

# Interactive Adaptive Learning ECML PKDD Tutorial

A. Calma, A. Holzinger, D. Kottke, G. Kreml, V. Lemaire

## Part 2: From Interactive ML to Explainable AI (ex-AI)

**Andreas Holzinger**

Human-Centered AI Lab (Holzinger Group)

Institute for Medical Informatics/Statistics, Medical University Graz, Austria  
and

Explainable AI-Lab, Alberta Machine Intelligence Institute, Edmonton, Canada



**HCAI**  
HUMAN-CENTERED.AI



**@aholzin #KandinskyPatterns**

In this introductory tutorial we ...

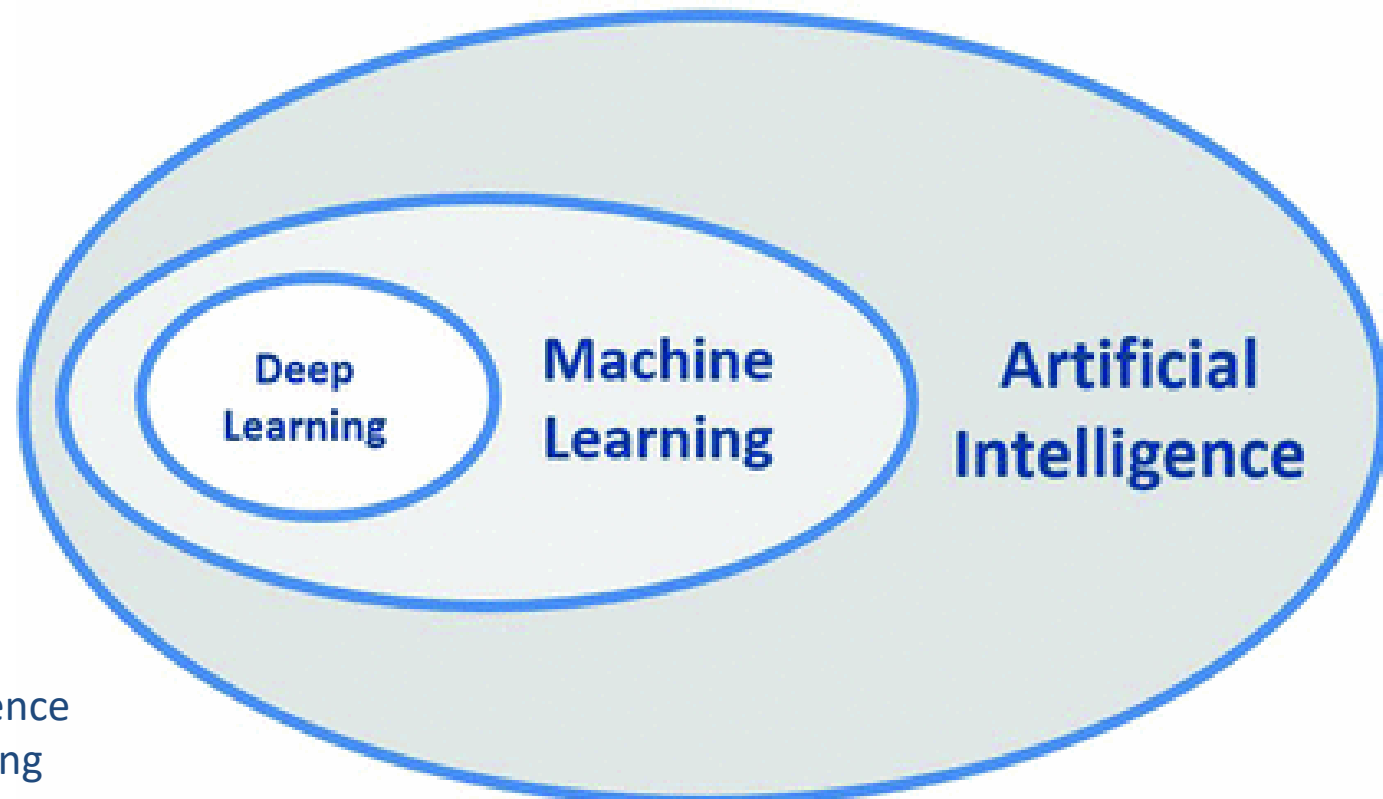
- 1) learn a few basics about “explainable AI” and get an overview on a few current state-of-the art methods,
- 2) see the necessity to go “beyond explainable AI” which we need as a basis for future Human-AI interfaces, and
- 3) will work with an open explanation environment [1], the so-called #KandinskyPatterns [2], which can also be used as an IQ-Test for machines [3].

*But let us start first with a preamble and motivation why we would need transparency, explainability, and interpretability*

[1] Links to GitHub etc. via the Project Homepage: <https://human-centered.ai/project/kandinsky-patterns>

[2] Heimo Müller & Andreas Holzinger 2019. Kandinsky Patterns. arXiv:1906.00657

[3] Andreas Holzinger, Michael Kickmeier-Rust & Heimo Mueller 2019. KANDINSKY Patterns as IQ-Test for machine learning. Springer Lecture Notes LNCS 11713, pp. 1-14, doi:10.1007/978-3-030-29726-8\_1



AI = Artificial Intelligence

ML = Machine Learning

DL = Deep Learning

aML = automatic (autonomous) ML

iML = interactive ML

HCI = Human-Computer Interaction

HCAI = Human-centered AI

HAI = Human AI Interfaces

KDD = Knowledge Discovery from Data

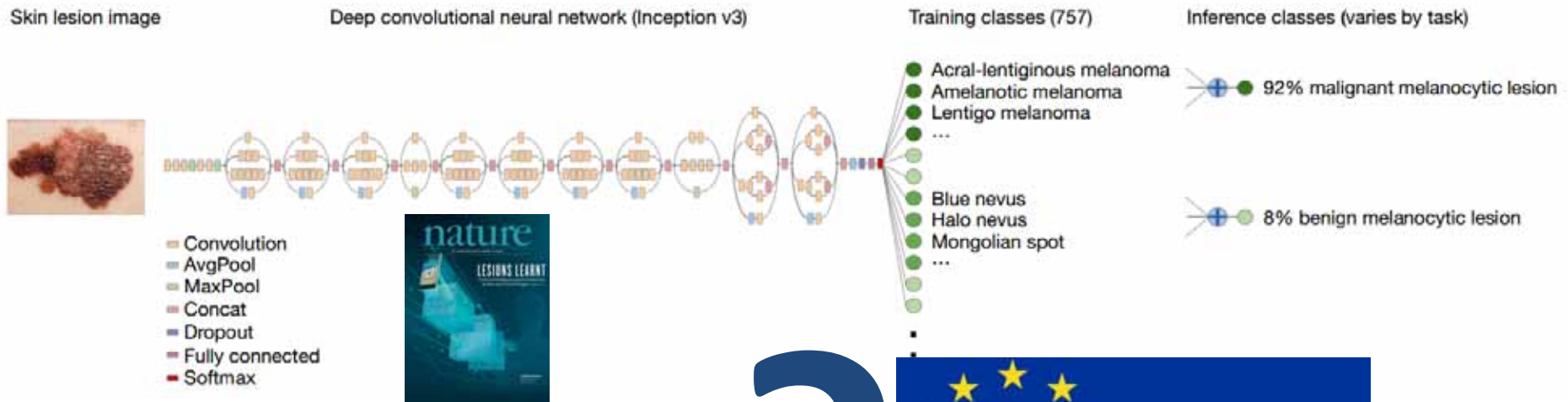
Ex-AI = explainable AI (also XAI)

Andreas Holzinger, Peter Kieseberg, Edgar Weippl & A Min Tjoa 2018.  
 Current Advances, Trends and Challenges of Machine Learning and  
 Knowledge Extraction: From Machine Learning to Explainable AI.  
 Springer Lecture Notes in Computer Science LNCS 11015. Cham:  
 Springer, pp. 1-8, doi:10.1007/978-3-319-99740-7\_1

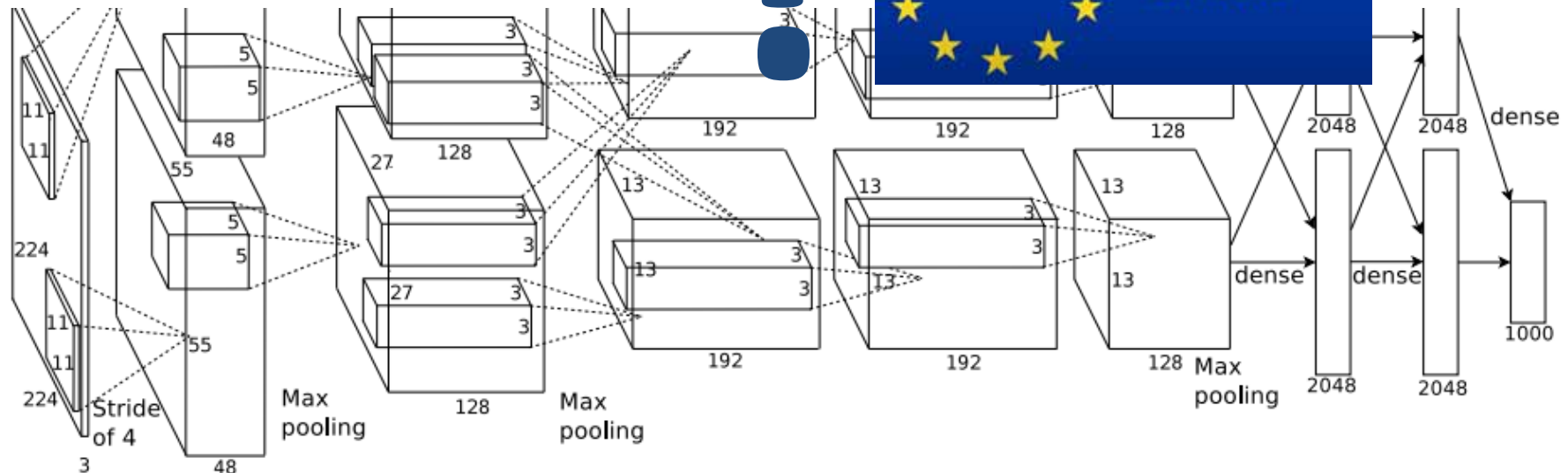


- **High-dimensional**
- **Non-convex**
- **Resource and data hungry**
- **Lacking interpretability ...**

Remember Richard Feynman on the universe: It's not complicated, it's just a lot ...



Esteva, A., Kuprel, B., Novoa, R. A., Ko, J., Swetter, S. M., Blau, N. P., Thaler, H. M. (2016). Classification of skin cancer with deep neural networks. *Nature*, 542, 393-396. doi:10.1038/nature21056.





