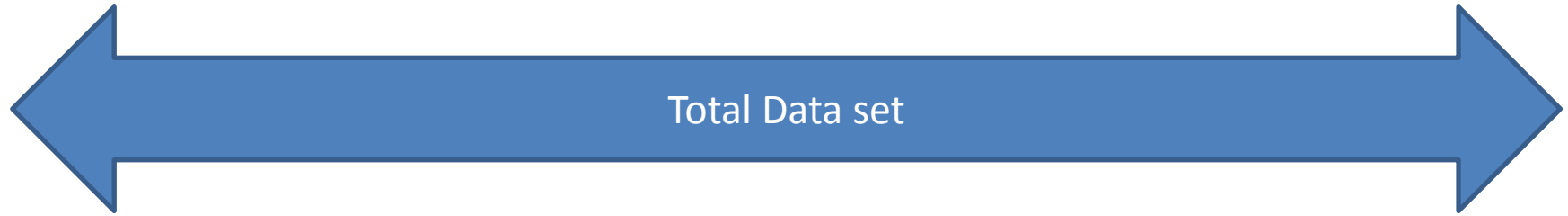
The classical view



Same model



Crash Course Complexity Theory



mutual model information? Computational View



Any two unrelated human beings differ by about 3 million distinct DNA variants.

The image shows a screenshot of a Facebook profile for Pieter Adriaans. The profile includes a cover photo of a man in a hat, a profile picture of the same man, and a bio that reads: "Studeerde Informatica aan University of Amsterdam", "Woonst in Kookengen", "Ligt Hengelo, Overijssel, Netherlands", and "Voeg je werkgever toe". The profile also shows a list of friends, a recent activity post, and a section for friends. A large, semi-transparent watermark with the text "+ Complete Personal History" is overlaid on the profile. The browser window shows the URL "www.facebook.com/pieter.adriaans" and the browser name "Firefox".

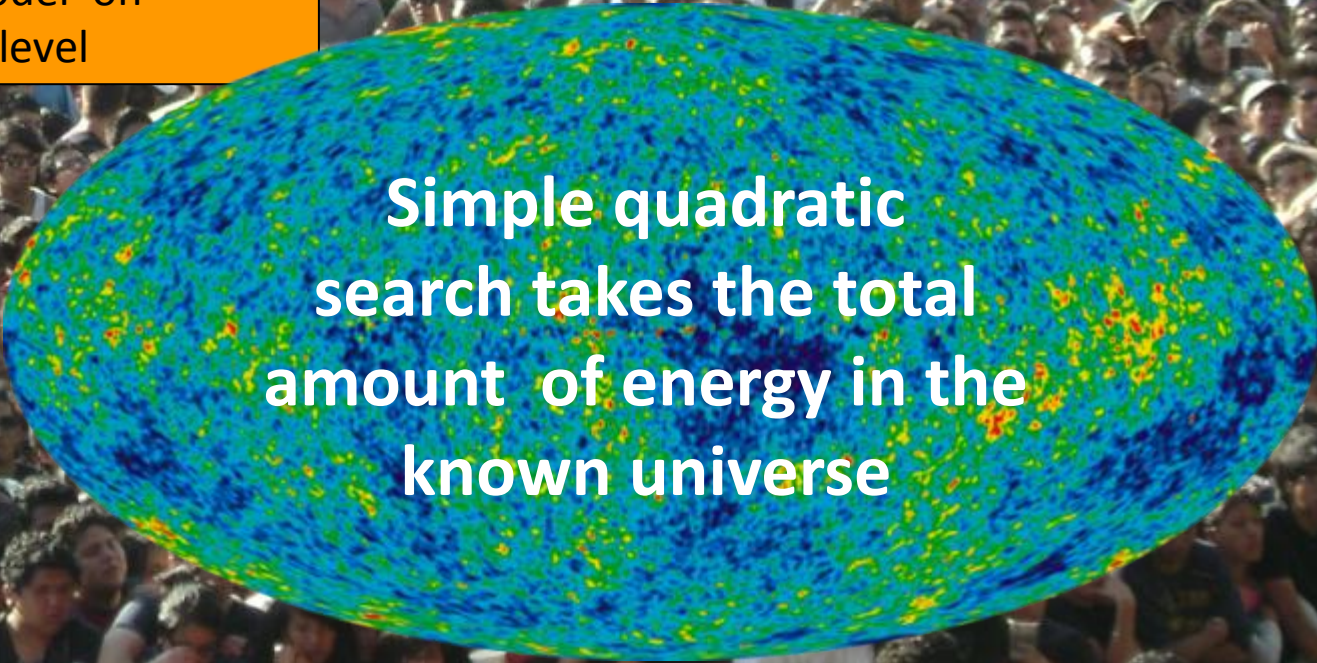
This block contains a collage of scientific and personal images. At the top is a red-tinted globe. Below it is a DNA microarray showing a grid of colored spots (yellow, green, red) on a black background. At the bottom is a close-up of a person's face, likely related to the DNA data shown.



