**Human-Centered AI**
**Research Seminar**

# Module 2: The Fundamentals: Theory of Science

**Andreas Holzinger**
**Human-Centered AI (Holzinger Group)**
**Institute for Interactive Systems and Data Science, TU Graz**
**Institute for Medical Informatics, Statistics & Documentation, Medical University Graz**
**and**
**Explainable AI-Lab, Alberta Machine Intelligence Institute, Edmonton, Canada**

**HCAI**
**HUMAN-CENTERED.AI**

**@aholzin #KandinskyPatterns**

Course Homepage: **https://human-centered.ai/hcai-research-seminar-2020**

---

## Remark

# This is the version for printing and reading. The lecture version is didactically different.

---

## Agenda

- In Module 2 we discuss questions including "What is science?", "*Why* contributing to the international scientific community?" and some methodological issues, whereas

- in Module 3 we discuss questions including "*How to* contribute to the international scientific community?" and learn the basic mechanics of science, the "know-how",

- Of course always from our human-centered AI and machine learning perspective …

---

# 01 Why should we contribute to the international scientific community?

- Science can provide explanations *)
- Science can make predictions
- *) these are the questions of "why" something is the cause (cf. Judea Pearl – see course 706.315)
- Why is a person doing this – anthropological explanations formed ancient philosophy – understanding and explaining nature …
- An explanation is a type of insight (sensemaking)
- When's a good explanation good? Obviously, when we **"feel"**\*) that something is satisfactorily explained – when we do not have any more questions – we understand it!

*) please watch the video shown on slide 7 by Richard Feynman

---

- 1) Causal explanation: If something causes X, it also explains X
- 2) Functional explanation: X has a reliable function f(X), thus X is explained
- 3) Purposeful explanation: X was wanted by Y
- 4) Pragmatic explanation: The explanation is adapted to the type of answer the questioner wants to hear!

---

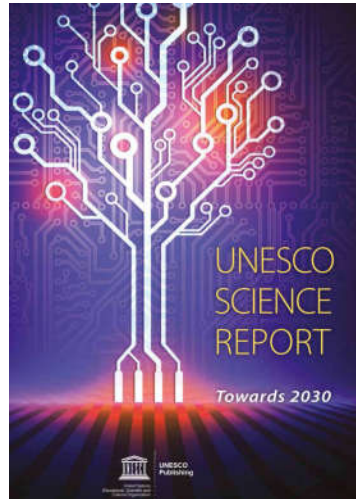Richard Feynman Magnets
708,434 views · 15 Apr 2009

This was one motivation to develop the Kandinsky Patterns, our "Swiss-Knife" for the study of explainable AI: https://human-centered.ai/project/kandinsky-patterns

https://www.youtube.com/watch?v=MO0r930Sn_8

---

https://en.unesco.org/themes/science-society

https://news.stanford.edu/2018/05/15/how-ai-is-changing-science

---

Friday 14 October lecture on 'Variational Inference.'

**Final Project**

In the second half of the course, you will complete a project. The ideal outcome of this project would be a paper that could be submitted to a top-tier machine learning conference such as NIPS, ICML, UAI, AISTATS, or KDD. There are different ways to approach this project, which are discussed in a more comprehensive document that is available from the course website under the Files tab. There are four separate components of the project:

---

- Progress is driven by the explosion in the availability of **big data** and **low-cost computation.**

- **Health is amongst the biggest challenges**

Jordan, M. I. & Mitchell, T. M. 2015. Machine learning: Trends, perspectives, and prospects. Science, 349, (6245), 255-260.

---

"An ultra-intelligent machine could design even better machines; there would then unquestionably be an **"intelligence explosion*"** and the intelligence of man would be left far behind …

It is curious that this point is made so seldom … outside of science fiction."

Irving John Good, Trinity College, Oxford, 1965
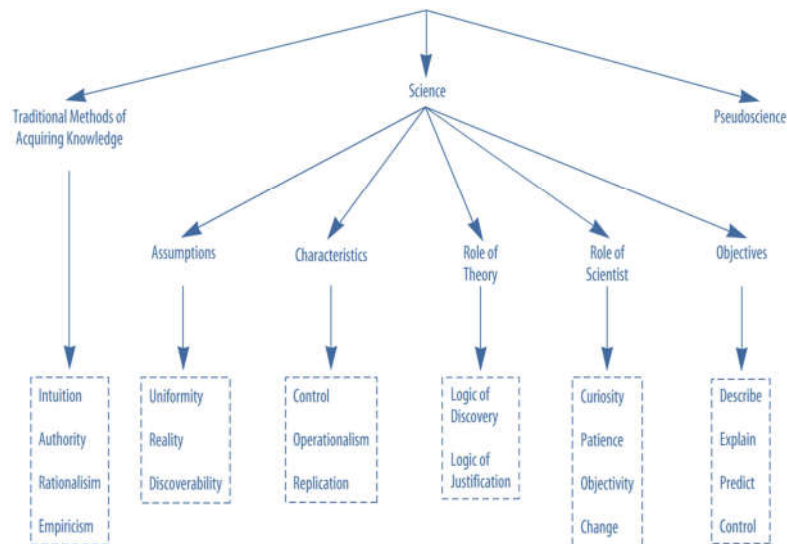Colleague of Alan Turing in Bletchley Park

Irving John Good 1966. Speculations Concerning the First Ultraintelligent Machine. Advances in Computers, 6, 31-88, doi: http://dx.doi.org/10.1016/S0065-2458(08)60418-0

Based on talks given in a Conference on the Conceptual Aspects of Biocommunications, Neuropsychiatric Institute, University of California, Los Angeles, October 1962; and in the Artificial Intelligence Sessions of the Winter General Meetings of the IEEE, January 1963 [1, 46]. The first draft of this monograph was completed in April 1963, and the present slightly amended version in May 1964
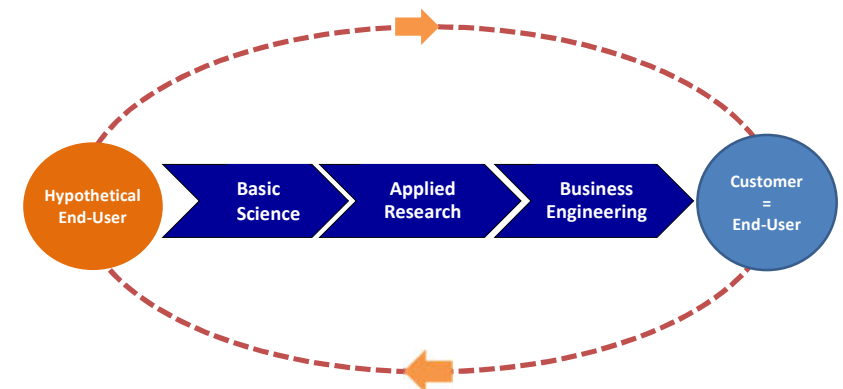
*) https://intelligence.org/ie-faq