+ .007 ×



=



$x$

“panda”

57.7% confidence

$\text{sign}(\nabla_x J(\theta, x, y))$

“nematode”

8.2% confidence

$x + \epsilon \text{sign}(\nabla_x J(\theta, x, y))$

“gibbon”

99.3 % confidence

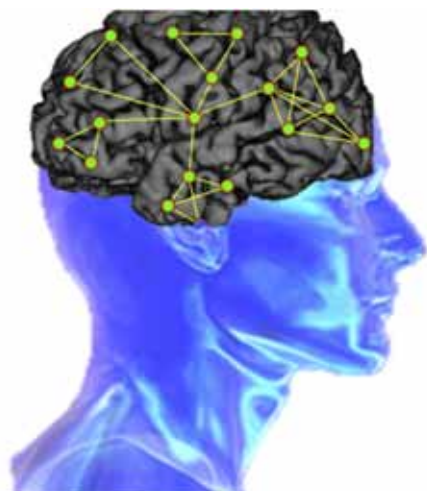
Ian Goodfellow, Jonathon Shlens & Christian Szegedy 2014. Explaining and harnessing adversarial examples. arXiv:1412.6572

# Urgent need for explainable AI !

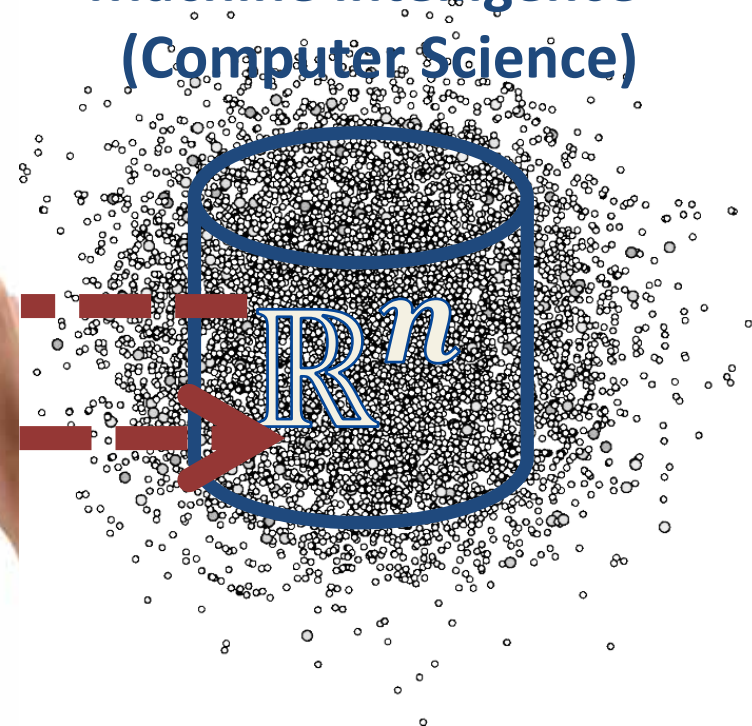
- **01 Introduction – probabilistic**
- **02 Few words about the application area health**
- **03 Probabilistic Learning**
- **04 Automatic ML (human-out-of-loop)**
- **05 interactive machine learning (human-in-loop)**
- **06 Methods of “explainable AI” (selection)**
- **07 Towards Human Interpretable Models**
- **08 Digression: Causality Learning**
- **09 IQ-Tests for Machines: #KandinskyPatterns**

# Our goal is that human values are aligned to ensure responsible machine learning

Human intelligence  
(Cognitive Science)



Machine intelligence  
(Computer Science)



Andreas Holzinger 2013. Human-Computer Interaction and Knowledge Discovery (HCI-KDD): What is the benefit of bringing those two fields to work together? In: Lecture Notes in Computer Science LNCS 8127. pp. 319-328, doi:10.1007/978-3-642-40511-2\_22.



# “Solve intelligence – then solve everything else”



Demis Hassabis, 22 May 2015

The Royal Society,  
Future Directions of Machine Learning Part 2



<https://youtu.be/XAbLn66iHcQ?t=1h28m54s>

- 1) learn from prior data
- 2) extract knowledge
- 3) generalize, i.e. guessing where a probability mass function concentrates
- 4) fight the curse of dimensionality
- 5) disentangle **underlying explanatory factors of data**, i.e.
- 6) **understand** the data in the **context** of an application domain

# Our goal is: Understanding Context !

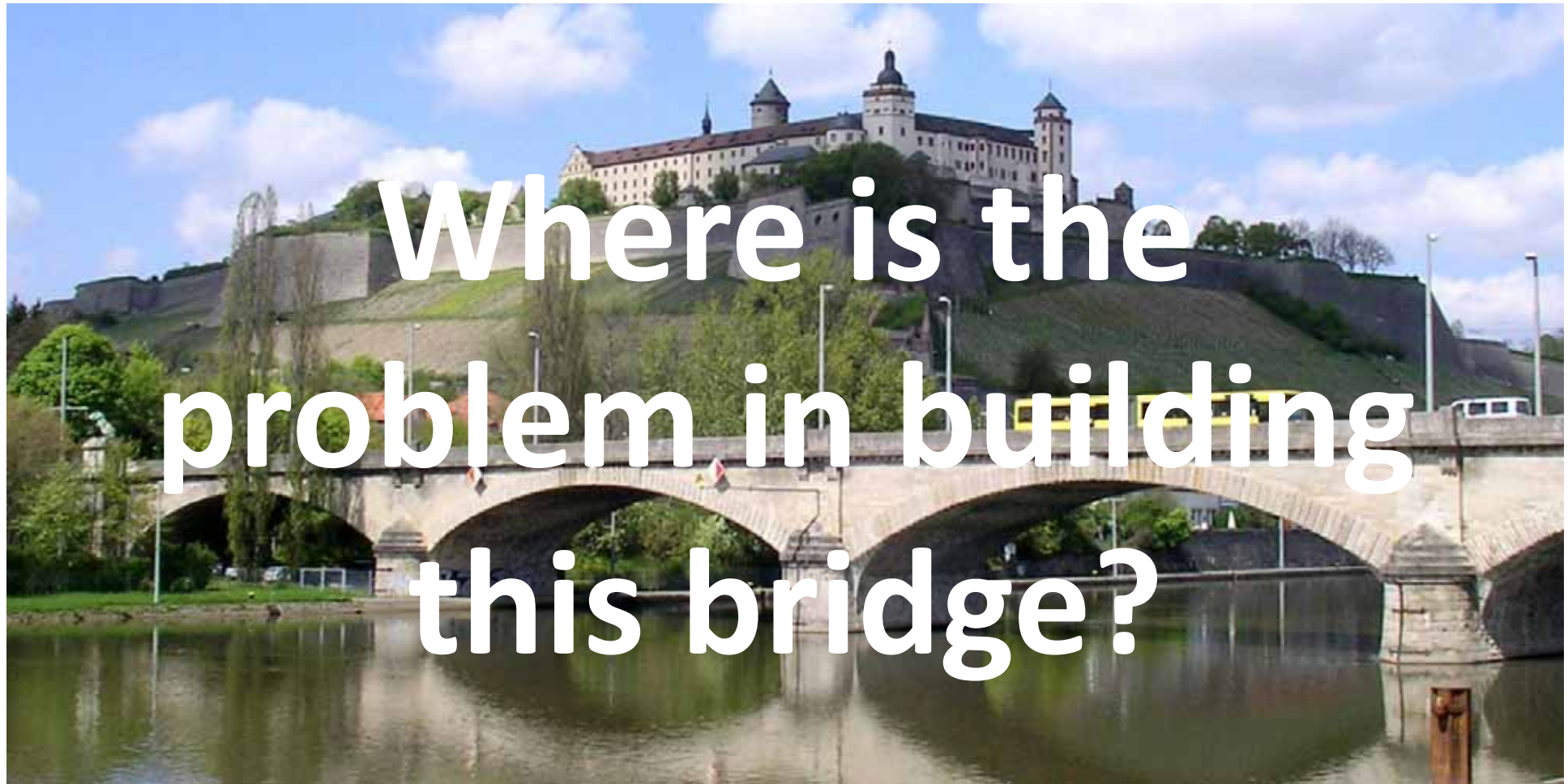


# Why is this application area complex ?



# Our central hypothesis: Information may bridge this gap

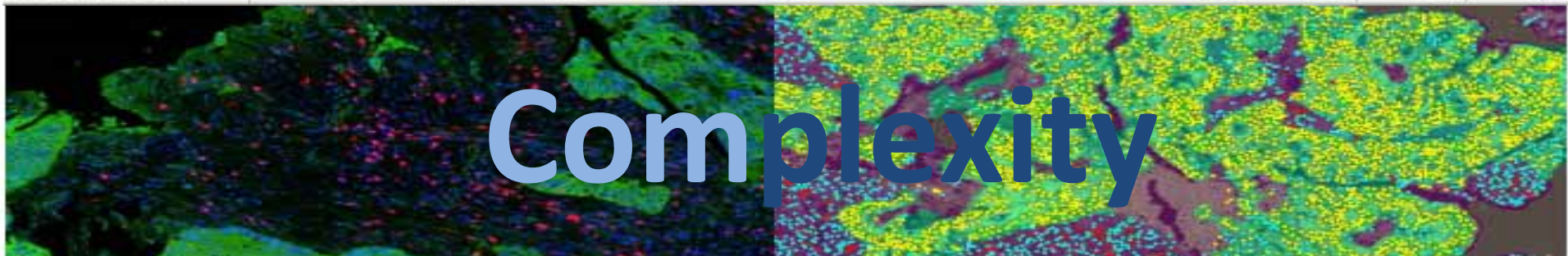
Andreas Holzinger & Klaus-Martin Simonic (eds.) 2011. Information Quality in e-Health. Lecture Notes in Computer Science LNCS 7058, Heidelberg, Berlin, New York: Springer, doi:10.1007/978-3-642-25364-5.