Loading model
in to memory takes the
amount of energy
stored in a fully fuelled
Boeing 747-400

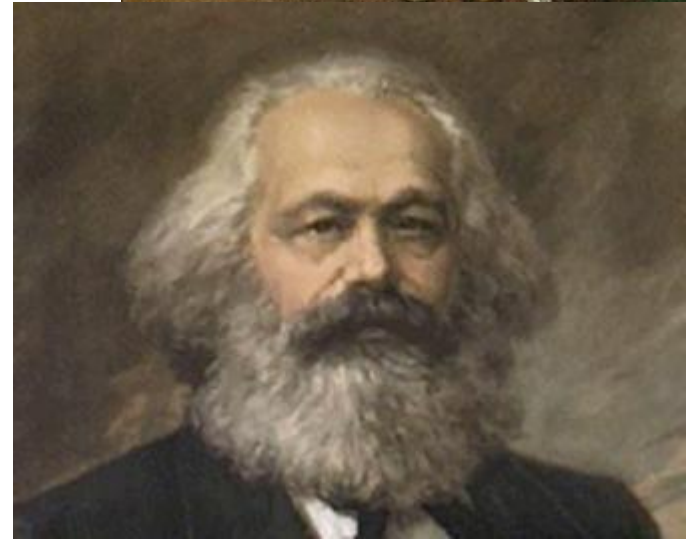
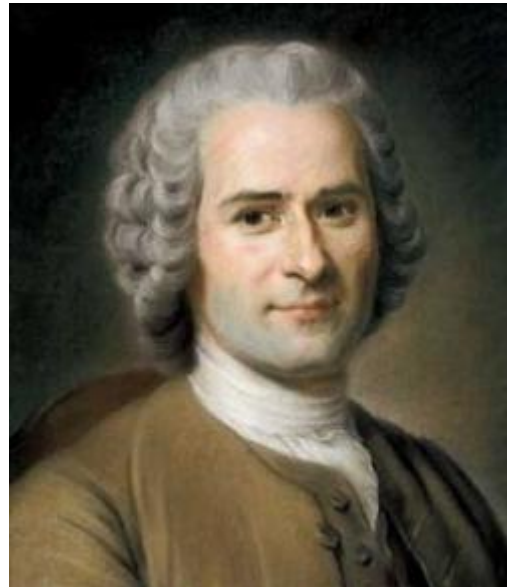
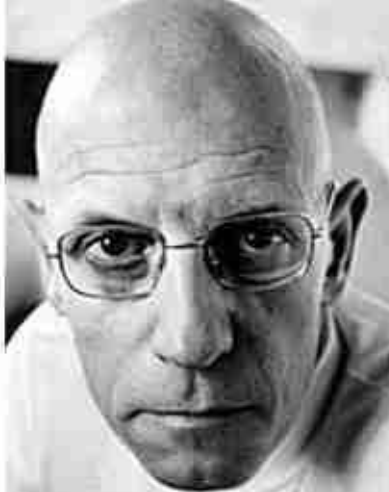
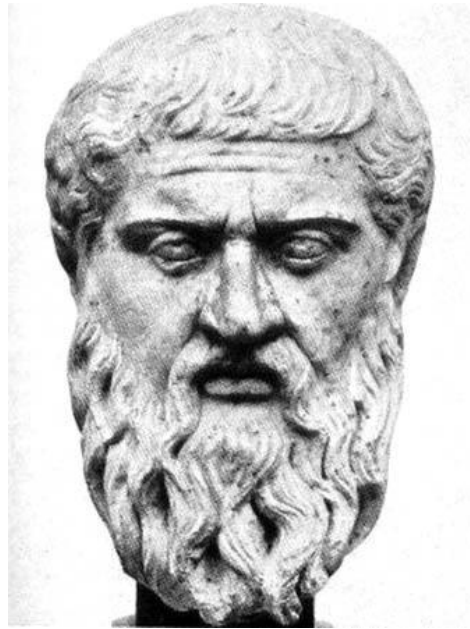


Individual model on
quantum level



Simple quadratic
search takes the total
amount of energy in the
known universe

Breaking a long tradition



Conclusions



- Data mining is developing in to the study of Big Data and e-Science
- We need to rethink scientific methodology for the 21st century.
- Our universe contains classes of entities that have a population size that is much smaller than their relevant model complexity measured in bits.
- The individual is the species!
- Consequences for:
 - Political organization of society
 - Ethics and human rights
 - Scientific Methodology
 - Medical policies
(Caritas versus evidence based medicine)