## Slide 13

# 02 What is Science ?
# What is Engineering?
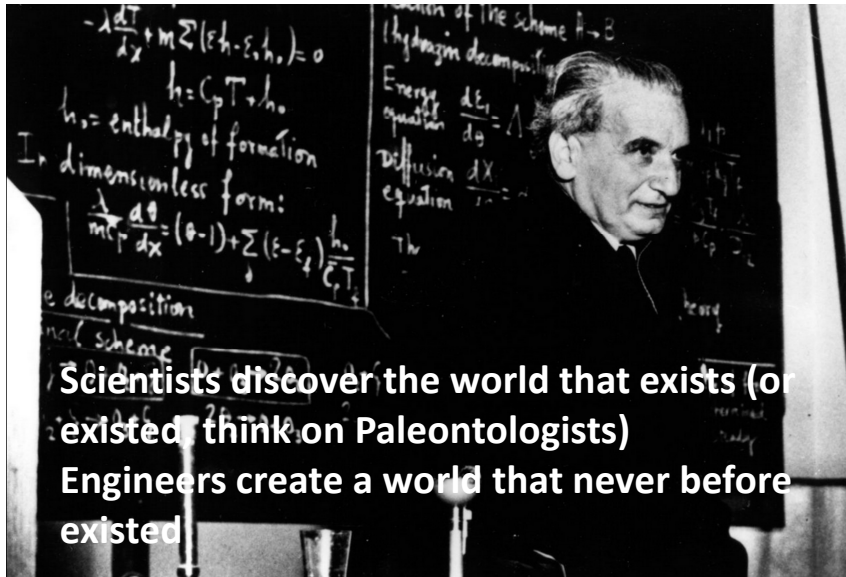
## Slide 14

- επιστημη (ancient Greek: episteme = Wissen)
- Scientia (latin = Wissen, engl. "knowledge")
- <u>systematic</u> and objective process to attain and organize new knowledge in the form of testable explanations and predictions about our universe
- Computer Science:
  - A) theoretical
  - B) experimental
  - Three pillars:
    - 1) language = information;
    - 2) process = algorithms;
    - 3) technology = computer (universal programable Machine)

## Slide 15

### Naïve Assumptions and scientific methods

## Slide 16

### Motto of the Holzinger Group

# Science is testing crazy ideas – Engineering is putting these ideas into Business



Holzinger, A. 2011. Successful Management of Research and Development, Norderstedt: BoD.

**Scientists discover the world that exists (or existed, think on Paleontologists)**
**Engineers create a world that never before existed**

Attributed to Theodore von Karman (1881-1963)

---

- Computer Science is asking:
  - P = NP ? *)
  - What is information?
  - What is intelligence?
  - What is computable?
- Computer Engineering is asking:
  - (How) can we build (intelligent) information systems (simply)?

  - *) a proof would have profound impact (think on cryptography, TSP, subgraph isomorphism, …

Janet M. Wing 2006. Computational thinking. Communications of the ACM, 49, (3), 33-35, doi:10.1145/1118178.1118215.

---

According to Aristoteles (384-322)

|  | Episteme | Techne |
|---|---|---|
| Objects | "unchangeable" | Changeable (plastic) |
| Goal | General knowledge | Specific knowledge |
| Activities | Building theoria | Building poiesis |
| Method | Abstraction | Concrete (Modeling) |
| Process | Conceptualizing | Optimizing |
| Innovation in form of | Discovery | Invention |
| Results | Law-like | Rule-like |

---

P versus NP and the Computational Complexity Zoo, please have a look at
https://www.youtube.com/watch?v=YX40hbAHx3s

$$\left( -\frac{\hbar^2}{2m} \Delta + U(\vec{r}, t) \right) \psi(\vec{r}, t) = i\hbar \frac{\partial}{\partial t} \psi(\vec{r}, t)$$

---

to reproduce …

to grow …

to evolve …

to self-replicate …

to generate/utilize energy …

# to process information …

Schrödinger, E. (1944) *What Is Life? The Physical Aspect of the Living Cell. Dublin Institute for Advanced Studies.*

---

Tony Hey, Stewart Tansley & Kristin Tolle 2009. The fourth paradigm: data-intensive scientific discovery, Redmond (WA), Microsoft Research.

https://www.microsoft.com/en-us/research/publication/fourth-paradigm-data-intensive-scientific-discovery

---

FIGURE 2.
While it took 2,300 years after the first report of angina for the condition to be commonly taught in medical curricula, modern discoveries are being disseminated at an increasingly rapid pace. Focusing on the last 150 years, the trend still appears to be linear, approaching the axis around 2025.

---

- **1000 years ago – Experimental Science:** observing and describing natural phenomena

- **300 years ago – Theoretical Science:** Newton's Laws, Maxwell's Equation, Einstein, …

- **70 years ago – Computational Science:** using universal programmable machines for simulating complex phenomena via math models

- **Today – Data-Intensive Science:** data captured by instruments, generated by simulations, sensor nets, …

---

Sauro Succi & Peter V. Coveney 2019. Big data: the end of the scientific method? Philosophical Transactions of the Royal Society A, 377, (2142), 20180145, doi:10.1098/rsta.2018.0145.

Available online: https://royalsocietypublishing.org/doi/full/10.1098/rsta.2018.0145

---

# 03 Basics of the Theory of Science (Wissenschaftstheorie)

- **Classic Newtonian** approach:
  - Ask question > develop theory > form a hypothesis to proof/disproof theory > conduct experiments > compare data with hypothesis > accept/reject theory
- **Computer Science** approach:
  - Find open problems to solve > form hypothesis how to solve the problem > experiment > evaluate > present new solution to the problem
- **Modern Machine Learning** approach:
  - Setting up experiments to answer questions including: How does model $m$ perform on data $d$ from domain $D$? Which of these models have the best performance? Much is feature engineering and precision and recall are your best friends! Now questions of causality to answer the questions of why are becoming important!

---

|  | Objects | Primary Method |
|---|---|---|
| Logic & Mathematics | Simple abstract objects: numbers, propositions, … | Deduction |
| Natural Sciences | Natural objects: physical objects, fields, organisms, … | Hypothetico-deductive |
| Social Sciences | Social objects: human individuals, collectives, society, … | Hypothetico-deductive & Hermeneutics |
| Humanities (e.g. Law) | Complex cultural objects: human ideas, principles, actions, relationships, language, artefacts, … | Hermeneutics |
|  |  |  |

---

- **Deductive Reasoning** = Hypothesis > Observations > Logical Conclusions
  - DANGER: Hypothesis must be correct! DR defines whether the truth of a conclusion can be determined for that rule, based on the truth of premises: A=B, B=C, therefore A=C
- **Inductive reasoning** = makes broad generalizations from specific observations
  - DANGER: allows a conclusion to be false if the premises are true
  - generate hypotheses and use DR for answering specific questions
- **Abductive reasoning** = inference = to get the best explanation from an incomplete set of preconditions.
  - Given a true conclusion and a rule, it attempts to select some possible premises that, if true also, may support the conclusion, though not uniquely.
  - Example: "When it rains, the grass gets wet. The grass is wet. Therefore, it might have rained." This kind of reasoning can be used to develop a hypothesis, which in turn can be tested by additional reasoning or data.

---

- := information provided by direct observation (empirical evidence) in contrast to information provided by inference
  - Empirical evidence = information acquired by observation or by experimentation in order to verify the truth (fit to reality) or falsify (non-fit to reality).
  - Empirical inference = drawing conclusions from empirical data (observations, measurements)
  - Causal inference = drawing a conclusion about a causal connection based on the conditions of the occurrence of an effect.
    - Causal inference is an example of causal reasoning.

https://human-centered.ai/project/kandinsky-patterns

A) True



B) False



C) Counterfactual

# 04 Hypothetico-Deductive Method

Karl Popper's Falsification

https://www.youtube.com/watch?v=wf-sGqBsWv4

- The scientific method is the logical scheme used by scientists searching for answers to the questions posed within science.
- The Scientific method is used to produce scientific theories, including both scientific meta-theories (theories about theories) as well as the theories used to design the tools for producing theories (instruments, algorithms, etc).