Sequence alignment showing protein families: MARHY0478, FadL, TbuX, and TodX. The alignment is annotated with domain markers  $\rho_{12}$ ,  $\rho_{13}$ ,  $\rho_{14}$ ,  $\rho_{15}$ , and  $\rho_{16}$ . Below the alignment is a summary table of counts for various metrics.

|                     |     |      |     |     |     |     |     |     |     |      |      |      |      |         |     |
|---------------------|-----|------|-----|-----|-----|-----|-----|-----|-----|------|------|------|------|---------|-----|
| Total pos/pS        | 16  | 16   | 5   | 21  | 21  | 21  | 21  | 5   | 26  | 26   | 26   | 26   | 5    | 31      |     |
| Total Infusionen    | 8   | 116  | 8   | 125 | 125 | 125 | 125 | 42  | 166 | 166  | 17   | 183  | 0    | 191     |     |
| Total Meds (pos+iv) | 4   | 4    | 4   | 4   | 4   | 4   | 4   | 2   | 6   | 6    | 6    | 6    | 0    | 6       |     |
| Total Perfusoren    | 1   | 9    | 1   | 10  | 10  | 10  | 10  | 5   | 15  | 15   | 2    | 17   | 1    | 18      |     |
| Total Meds+Perfusor | 1   | 13   | 1   | 14  | 14  | 14  | 14  | 7   | 21  | 21   | 2    | 23   | 1    | 24      |     |
| Total Blut          |     |      |     |     |     |     |     |     |     |      |      |      |      |         |     |
| Total Harn          | 43  | 43   | 43  | 43  | 43  | 43  | 43  | 15  | 15  | 15   | 2    | 134  | 2    | 134     |     |
| Harnmenge/Zeit      |     |      | 10  | 4   | 4   | 4   | 4   | 19  | 4   | 4    |      |      |      | 134/ 24 |     |
| Harn/kg/Std         |     |      |     |     |     |     |     |     |     |      |      |      |      | 2,0     |     |
| Total Na-Darm       | 6   | 6    | 6   | 6   | 6   | 6   | 6   | 0   | 6   | 6    | 6    | 6    | 6    | 6       |     |
| Total Blut          |     |      |     |     |     |     |     |     |     |      |      |      |      |         |     |
| Total Ein           | 9   | 145  | 9   | 154 | 5   | 159 | 159 | 159 | 54  | 213  | 213  | 19   | 232  | 9       | 241 |
| Total Aus           | 49  | 49   | 49  | 40  | 89  | 89  | 89  | 89  | 29  | 118  | 118  | 118  | 22   | 140     |     |
| Nettobilanz 24h     | +96 | +105 | +70 | +70 | +70 | +70 | +70 | +95 | +95 | +114 | +101 | +101 | +106 | +18     |     |

**Heterogeneity**  
**Dimensionality**



**Complexity**  
**Uncertainty**

Andreas Holzinger, Matthias Dehmer & Igor Jurisica 2014. Knowledge Discovery and interactive Data Mining in Bioinformatics - State-of-the-Art, future challenges and research directions. Springer/Nature BMC Bioinformatics, 15, (S6), I1, doi:10.1186/1471-2105-15-S6-I1.



# 03 Probabilistic Learning

The true logic of this world is  
in the calculus of  
probabilities.

James Clerk Maxwell



Maxwell, J. C. (1850). Letter to Lewis Campbell;  
reproduced in L. Campbell and W. Garrett, *The  
Life of James Clerk Maxwell*, Macmillan, 1881.

**Probability  
theory is  
nothing but  
common  
sense reduced  
to calculation  
...**



**Pierre Simon de Laplace (1749-1827)**

- **1763:** Richard Price publishes post hum the work of Thomas Bayes (see next slide)
- **1781:** Pierre-Simon Laplace: Probability theory is nothing but common sense reduced to calculation ...
- **1812:** *Théorie Analytique des Probabilités*, now known as Bayes' Theorem
- **Hypothesis**  $h \in \mathcal{H}$  (uncertain quantities (Annahmen))
- **Data**  $d \in \mathcal{D}$  ... measured quantities (Entitäten)
- **Prior probability**  $p(h)$  ... probability that  $h$  is true
- **Likelihood**  $p(d|h)$  ... “how probable is the prior”
- **Posterior Probability**  $p(h|d)$  ... probability of  $h$  given  $d$



Pierre Simon de Laplace (1749-1827)

$$p(h|d) \propto p(d|h) * p(h) \qquad p(h|d) = \frac{p(d|h)p(h)}{p(d)}$$

$d$  ... data

$\mathcal{H} \dots \{H_1, H_2, \dots, H_n\}$

$\forall h, d \dots$

$h$  ... hypotheses

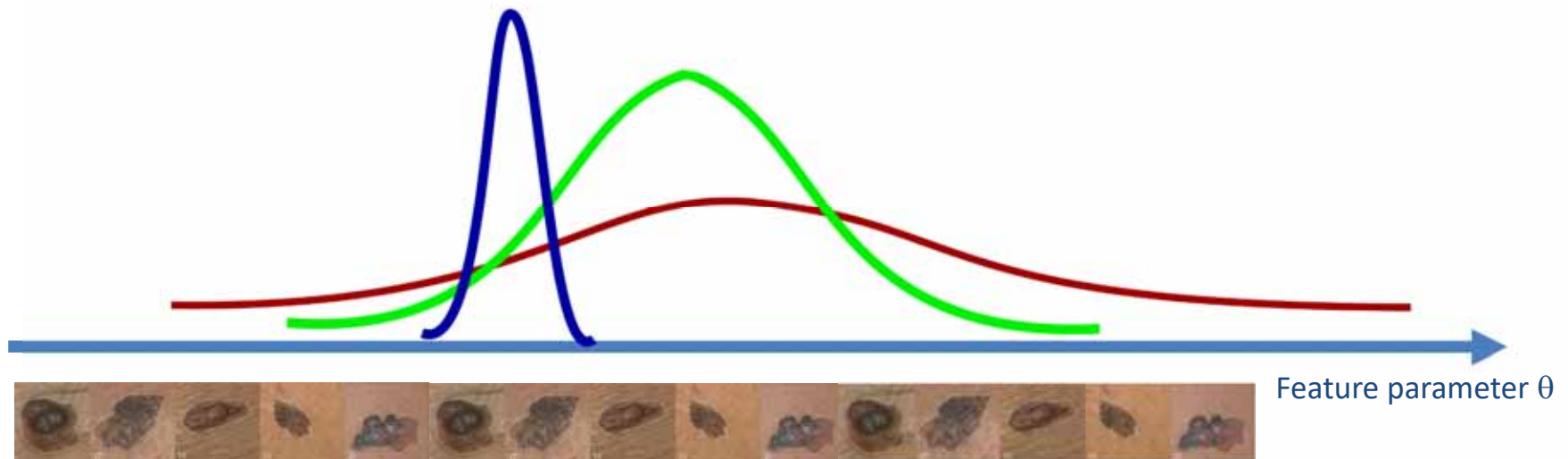
$$p(h|d) = \frac{p(d|h) * p(h)}{\sum_{h \in \mathcal{H}} p(d|h') p(h')}$$

Likelihood

Prior Probability

Posterior Probability

Problem in  $\mathbb{R}^n \rightarrow$  complex

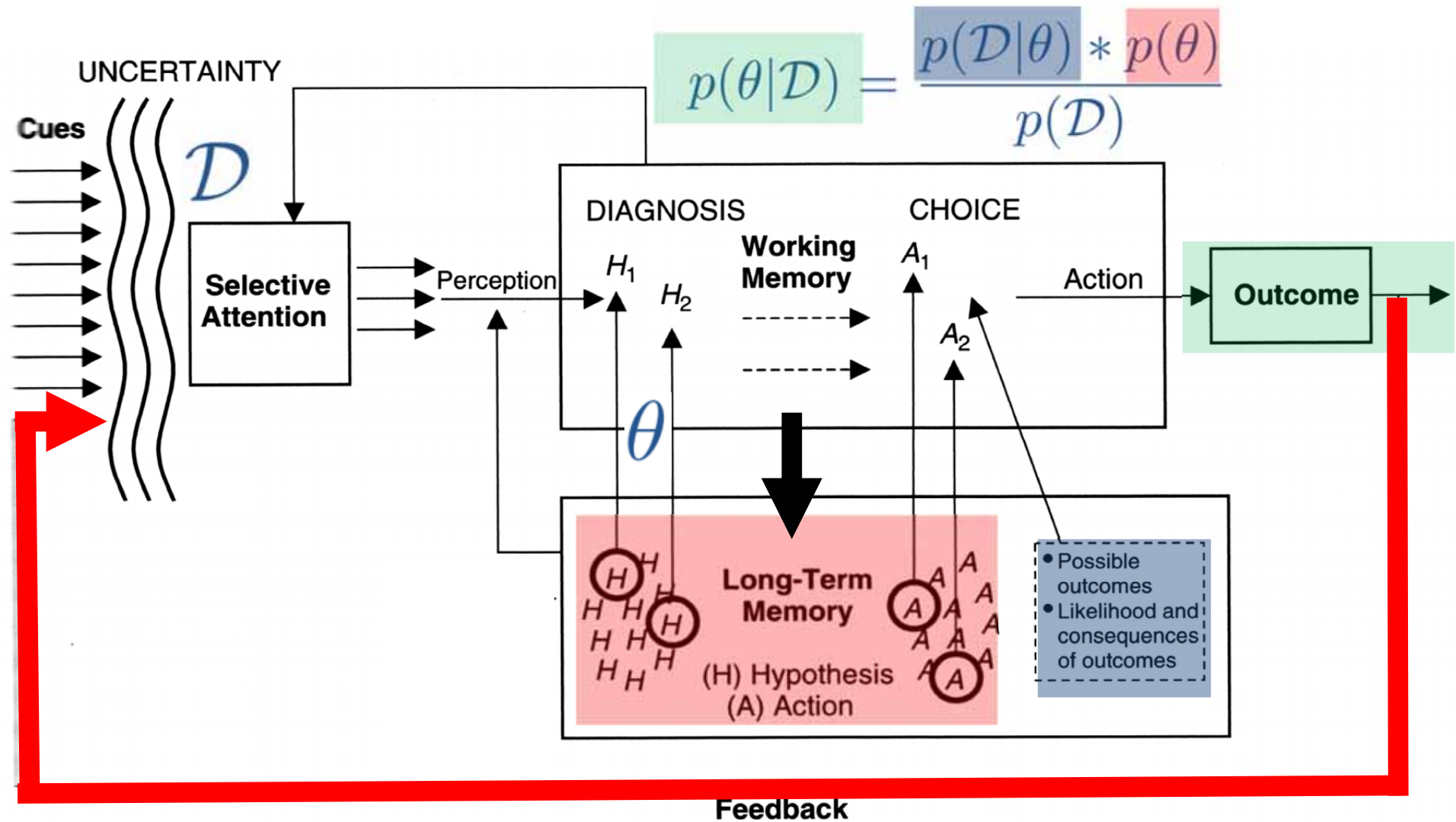


# Why is this relevant for medicine?

- Take patient information, e.g., observations, symptoms, test results, -omics data, etc. etc.
- Reach conclusions, and **predict** the future, e.g. how likely will the patient ...
- Prior = belief before making a particular observation
- Posterior = belief after making the observation and is the prior for the next observation – intrinsically incremental