- Observe an event
- Ask a question ("Originäre Fragestellung") and check with the state-of-the-art whether and to what extent this question has already been answered!!
- Develop a model (or hypothesis) which makes a prediction to explain this event
- Test your prediction with (new) data
- Observe the result
- Proof the hypothesis or revise appropriately
- repeat as long as needed
- A successful model (or hypothesis) becomes a scientific theory !

---

---

1. Formulate a research question in the context of existing knowledge (theory & observations) – you must know who did what when and where!
2. Formulate a hypothesis as a tentative answer to this question
3. Deduce consequences and make predictions
4. Test the hypothesis in a specific experiment/theory field. The new hypothesis must prove to fit-in the existing world-view, think about what Sir Karl Popper said!

---

- In case the hypothesis leads to contradictions and demands a radical change in the existing theoretical background, test it carefully again!
- If you succeed and it replaces the existing scientific paradigm – this is called "scientific revolution" and it happens very rarely and cannot be planned …
- Repeat the process with modifications of the hypothesis until you reach an agreement which leads to a satisfiable result
- If you find major discrepancies, you must start the process from the beginning, or you state an alternative research question!

- When consistency is obtained the hypothesis becomes a theory and provides a coherent set of propositions that define a new class of phenomena or a new theoretical concept.
- The results have to be published and is subject of process of "natural selection" among competing theories ... reviewer give you a hard time!
- A theory is then becoming a framework within which observations/theoretical facts are explained and predictions can be made.

- Science undergoes always periodic paradigmatic changes
- These paradigm shifts open up new approaches
- Scientists can (of course) never separate their subjective perspective from their work
- thus, our comprehension of science can never rely on full objectivity according to Kuhn

https://en.wikipedia.org/wiki/Commensurability_(philosophy_of_science)
For further information on the work of Thomas Kuhn et al.

- There is no logical path leading to [the highly universal laws of science]. They can only be reached by intuition, based upon something like an intellectual love of the objects of experience'.
- Science starts with problems (yes, we engineers call it challenges !).
- Coping with a problem, the scientist makes observations.
- Observations are selectively designed to test if a theory functions as a satisfactory solution to a given problem.
- Read more here: http://plato.stanford.edu/entries/popper

- Science is interested in universal affirmative conclusions.
- However, such conclusions could never be verified.
- But, they could be falsified by the discovery of a counterexample!
- Science should aim not to verify or confirm hypotheses but to falsify them.
- According to Popper, there can be never a confirmation (Bestätigung) of a hypothesis.
- It can only be a corroboration (Bekräftigung, Erhärtung)

https://www.youtube.com/watch?v=0KmimDq4cSU

# 05 Occam's Razor

- Occam's razor (novacula Occami; or law of parsimony: Latin: lex parsimoniae)
- is the problem-solving principle that states "Entities should not be multiplied without necessity."
- "Pluralitas non est ponenda sine necessitate"
- "All else equal – prefer the simplest theory"
- "The simplest solution is most likely the right one."
- Occam's razor says that when presented with competing hypotheses that make the same predictions, one should select the solution with the fewest assumptions



William of Ockham, or Occam.
This image is in the public domain

Problem is Complexity:
Ad-hoc hypotheses
Multiple mechanisms
Coincidences
Many free parameters, …

- Explanation is difficult:
  - A) Is the minimum explanation the simplest ?
  - B) Is the simplest explanation the best explanation ?
  - C) When is it enough? ("where is the saturation point")
  - D) How can an explanation be adapted to different previous knowledge?

Pedro Domingos 1999. The Role of Occam's Razor in Knowledge Discovery. Data Mining and Knowledge Discovery, 3, (4), 409-425, doi:10.1023/a:1009868929893.

Anselm Blumer, Andrzej Ehrenfeucht, David Haussler & Manfred K Warmuth 1987. Occam's razor. Information processing letters, 24, (6), 377-380.

- Note: it means to prefer the most obvious explanation, e.g. you cheese is vanished:
- A) most likely it has been eaten by a mouse
- B) most unlikely it has been taken by a Marsian
- Danger: "best explanation" in what sense?
- Occam originally emphasized that you should always take an explanation with the fewest assumptions
- Counterexample: Quantum Mechanics

"Nature operates in the shortest way possible" – Aristotle.

"Entities should not be multiplied without necessity" – William of Occam.

"Everything should be made as simple as possible, but not simpler" – Albert Einstein.

# 06 Reichenbach's Principle

## Hans Reichenbach (1891-1953)

**Proponent of logical empiricism.**

- We can use probability to decide if A is a possible cause of B.
- If $P(A\&B) > P(A)P(B)$ we say that A and B are (positively) correlated. Then A could be a cause of B.
- The condition is equivalent to $P(B|A) > P(B)$.
- The condition is symmetric in A and B: If A is a possible cause of B, then B is a possible cause of A.
- "Which came first? The hen or the egg."

"If an improbable coincidence has occurred, there must exist a common cause" (p. 157)

**Two examples:** "Suppose both lamps in a room go out suddenly. We regard it as improbable that by chance both bulbs burned out at the same time and look for a burned out fuse or some other interruption of the common power supply. The improbable coincidence is thus explained as the product of a common cause."

"Or suppose several actors in a stage play fall ill showing symptoms of food poisoning. We assume that the poisoned food stems from the same source – for instance, that it was contained in a common meal – and then look for an explanation of the coincidence in terms of a common cause."

THE DIRECTION OF TIME — HANS REICHENBACH

Classical probability measure space: $(\Omega, \Sigma, p)$
Positive correlation: $A, B \in \Sigma$

$$p(AB) > p(A)\,p(B)$$

Reichenbachian common cause: $C \in \Sigma$

$$p(AB|C) = p(A|C)p(B|C)$$
$$p(AB|C^\perp) = p(A|C^\perp)p(B|C^\perp)$$
$$p(A|C) > p(A|C^\perp)$$
$$p(B|C) > p(B|C^\perp)$$

---

- **Reichenbach's common cause principle:**
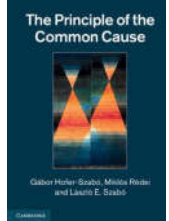  Assume that *X not independent Y* .

- Then
  *X* causes *Y* ,
  *Y* causes *X*,
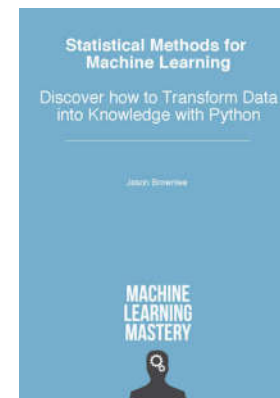  there is a hidden common cause or combination.
  For example:

- X = You like this lecture hall
  Y = You like this theory
  Z = You attend this course

The Principle of the Common Cause — Gábor Hofer-Szabó, Miklós Rédei and László E. Szabó · CAMBRIDGE

---

# 07 Experiments in Machine Learning

---

Statistical Methods for Machine Learning — Discover how to Transform Data into Knowledge with Python — Jason Brownlee — MACHINE LEARNING MASTERY

The COMP61011 Not-Very-Scary Guide to … Machine Learning — Professor Gavin Brown, University of Manchester (Sept 2019) — http://studentnet.cs.manchester.ac.uk/pgt/COMP61011/

https://machinelearningmastery.com　　　http://syllabus.cs.manchester.ac.uk/pgt/2019/COMP61011

*"If you can't measure it,*
*nor assign it an exact numerical value, nor express*
*it in numbers,*
*then your knowledge is of a meager and*
*unsatisfactory kind"*

(attributed to William Thomson (1824-1907), aka Lord Kelvin)

---

# What to measure?
# How to measure?
# How to interpret?

---

- The answers to the question of which algorithm works best on your specific data set or problem, or which input features to use, can only be found through rigorous experiments.

- This is because many ML algorithms are too complex for formal analysis, at least at the level of generality assumed by most theoretical treatments. As a result, empirical studies of the behaviour of machine learning algorithms must retain a central role.
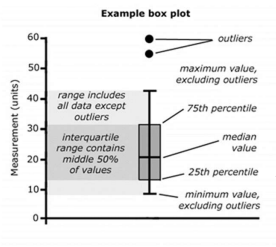
Excellent source: https://machinelearningmastery.com/controlled-experiments-in-machine-learning

---
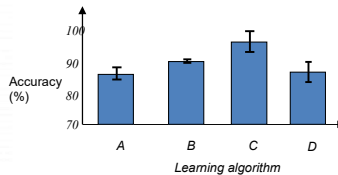
- This is a challenge for beginners who must learn some rigor.
- Three types of experiments:
- Choose-Features Experiments. When determining what data features (input variables) are most relevant to a model, the independent variables may be the input features and the dependent variable might be the estimated skill of the model on unseen data.
- Tune-Model Experiments. When tuning a machine learning model, the independent variables may be the hyperparameters of the learning algorithm and the dependent variable might be the estimated skill of the model on unseen data.
- Compare-Models Experiments. When comparing the performance of machine learning models, the independent variables may be the learning algorithms themselves with a specific configuration and the dependent variable is the estimated skill of the model on unseen data.
- What makes the experimental focus of applied machine learning so exciting is two-fold:
  - 1) Discovery. You can discover what works best for your specific problem and data.
  - 2) Contribution. You can make broader discoveries in the field, without any specialized knowledge other than rigorous and systematic experimentation.

Excellent source: https://machinelearningmastery.com/controlled-experiments-in-machine-learning

- Cross-Validation is a way to monitor stability
- Check always confidence intervals
- ROC-Analysis, particularly for imbalanced data
  - (We rarely have I.I.D. data !!!)
  - Check: Accuracy, Training Time, Space complexity (how much memory is needed), Interpretability = how can we explain why it does what it does!
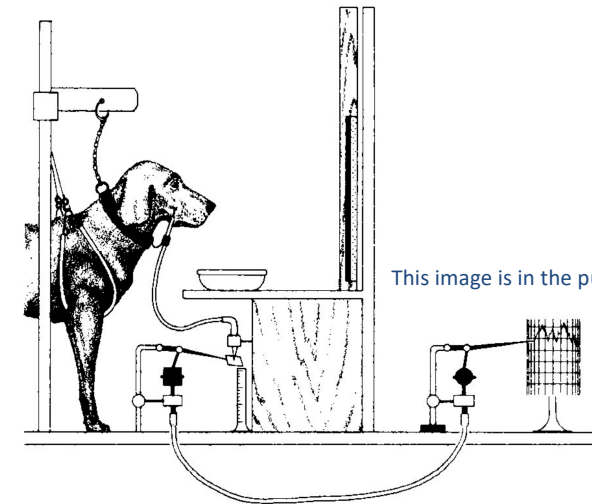  - Mutual Information, see MacKay, Section 44.5



Example box plot

Confusion Matrix

This image is in the public domain

- Machine learning experiments pose questions about models that we try to answer by means of measurements on data. We may ask questions including:
- Which learning algorithm provides the best model for data from Domain D?
- How does the model m perform on data from D?
- What model has the best performance on D?
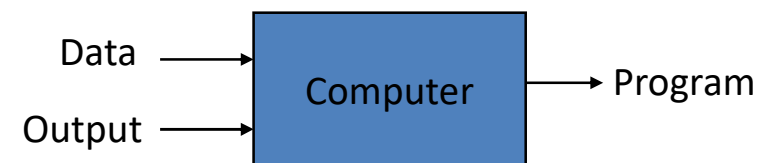- How do you benchmark? What is the ground truth?

## Traditional Programming



Data → Computer → Output
Program →

## Machine Learning = Learning from Data



Data → Computer → Program
Output →

# There is
# no free lunch!

Wolpert, D. H. & Macready, W. G. 1997. No free lunch theorems for optimization. IEEE Transactions on Evolutionary Computation, 1, (1), 67-82.

- Scalability
- Predictive accuracy = Hit rate
- Weighted (cost-sensitive) accuracy
- Speed (on model building and predicting)
- Robustness (one weakness in iML-approach)
- Precision/Recall (F-Measure, Break Even Point)
- Area under the ROC (see next slides)

Japkowicz, N. & Shah, M. 2011. Evaluating learning algorithms: a classification perspective, Cambridge University Press.

- There are many datasets for testing machine learning algorithms, just some examples:
- https://www.kaggle.com
- http://archive.ics.uci.edu/ml/datasets.html (UCI Machine Learning Repository)
- http://image-net.org
- http://yann.lecun.com/exdb/mnist (handwritten digit database)
- https://data.medicare.gov/

http://hci-kdd.org/open-data-sets/

- **Question: is 99% accuracy good?**
- **Answer: It depends on the problem!**

- **Accuracy** = error rate of correct/incorrect predictions made by the model over a data set (cf. coverage).

- **Precision** = precision (positive predictive value) is the fraction of retrieved instances that are relevant, while **Recall** (aka sensitivity) is the fraction of relevant instances that are retrieved

- **Reliability** = basically the "consistency" or "repeatability"