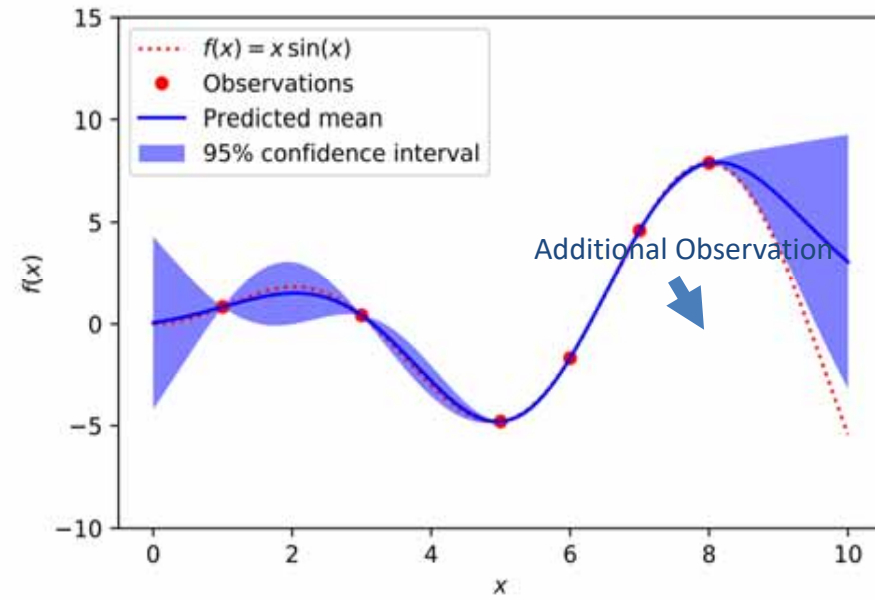
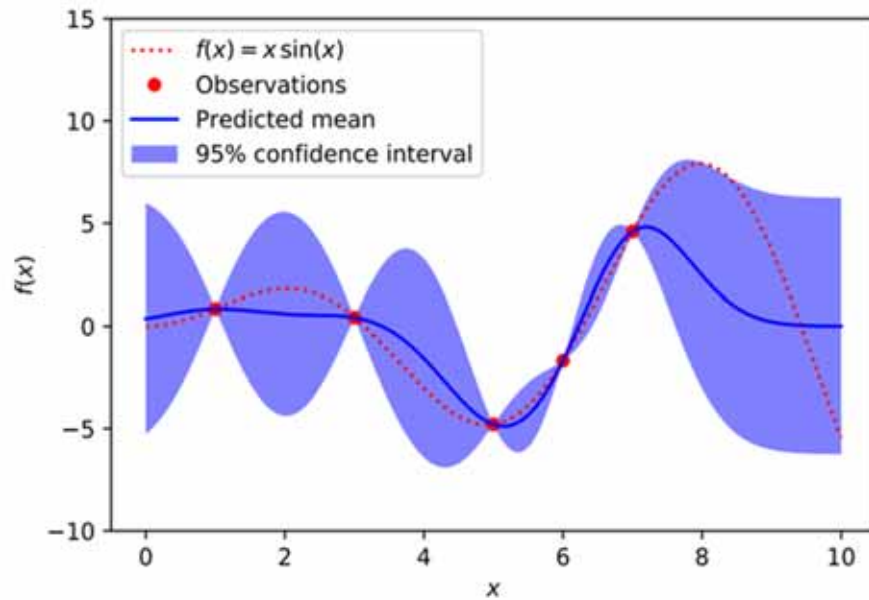

$$p(x_i|y_j) = \frac{p(y_j|x_i)p(x_i)}{\sum p(x_i, y_j)p(x_i)}$$

# Why is this relevant for decision support?



Wickens, C. D. (1984) *Engineering psychology and human performance*. Columbus (OH), Charles Merrill, modified by Holzinger, A.

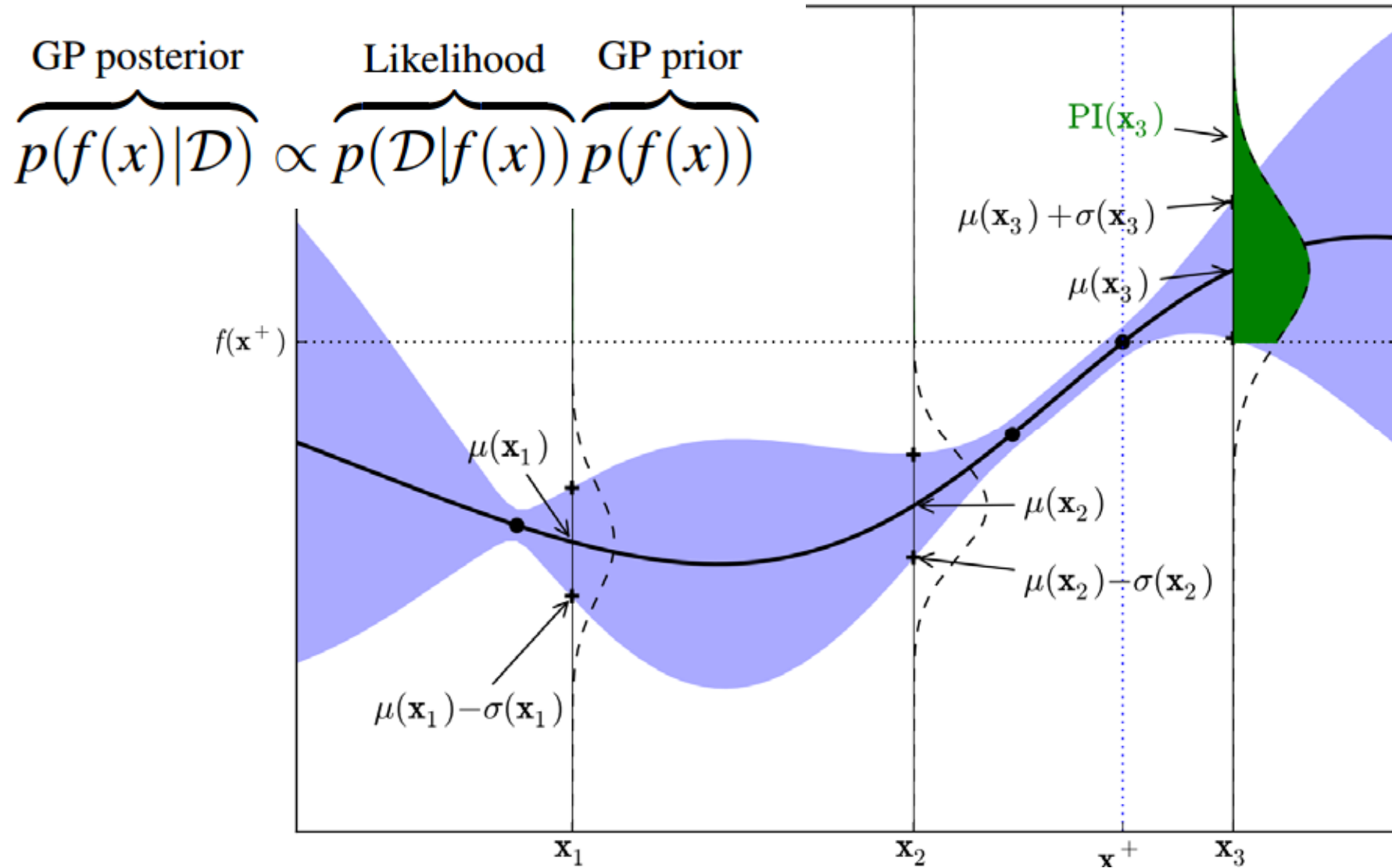




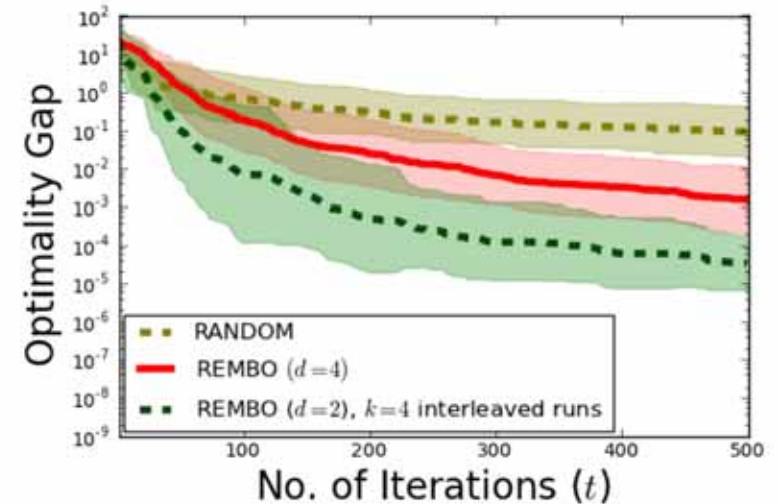
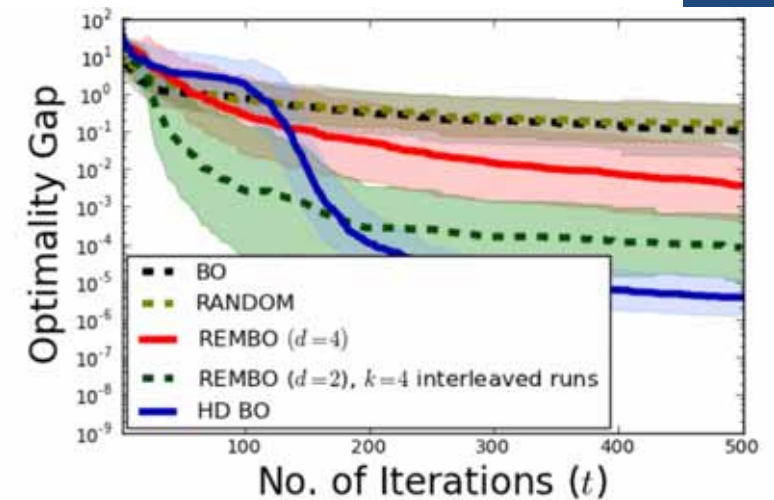
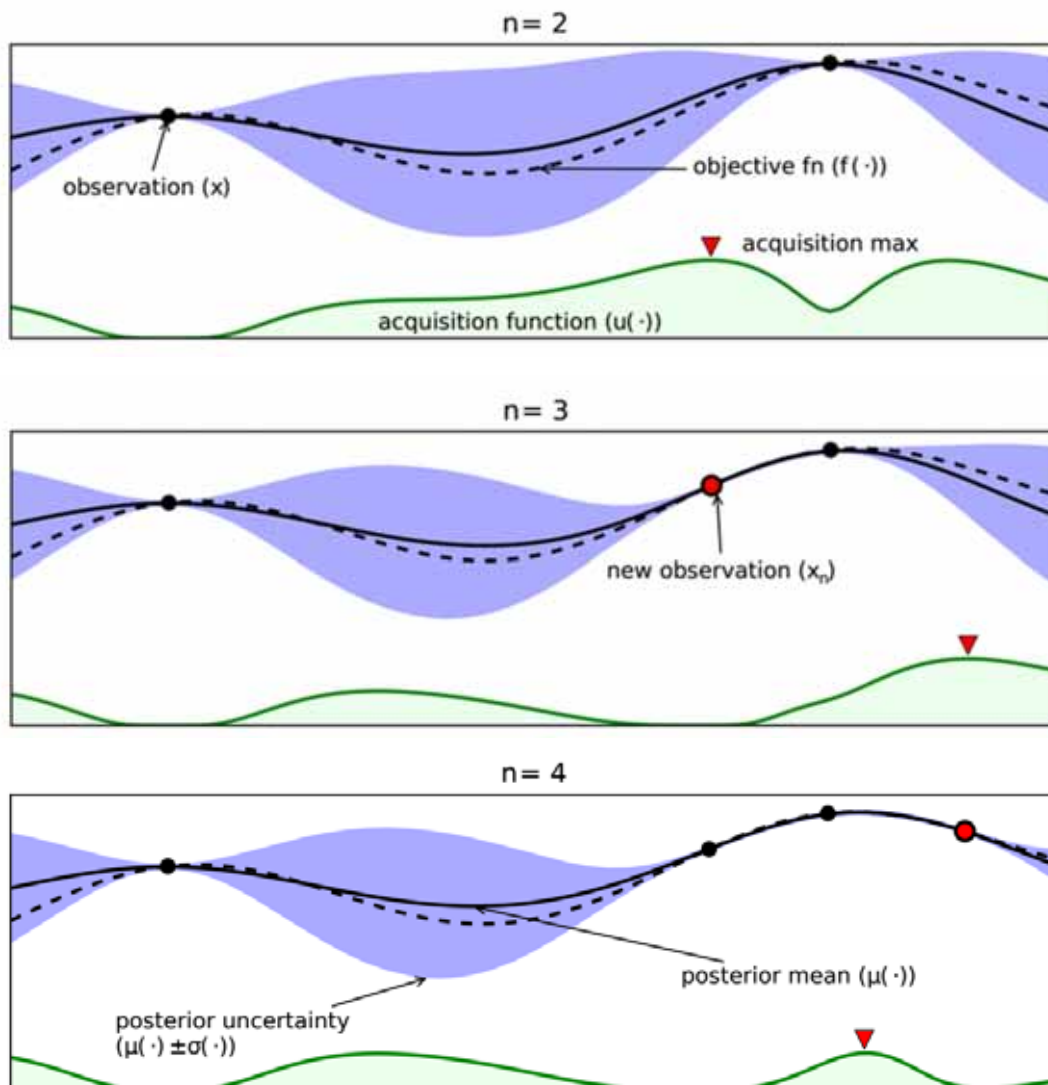
$$\mathbb{E}[f] = \int p(x) f(x) dx$$