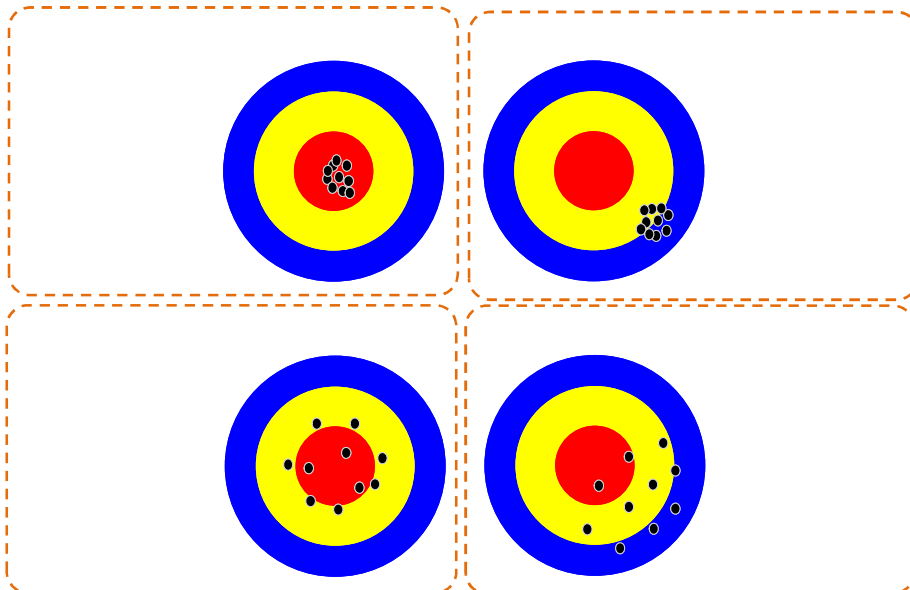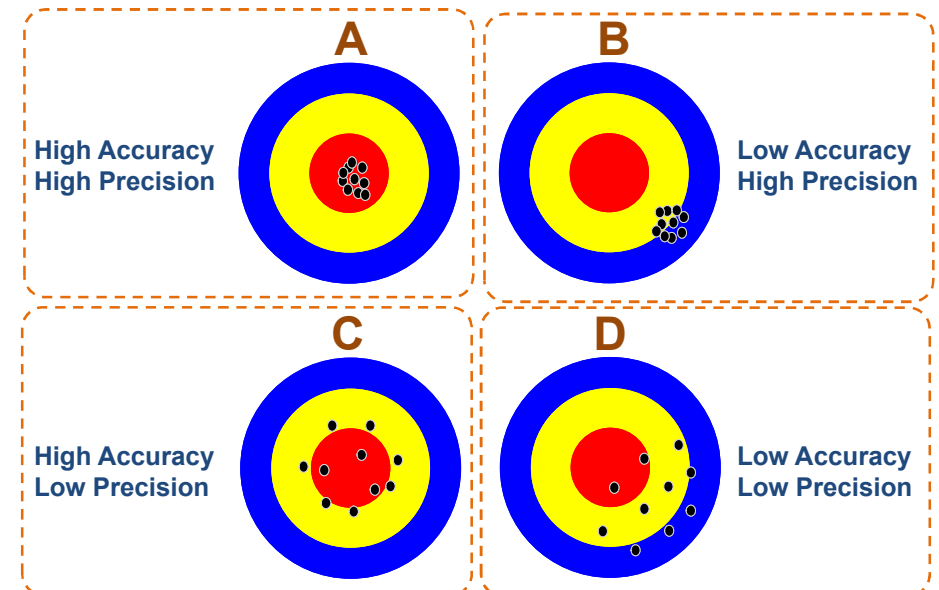- **Validity** = generally, to get valid conclusions
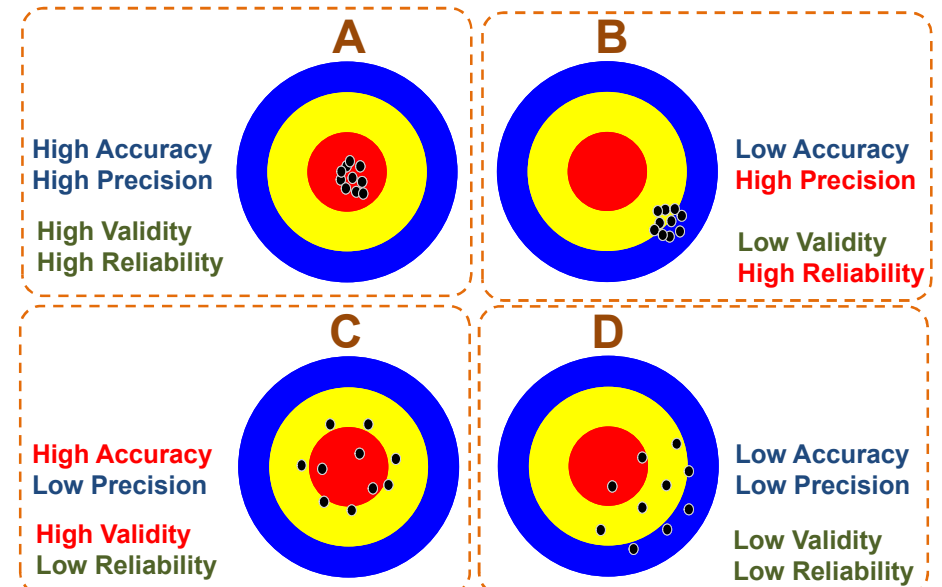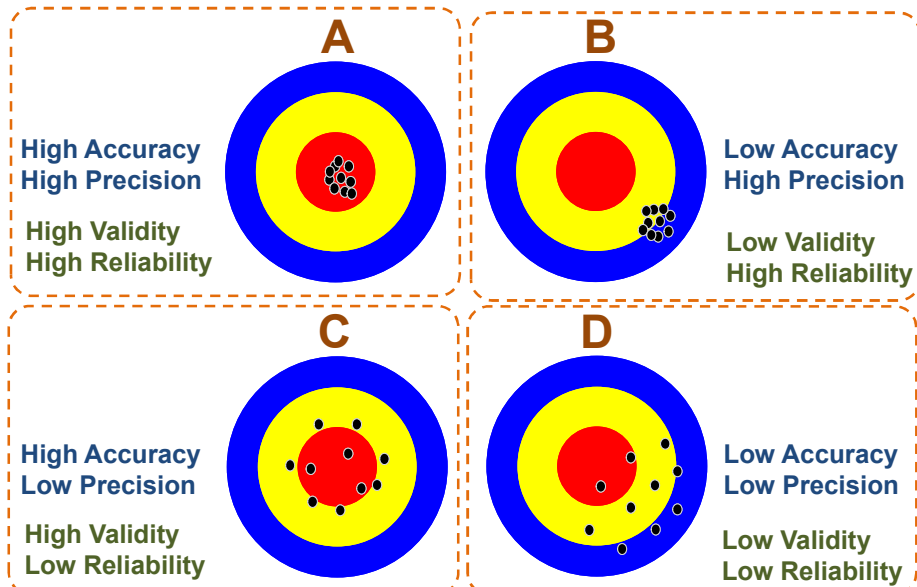
---

## Accuracy          Validity

## Precision          Reliability

---

---

A — High Accuracy High Precision
B — Low Accuracy High Precision
C — High Accuracy Low Precision
D — Low Accuracy Low Precision

**A**

High Accuracy
High Precision

High Validity
High Reliability

**B**

Low Accuracy
High Precision

Low Validity
High Reliability

**C**

High Accuracy
Low Precision

High Validity
Low Reliability

**D**

Low Accuracy
Low Precision

Low Validity
Low Reliability

**A**

High Accuracy
High Precision

High Validity
High Reliability

**B**

Low Accuracy
High Precision

Low Validity
High Reliability

**C**

High Accuracy
Low Precision

High Validity
Low Reliability

**D**

Low Accuracy
Low Precision

Low Validity
Low Reliability

|  |  | True Class | |
|---|---|---|---|
|  |  | Positive | Negative |
| Predicted Class | Positive | True Positive Count (TP) | False Positive Count (FP) |
|  | Negative | False Negative Count (FN) | True Negative Count (TN) |

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

$$True\ Positive\ Rate = \frac{TP}{TP + FN}$$

$$True\ Negative\ Rate = \frac{TN}{TN + FP}$$

$$Precision = \frac{TP}{TP + FP} \qquad Recall = \frac{TP}{TP + FN}$$

Turban, E., Sharda, R., Delen, D. & Efraim, T. 2007. Decision support and business intelligence systems, Pearson Education.

Bradley, A. P. 1997. The use of the area under the ROC curve in the evaluation of machine learning algorithms. *Pattern Recognition,* 30, (7), 1145-1159.

For a detailed explanation refer to: Fawcett, T. 2006. An introduction to ROC analysis. Pattern recognition letters, 27, (8), 861-874.

# 08 #KandinksyPatterns – our "Swiss-Knife" for the study of explainable AI

Andreas Holzinger, Georg Langs, Helmut Denk, Kurt Zatloukal & Heimo Mueller 2019. Causability and Explainability of AI in Medicine. Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery, doi:10.1002/widm.1312.

Holzinger, A., Geierhofer, R. & Errath, M. 2007. Semantische Informationsextraktion in medizinischen Informationssystemen. *Informatik Spektrum, 30, (2), 69-78.*
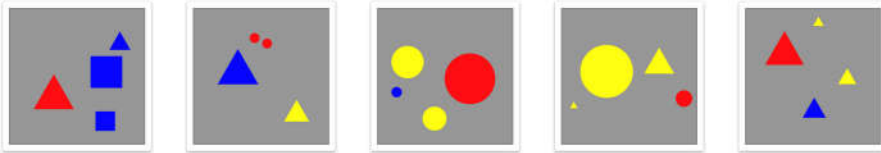
David H. Hubel & Torsten N. Wiesel 1962. Receptive fields, binocular interaction and functional architecture in the cat's visual cortex. The Journal of Physiology, 160, (1), 106-154, doi:10.1113/jphysiol.1962.sp006837.

Komposition VIII, 1923, Solomon R. Guggenheim Museum, New York.Source: https://de.wikipedia.org/wiki/Wassily_Kandinsky
This images are in the public domain.

- … a square image containing $1$ to $n$ geometric objects.
- Each object is characterized by its shape, color, size and position within this square.
- Objects do not overlap and are not cropped at the border.
- All objects must be easily recognizable and clearly distinguishable by a human observer.

- about a Kandinsky Figure $k$ is …
- either a mathematical function $s(k) \rightarrow B$; with $B$ $(0,1)$
- or a *natural language statement* which is true or false

- Remark: The evaluation of a natural language statement is always done in a specific context. In the followings examples we use **well known concepts from human perception** and linguistic theory.
- If $s(k)$ is given as an algorithm, it is essential that the function is a pure function, which is a computational analogue of a mathematical function.
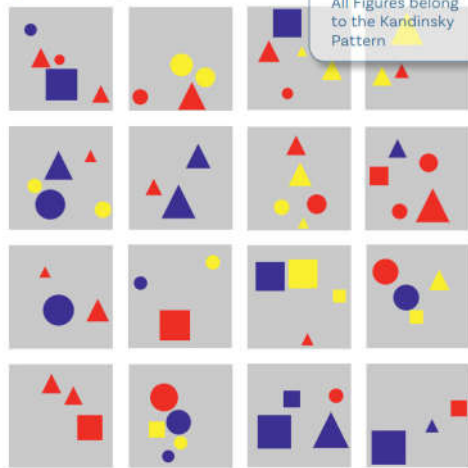
- … is defined as the subset of all possible Kandinsky Figures k with $s(k) \to 1$ or the natural language statement is true.
- $s(k)$ and a natural language statement are equivalent, if and only if the resulting Kandinsky Patterns contains the same Kandinsky Figures.
- $s(k)$ and the natural language statement are defined as the **Ground Truth** of a Kandinsky Pattern



*"… the Kandinsky Figure has two pairs of objects with the same shape, in one pair the objects have the same color, in the other pair different colors, two pairs are always disjunct, i.e. they don't share a object …".*

A Colour

B Shape

C Quantity

D Arrangement

E Gestalt

F Domain

A) True

B) False

C) Counterfactual

Part of the pattern

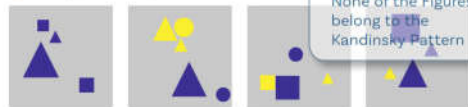All Figures belong to the Kandinsky Pattern

**Hypothesis 1**

It only contains circles and triangles.

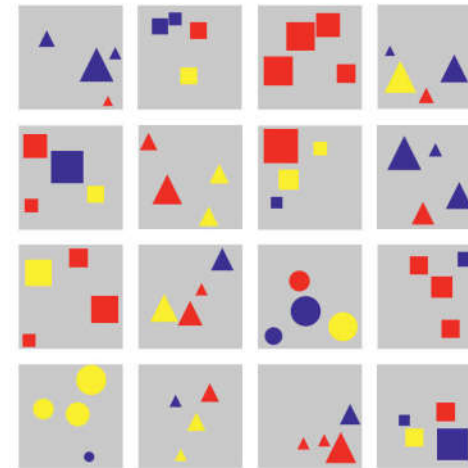**Hypothesis 2**

It contains at least a red object. ✓

Not part of the pattern

None of the Figures belong to the Kandinsky Pattern
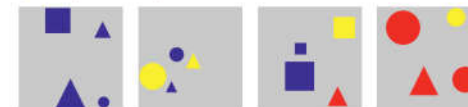
---

Part of the pattern

S2   **Basic Pattern 2**
Title: **All of Same Shape** ->
All objects have the same shape.
Hint: Don't be distracted by the colors

Not part of the pattern

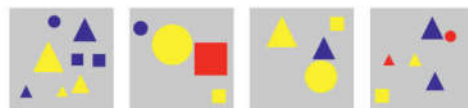---

Part of the pattern

S8   **Basic Pattern 8**
Title: **Mickey Mouse** ->
Every figure contains a pattern which is made out of a big yellow circle and two smaller blue ones and looks like a Mickey Mouse.

Not part of the pattern

---

# 09 Sample Questions and Conclusion

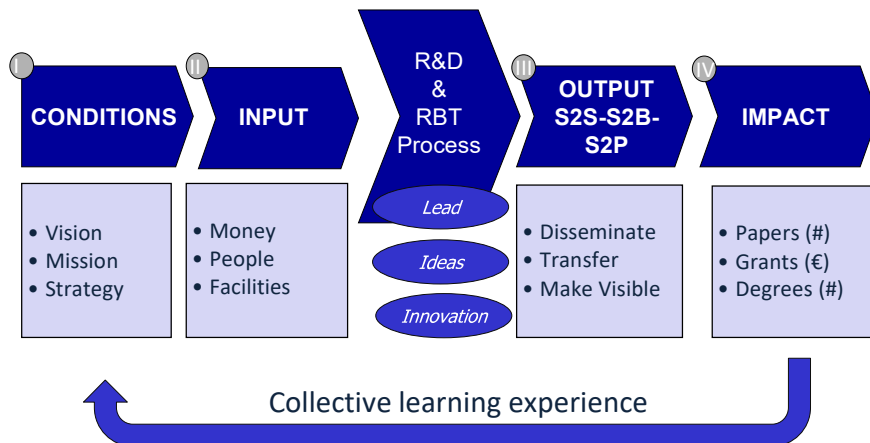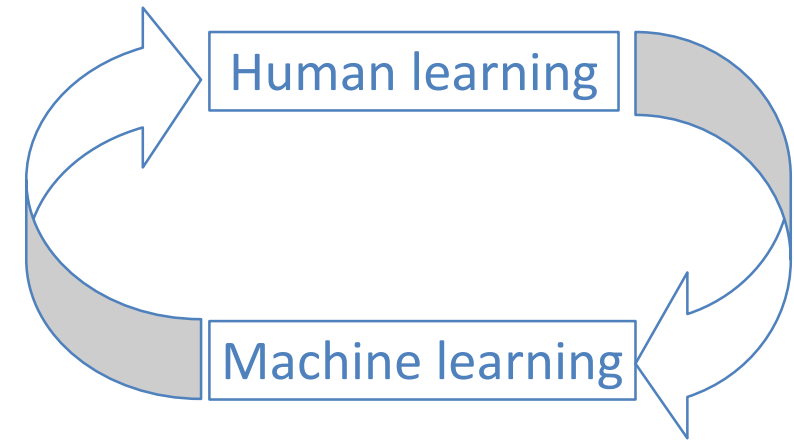# Human-AI collaboration will sustainably influence the way science is done