$$\mathbb{E}[f] \simeq \frac{1}{N} \sum_{n=1}^N f(x_n)$$

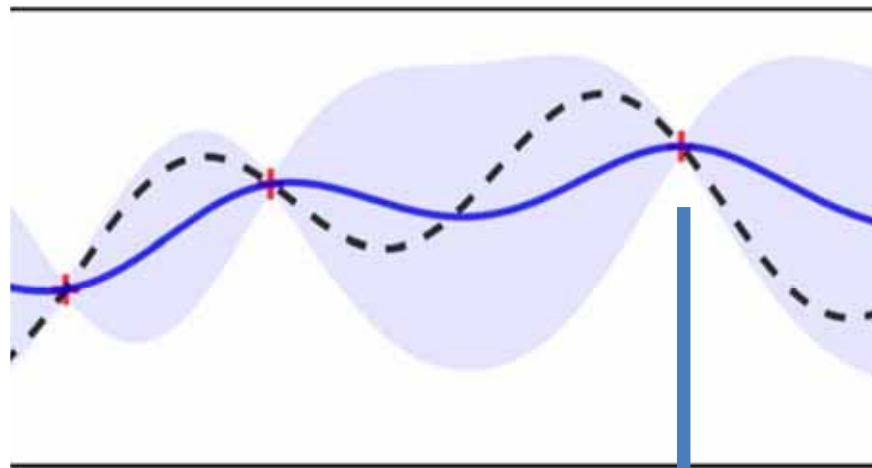
Holzinger, A. 2017. Introduction to Machine Learning and Knowledge Extraction (MAKE). Machine Learning and Knowledge Extraction, 1, (1), 1-20, doi:10.3390/make1010001.



Brochu, E., Cora, V. M. & De Freitas, N. 2010. A tutorial on Bayesian optimization of expensive cost functions, with application to active user modeling and hierarchical reinforcement learning. arXiv:1012.2599.

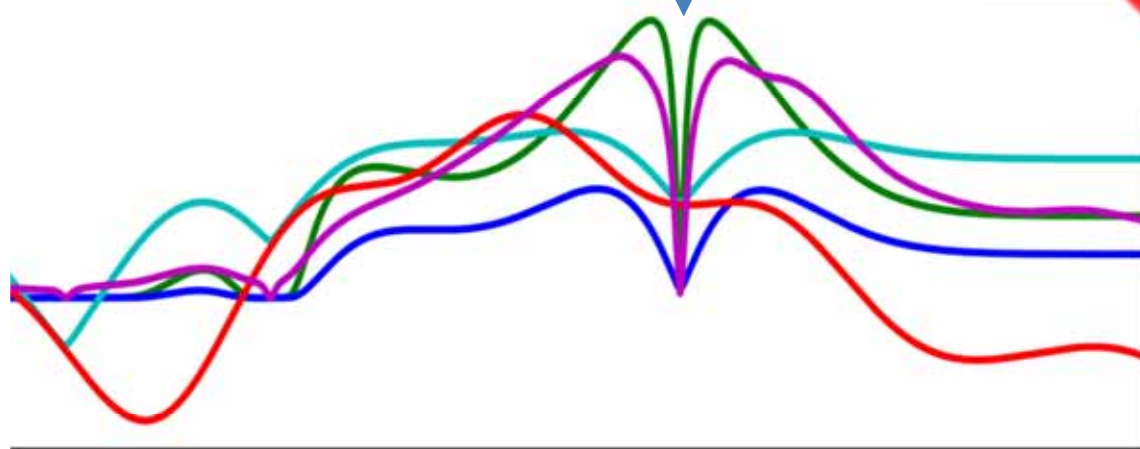
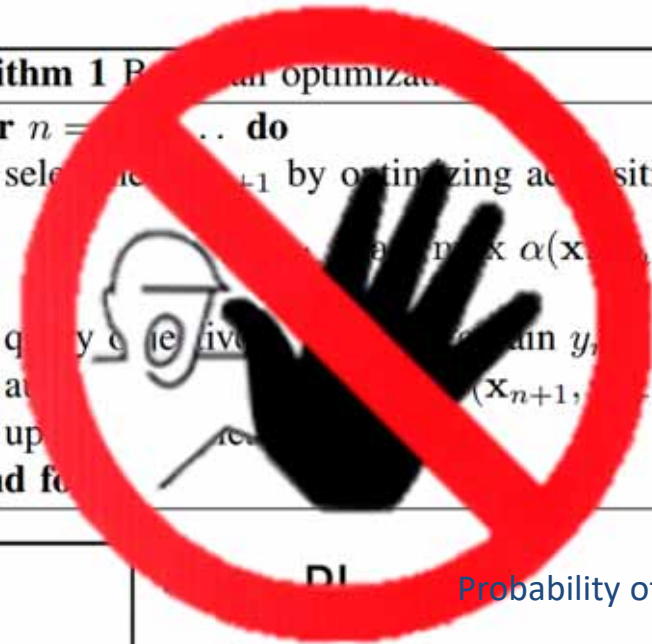


Wang, Z., Hutter, F., Zoghi, M., Matheson, D. & De Feitas, N. 2016. Bayesian optimization in a billion dimensions via random embeddings. Journal of Artificial Intelligence Research, 55, 361-387, doi:10.1613/jair.4806.



```

Algorithm 1 Bayesian optimization
1: for  $n = 1, \dots, N$  do
2:   select  $\mathbf{x}_{n+1}$  by optimizing acquisition function  $\alpha$ 
3:   query objective function to obtain  $y_n = f(\mathbf{x}_n)$ 
4:   add  $(\mathbf{x}_n, y_n)$  to the training set  $\mathcal{D}_n = \{(\mathbf{x}_{n+1}, y_{n+1})\}$ 
5:   update model
6: end for
    
```

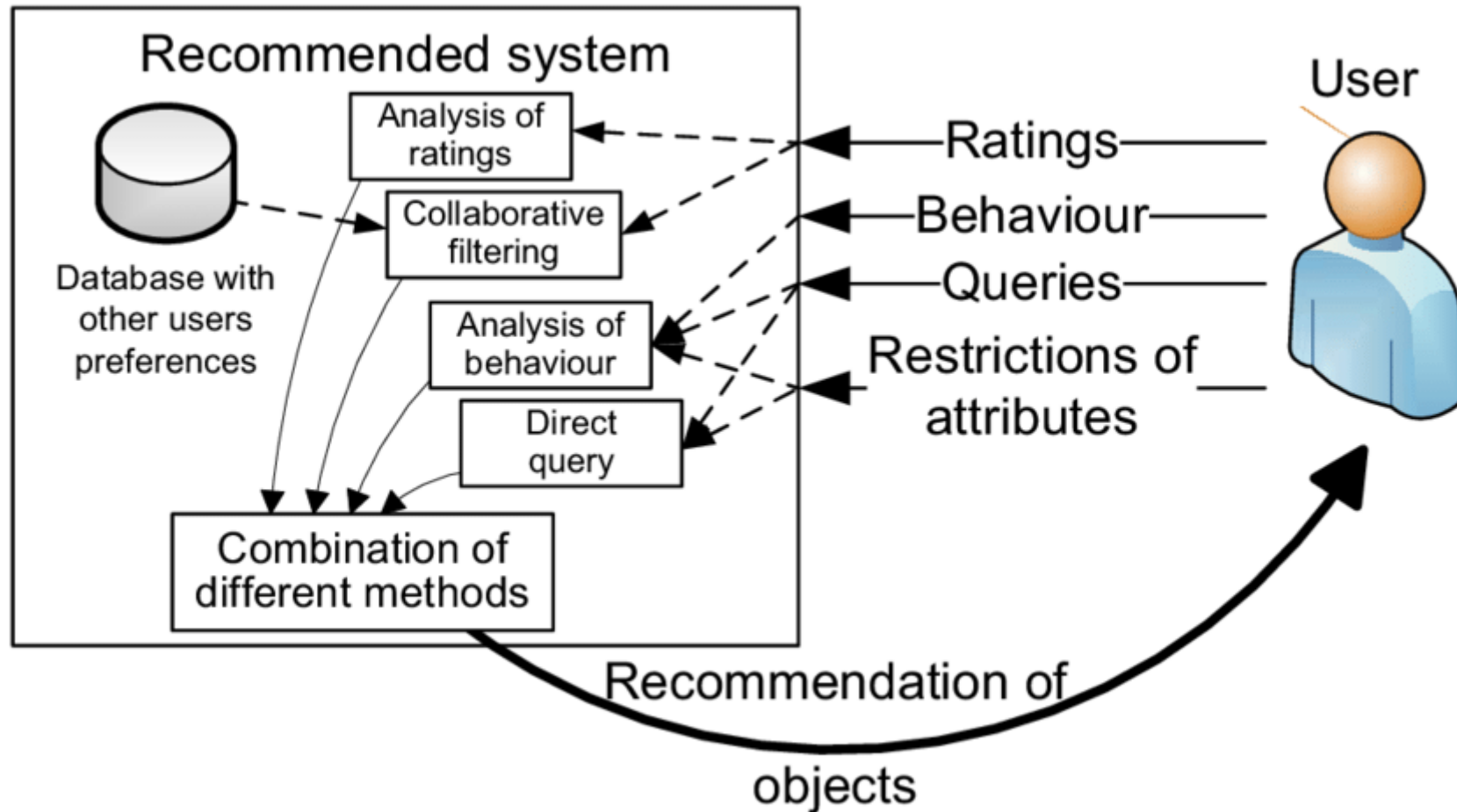


- EI      Expected Improvement
- UCB    Upper Confidence Bound
- TS      Thompson Sampling
- PES     Predictive Entropy Search

Shahriari, B., Swersky, K., Wang, Z., Adams, R. P. & De Freitas, N. 2016.  
**Taking the human out of the loop:** A review of Bayesian optimization.  
*Proceedings of the IEEE*, 104, (1), 148-175, doi:10.1109/JPROC.2015.2494218.

# 04 aML

# Best practice examples of aML ...



Alan Eckhardt 2009. Various aspects of user preference learning and recommender systems. DATESO. pp. 56-67.



| Human                       |                                  |                                   | Machine                 |                        |                         |
|-----------------------------|----------------------------------|-----------------------------------|-------------------------|------------------------|-------------------------|
| LEVEL 0                     | LEVEL 1                          | LEVEL 2                           | LEVEL 3                 | LEVEL 4                | LEVEL 5                 |
| No Active Assistance System | Longitudinal or Transverse Guide | Traffic Control                   | Awareness for Take Over | No Driver Intervention | No Driver               |
|                             | Longitudinal or Transverse Guide | Longitudinal and Transverse Guide | Take Over Request       | No Take Over Request   |                         |
| Hands On<br>Eyes On         | Hands On<br>Eyes On              | Hands Temp Off<br>Eyes Temp Off   | Hands Off<br>Eyes Off   | Hands Off<br>Mind Off  | Hands Off<br>Driver Off |
|                             |                                  |                                   |                         |                        |                         |

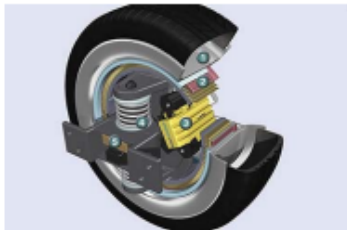




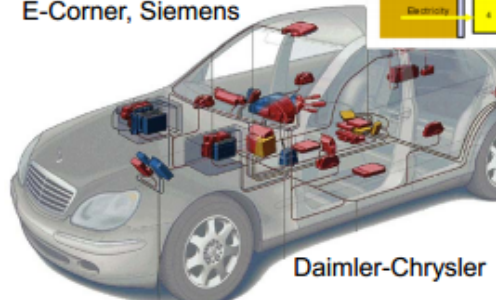
<https://www.businessinsider.sg/the-worlds-first-passenger-drone-makes-public-flight-in-china-and-you-could-soon-own-one>

## Cyber-Physical Systems (CPS): Tight integration of networked computation with physical systems

### Automotive

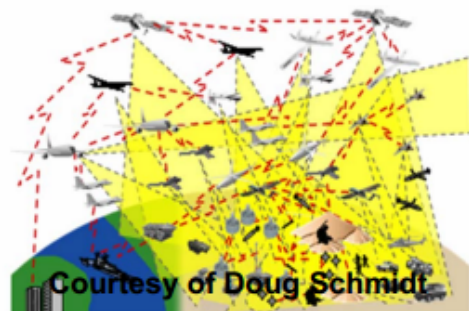


E-Corner, Siemens



Daimler-Chrysler

### Military systems:



Courtesy of Doug Schmidt

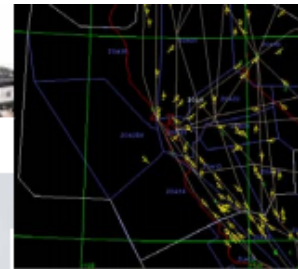
### Building Systems



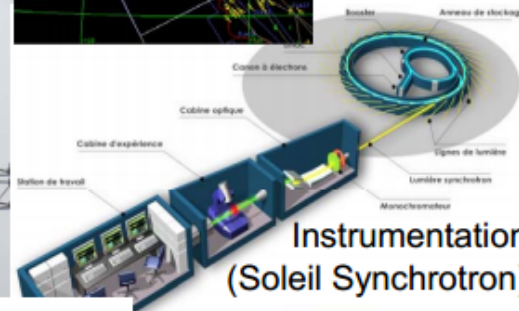
### Telecommunications



### Avionics



Transportation  
(Air traffic control at SFO)



Instrumentation  
(Soleil Synchrotron)

### Power generation and distribution



Courtesy of General Electric

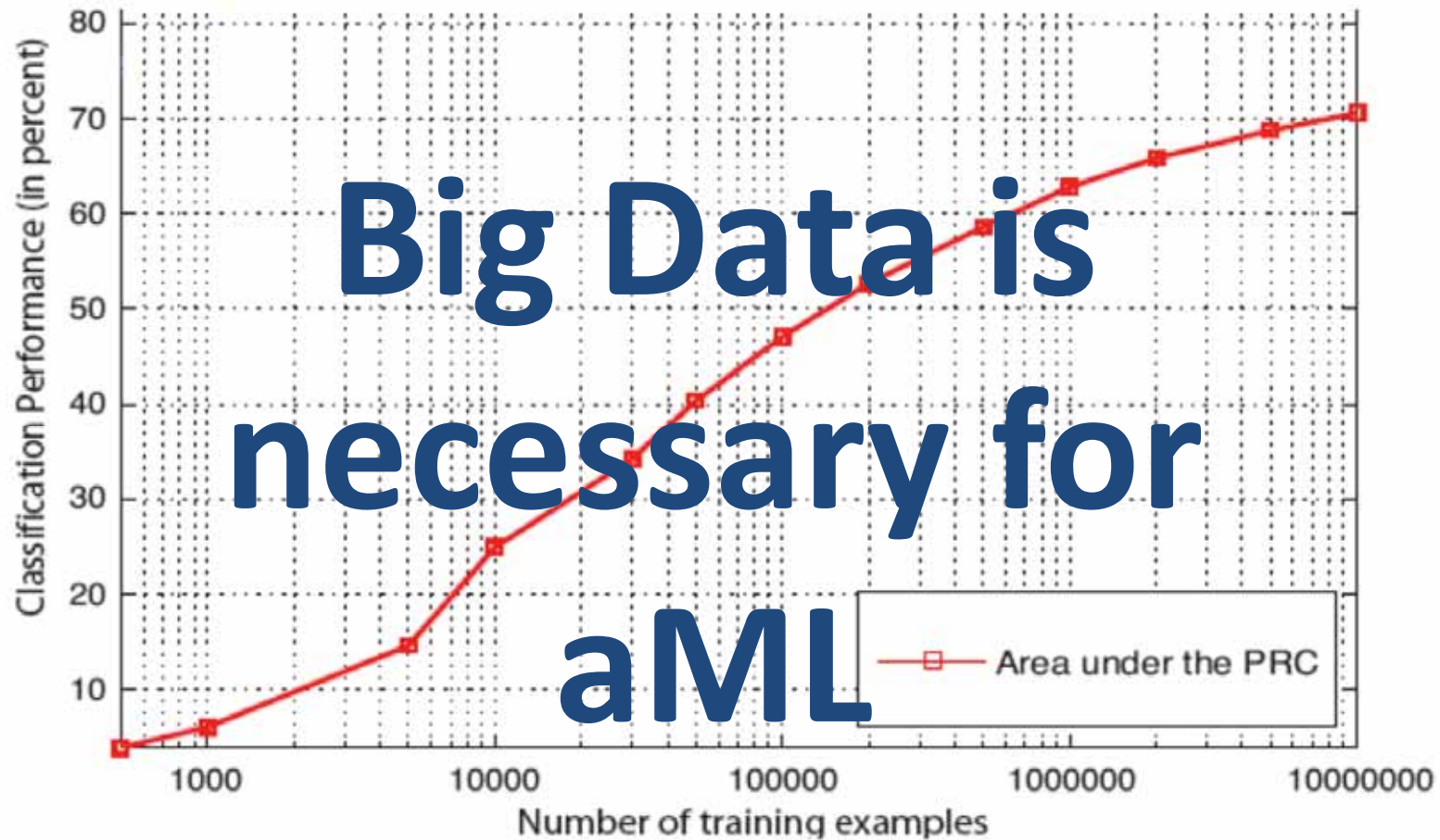
### Factory automation



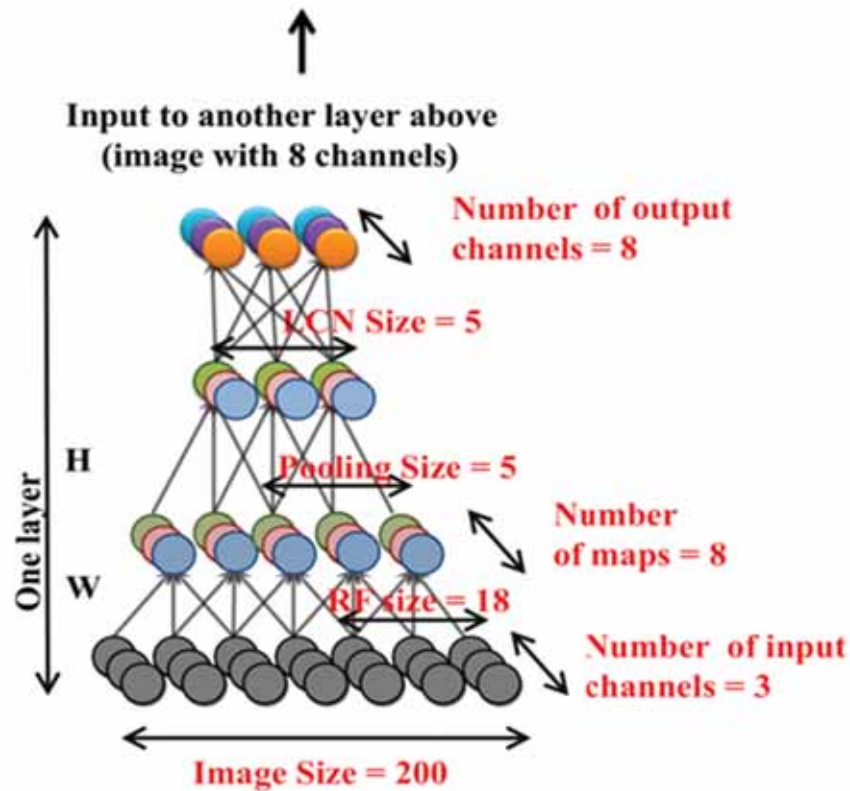
Courtesy of Kuka Robotics Corp.



Seshia, S. A., Juniwal, G., Sadigh, D., Donze, A., Li, W., Jensen, J. C., Jin, X., Deshmukh, J., Lee, E. & Sastry, S. 2015. Verification by, for, and of Humans: Formal Methods for Cyber-Physical Systems and Beyond. Illinois ECE Colloquium.