# How can I use visual representations of abstract data to amplify the acquisition of knowledge?

- 1) Given a set of (complex) data
- 2) Set a hypothesis
- 3) Extract information
- 4) Discover hidden knowledge
- 5) Support your previous set hypothesis
- **Machine intelligence + Human intelligence**
- = powerful methods for many sciences
- Application e.g. in many domains, e.g. medicine and health, education, psychology, industry, etc

- How can we transfer learned representations to improve learning in other tasks/domains?
- Which learning algorithm should be used when?
- Can machine learning theories help to understand human learning and vice versa?
- Machine Learning vs. Human Learning: role of motivation, emotion, forgetting, … ?
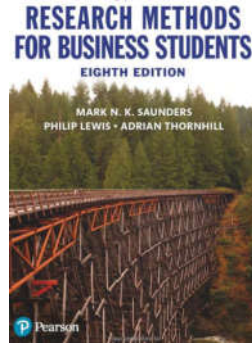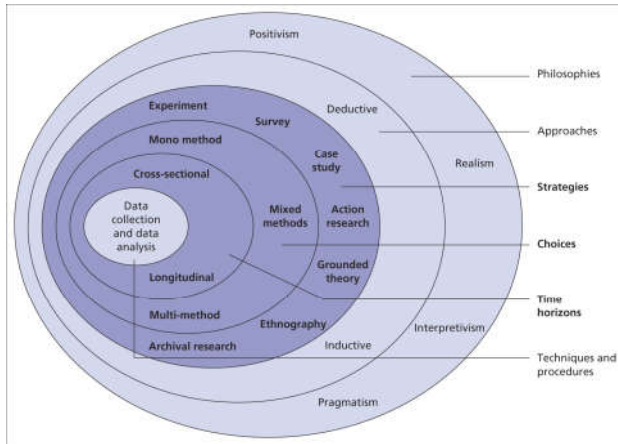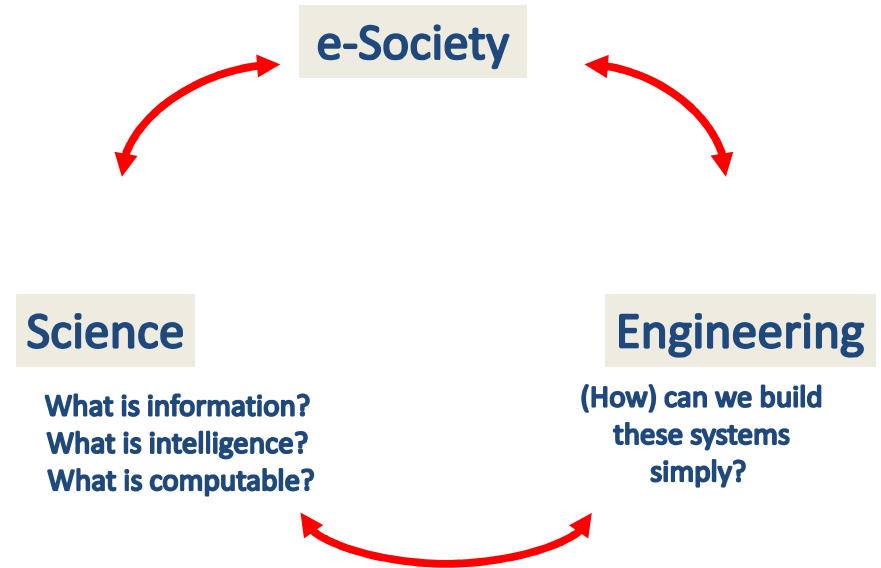- How can we use self-supervised learning with multiple sensory input?

https://ai100.stanford.edu/

# Appendix

---

## Human learning

## Machine learning

---

**CONDITIONS** → **INPUT** → **R&D & RBT Process** → **OUTPUT S2S-S2B-S2P** → **IMPACT**

- Lead
- Ideas
- Innovation

| CONDITIONS | INPUT | OUTPUT | IMPACT |
|---|---|---|---|
| • Vision | • Money | • Disseminate | • Papers (#) |
| • Mission | • People | • Transfer | • Grants (€) |
| • Strategy | • Facilities | • Make Visible | • Degrees (#) |

Collective learning experience

Holzinger, A. 2011. *Successful Management of Research and Development,* Norderstedt, BoD.

---

Essential Tools for Scientific Machine Learning and Scientific AI

http://www.stochasticlifestyle.com/the-essential-tools-of-scientific-machine-learning-scientific-ml

## Slide 101

## Slide 102

### e-Society

### Science

**What is information?**
**What is intelligence?**
**What is computable?**

### Engineering

**(How) can we build these systems simply?**

## Slide 103

| Kleine Einheiten | | | Große Einheiten | | |
|---|---|---|---|---|---|
| $10^{-3}$ | milli | m | $10^{3}$ | kilo | k |
| $10^{-6}$ | micro | µ | $10^{6}$ | mega | M |
| $10^{-9}$ | nano | n | $10^{9}$ | giga | G |
| $10^{-12}$ | pico | p | $10^{12}$ | tera | T |
| $10^{-15}$ | femto | f | $10^{15}$ | peta | P |
| $10^{-18}$ | atto | a | $10^{18}$ | exa | E |
| $10^{-21}$ | zepto | z | $10^{21}$ | zetta | Z |
| $10^{-24}$ | yocto | y | $10^{24}$ | yotta | Y |

## Slide 104

**Please note: Computers change constantly …**

- Old dream of mankind: using technology to augment human capabilities for structuring, retrieving and managing information and decision support:

- **… challenges for HCI …**



*Harper, Rodden, Rogers, Sellen (2008)*

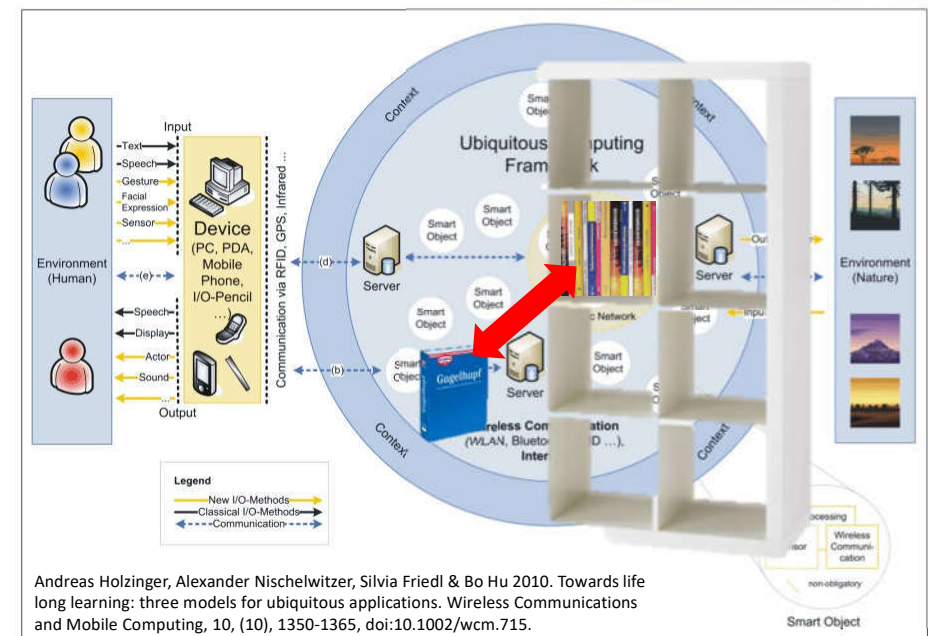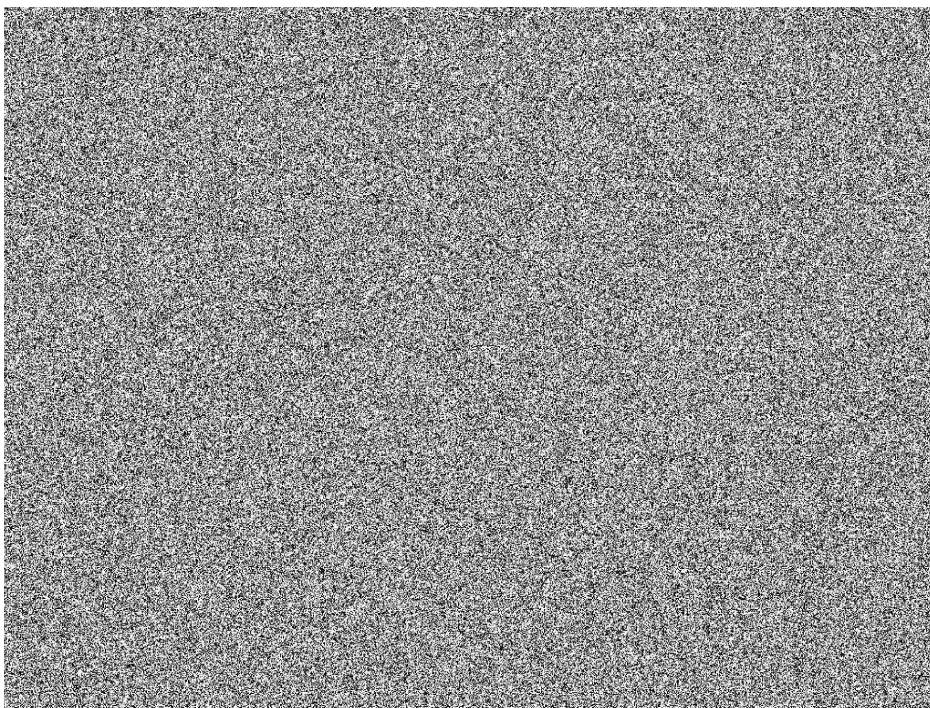Andreas Holzinger, Alexander Nischelwitzer, Silvia Friedl & Bo Hu 2010. Towards life long learning: three models for ubiquitous applications. Wireless Communications and Mobile Computing, 10, (10), 1350-1365, doi:10.1002/wcm.715.

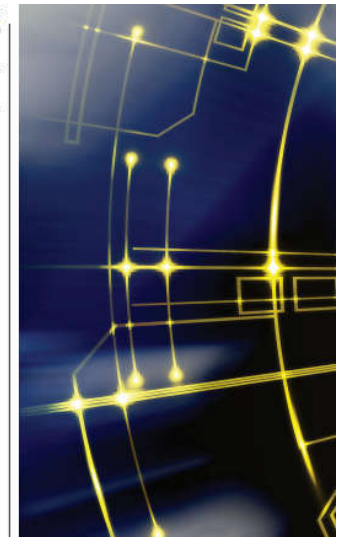---

## Let the user interactively manipulate the data

- **Focus Selection** = via direct manipulation and selection tools, e.g. multi-touch (in data space a n-dim location might be indicated); see a recent work by Randy Goebel
- **Attention Routing** = anomaly detection, draws people's attention to interesting areas to start their analyses;
- **Extent Selection** = specifying extents for an interaction, e.g. via a vector of values (a range for each data dimension or a set of constraints;
- **Interaction type selection** = e.g. a pair of menus: one to select the space, and the other to specify the general class of the interaction;
- **Interaction level selection** = e.g. the magnitude of scaling that will occur at the focal point (via a slider, along with a reset button);
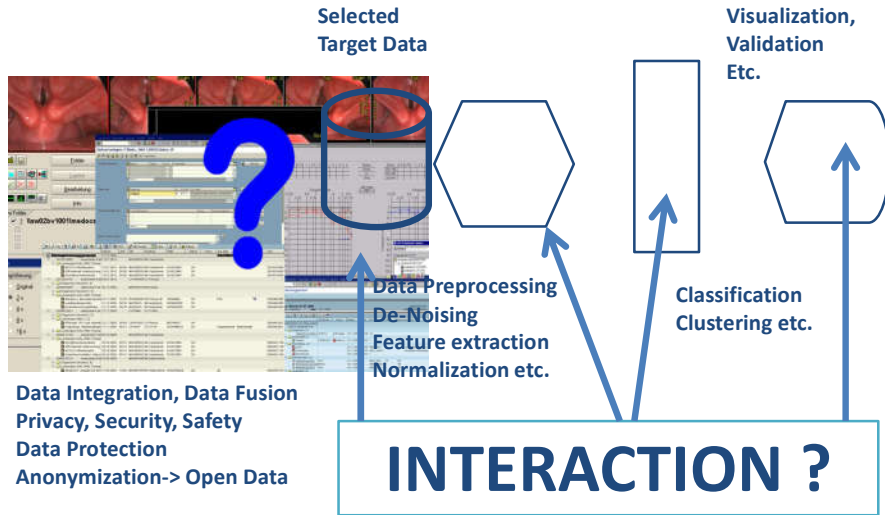
---

## Time is essential !!



DOI:10.1145/2500886

**Timing analysis for hard real-time systems.**

BY REINHARD WILHELM AND DANIEL GRUND

# Computation Takes Time, But How Much?

Wilhelm, R. & Grund, D. 2014. Computation takes time, but how much? *Communications of the ACM, 57, (2), 94-103.*

**Selected
Target Data**

**Visualization,
Validation
Etc.**



**?**

**Data Preprocessing
De-Noising
Feature extraction
Normalization etc.**

**Classification
Clustering etc.**

**Data Integration, Data Fusion
Privacy, Security, Safety
Data Protection
Anonymization-> Open Data**

**INTERACTION ?**

Holzinger, A. & Zupan, M. 2013. KNODWAT: A scientific framework application for testing
knowledge discovery methods for the biomedical domain. *BMC Bioinformatics, 14, (1), 191.*

---

# Our central hypothesis:
# Information bridges this gap

Simonic, K.-M. & Holzinger, A. (2010) Zur Bedeutung von Information in der Medizin. *OCG Journal, 35, **1, 8.***