Sonnenburg, S., Rätsch, G., Schäfer, C. & Schölkopf, B. 2006. Large scale multiple kernel learning. Journal of Machine Learning Research, 7, (7), 1531-1565.



$$x^* = \arg \min_x f(x; W, H), \text{ subject to } \|x\|_2 = 1.$$

Le, Q. V., Ranzato, M. A., Monga, R., Devin, M., Chen, K., Corrado, G. S., Dean, J. & Ng, A. Y. 2011. Building high-level features using large scale unsupervised learning. arXiv preprint arXiv:1112.6209.

Le, Q. V. 2013. Building high-level features using large scale unsupervised learning. *IEEE Intl. Conference on Acoustics, Speech and Signal Processing ICASSP*. IEEE. 8595-8598, doi:10.1109/ICASSP.2013.6639343.

- Sometimes we **do not have “big data”**, where aML-algorithms benefit.
- Sometimes we have
  - **Small amount of data sets**
  - **Rare Events – no training samples**
  - **NP-hard problems, e.g.**
    - Subspace Clustering,
    - k-Anonymization,
    - Protein-Folding, ...

**sometimes  
we need a  
human-in-the-loop  
Because humans are very good in  
learning abstract concepts!!  
(see #KandinskyPatterns)**

# 05 iML

Holzinger, A. 2016. Interactive Machine Learning (iML). Informatik Spektrum, 39, (1), 64-68, doi:10.1007/s00287-015-0941-6.

**sometimes  
we need a  
human-in-the-loop  
Because humans are very good in  
learning abstract concepts!!  
(see #KandinskyPatterns)**

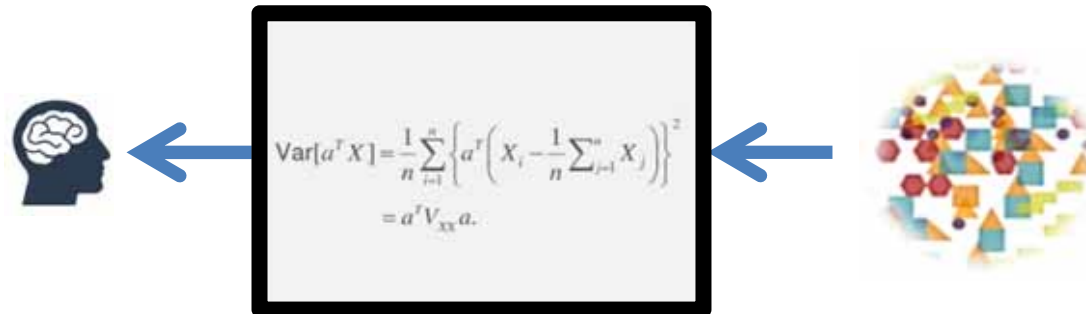


- **Example 1: Subspace Clustering**
- **Example 2: k-Anonymization**
- **Example 3: Protein Design**

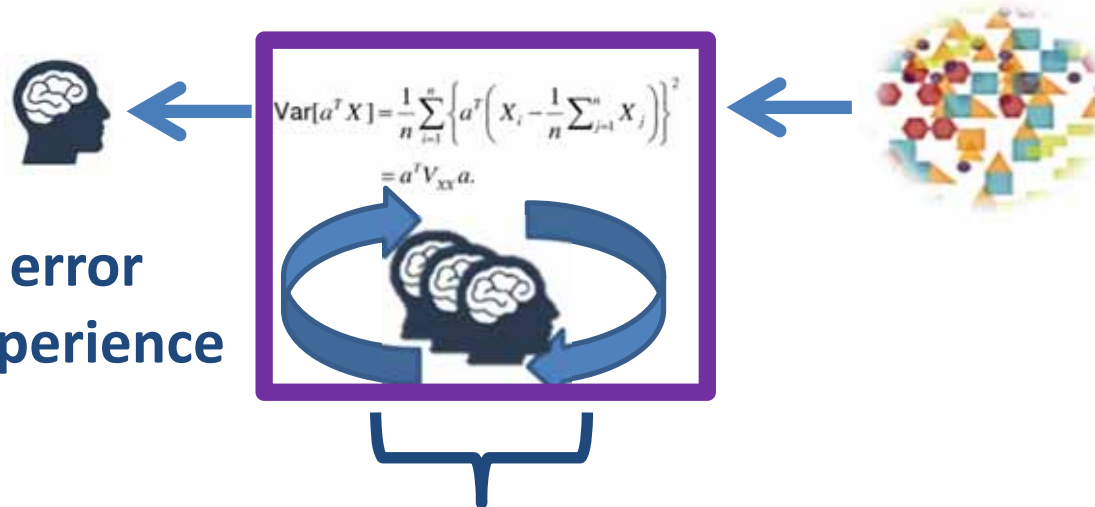
Hund, M., Böhm, D., Sturm, W., Sedlmair, M., Schreck, T., Ullrich, T., Keim, D. A., Majnaric, L. & Holzinger, A. 2016. Visual analytics for concept exploration in subspaces of patient groups: Making sense of complex datasets with the Doctor-in-the-loop. *Brain Informatics*, 1-15, doi:10.1007/s40708-016-0043-5.

Kieseberg, P., Malle, B., Fruehwirt, P., Weippl, E. & Holzinger, A. 2016. A tamper-proof audit and control system for the doctor in the loop. *Brain Informatics*, 3, (4), 269–279, doi:10.1007/s40708-016-0046-2.

Lee, S. & Holzinger, A. 2016. Knowledge Discovery from Complex High Dimensional Data. In: Michaelis, S., Piatkowski, N. & Stolpe, M. (eds.) *Solving Large Scale Learning Tasks. Challenges and Algorithms*, Lecture Notes in Artificial Intelligence LNAI 9580. Springer, pp. 148-167, doi:10.1007/978-3-319-41706-6\_7.



Generalization error



Generalization error  
plus human experience

**iML = human inspection – bring in human intuition**

Andreas Holzinger et al. 2018. Interactive machine learning: experimental evidence for the human in the algorithmic loop. Springer/Nature Applied Intelligence, doi:10.1007/s10489-018-1361-5.

# Why using human concept learning?

- “How do humans generalize from few examples?”
  - Learning relevant representations
  - Disentangling the explanatory factors
  - Finding the shared underlying explanatory factors, in particular between  $P(x)$  and  $P(Y|X)$ , with a causal link between  $Y \rightarrow X$

Bengio, Y., Courville, A. & Vincent, P. 2013. Representation learning: A review and new perspectives. IEEE transactions on pattern analysis and machine intelligence, 35, (8), 1798-1828, doi:10.1109/TPAMI.2013.50.

Tenenbaum, J. B., Kemp, C., Griffiths, T. L. & Goodman, N. D. 2011. How to grow a mind: Statistics, structure, and abstraction. Science, 331, (6022), 1279-1285, doi:10.1126/science.1192788.

# even Children can make inferences from little, noisy, incomplete data ...



This image is in the public domain, Source: freedesignfile.com

Brenden M. Lake, Ruslan Salakhutdinov & Joshua B. Tenenbaum 2015. Human-level concept learning through probabilistic program induction. *Science*, 350, (6266), 1332-1338, doi:[10.1126/science.aab3050](https://doi.org/10.1126/science.aab3050)

# Adversarial Examples that Fool both Computer Vision and Time-Limited Humans

**Gamaleldin F. Elsayed\***  
 Google Brain  
 gamaleldin.elsayed@gmail.com

**Shreya Shankar**  
 Stanford University

**Brian Cheung**  
 UC Berkeley

**Nicolas Papernot**  
 Pennsylvania State University

**Alex Kurakin**  
 Google Brain

**Ian Goodfellow**  
 Google Brain

**Jascha Sohl-Dickstein**  
 Google Brain  
 jaschasd@google.com

## Abstract

Machine learning models are vulnerable to **adversarial examples**: small changes to images can cause computer vision models to make mistakes such as identifying a school bus as an ostrich. However, it is still an open question whether humans are prone to similar mistakes. Here, we address this question by leveraging recent techniques that transfer adversarial examples from computer vision models with known parameters and architecture to other models with unknown parameters and architecture, and by matching the initial processing of the human visual system. We find that adversarial examples that strongly transfer across computer vision models influence the classifications made by time-limited human observers.

Gamaleldin F Elsayed, Shreya Shankar, Brian Cheung, Nicolas Papernot, Alex Kurakin, Ian Goodfellow & Jascha Sohl-Dickstein 2018. Adversarial Examples that Fool both Human and Computer Vision. arXiv:1802.08195.

v3 [cs.LG] 22 May 2018



a woman riding a horse on a dirt road



an airplane is parked on the tarmac at an airport



a group of people standing on top of a beach

Andrej Karpathy & Li Fei-Fei. Deep visual-semantic alignments for generating image descriptions. Proceedings of the IEEE conference on computer vision and pattern recognition, 2015. 3128-3137.

Image Captions by dee learning : [github.com/karpathy/neuraltalk2](https://github.com/karpathy/neuraltalk2)

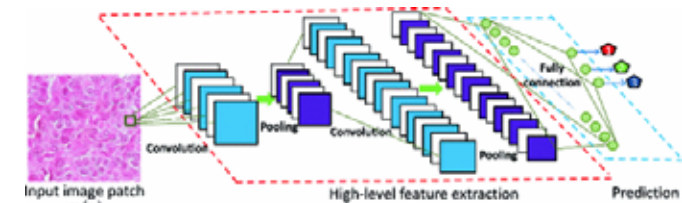
Image Source: Gabriel Villena Fernandez; Agence France-Press, Dave Martin (left to right)

# 06 Methods of Explainable AI



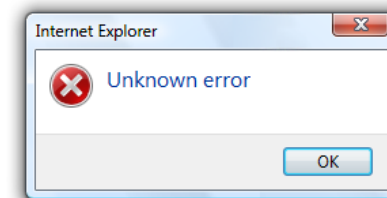
## Verify that algorithms/classifiers work as expected ...

Wrong decisions can be costly and dangerous ...



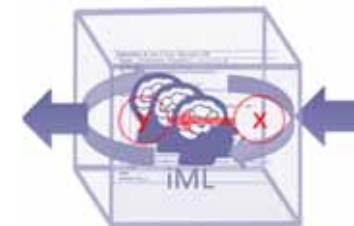
## Understanding the errors ...

Detection of bias, weaknesses, unknowns, ...



## Scientific replicability and causality ...

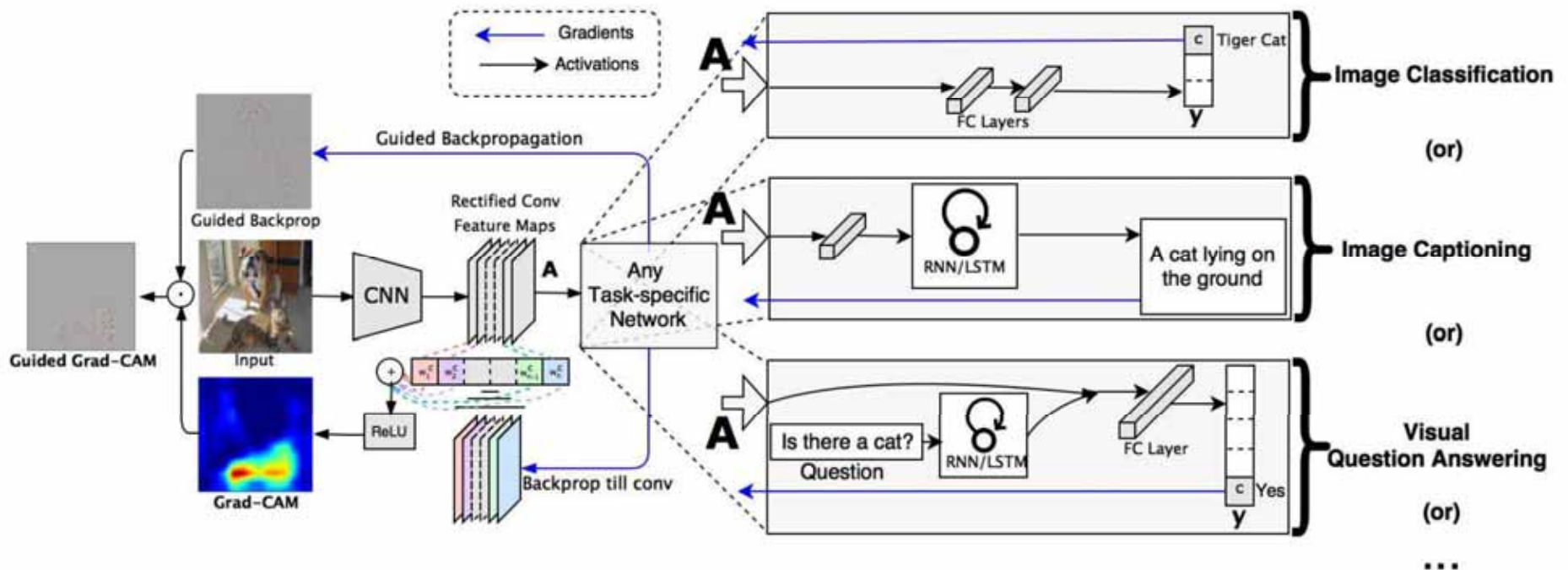
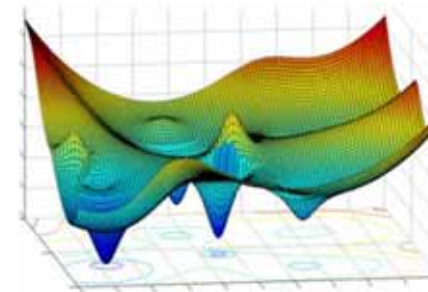
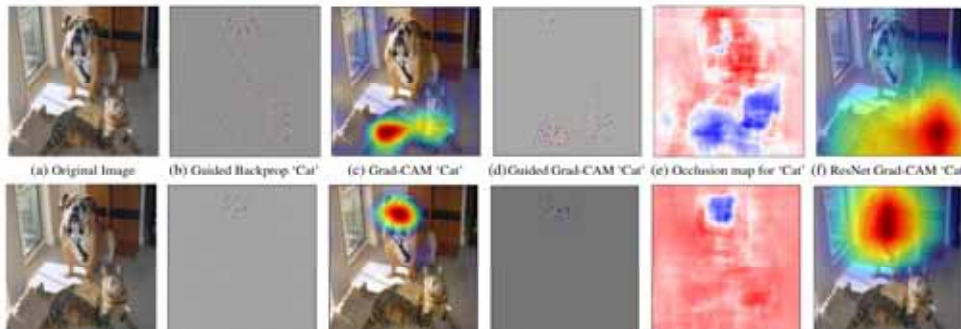
The “why” is often more important than the prediction ...



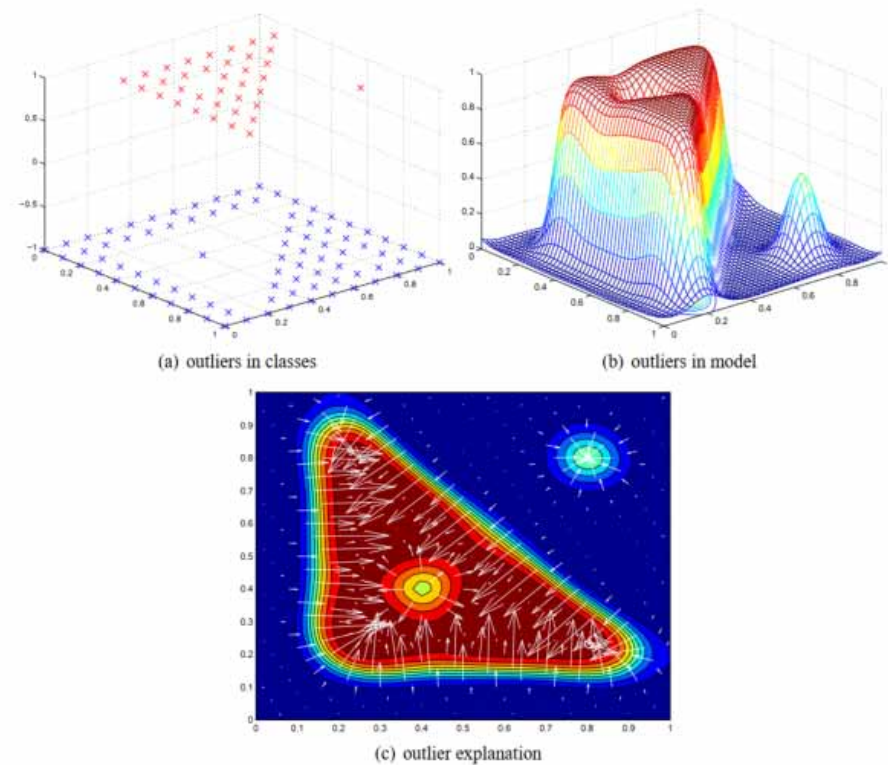
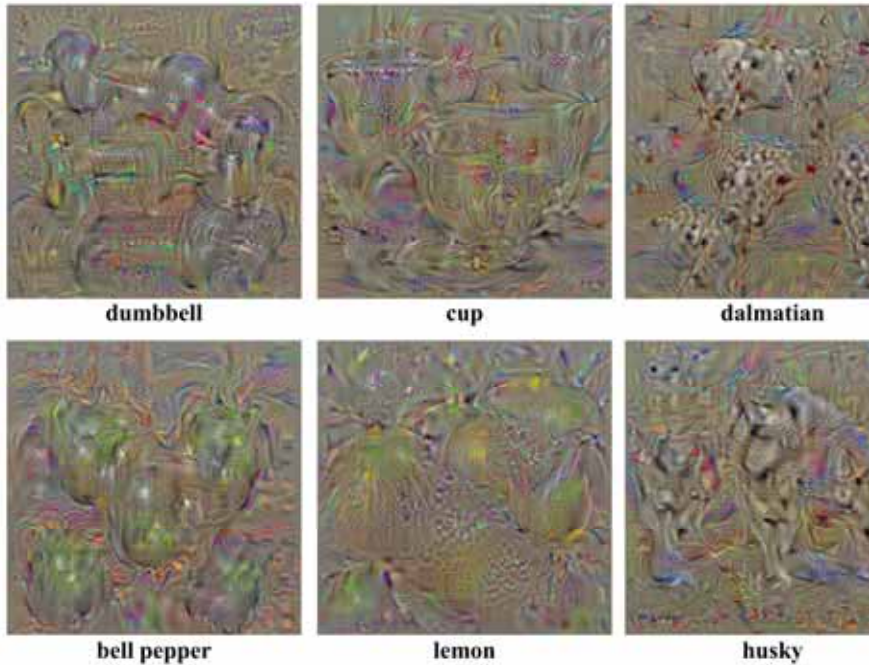
- 1) Gradients
- 2) Sensitivity Analysis
- 3) Decomposition Relevance Propagation  
(Pixel-RP, Layer-RP, Deep Taylor Decomposition, ...)
- 4) Optimization (Local-IME – model agnostic,  
BETA transparent approximation, ...)
- 5) Deconvolution and Guided Backpropagation
- 6) Model Understanding
  - Feature visualization, Inverting CNN
  - Qualitative Testing with Concept Activation Vectors TCAV
  - Network Dissection

Andreas Holzinger LV 706.315 From explainable AI to Causability, 3 ECTS course at Graz University of Technology

<https://human-centered.ai/explainable-ai-causability-2019> (course given since 2016)



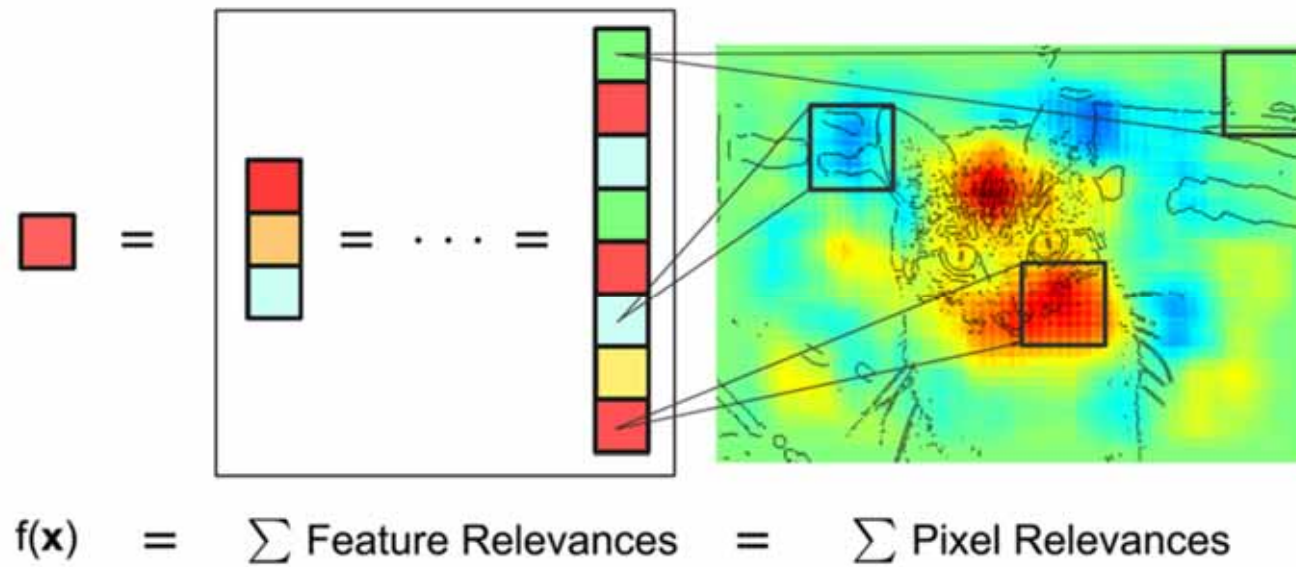
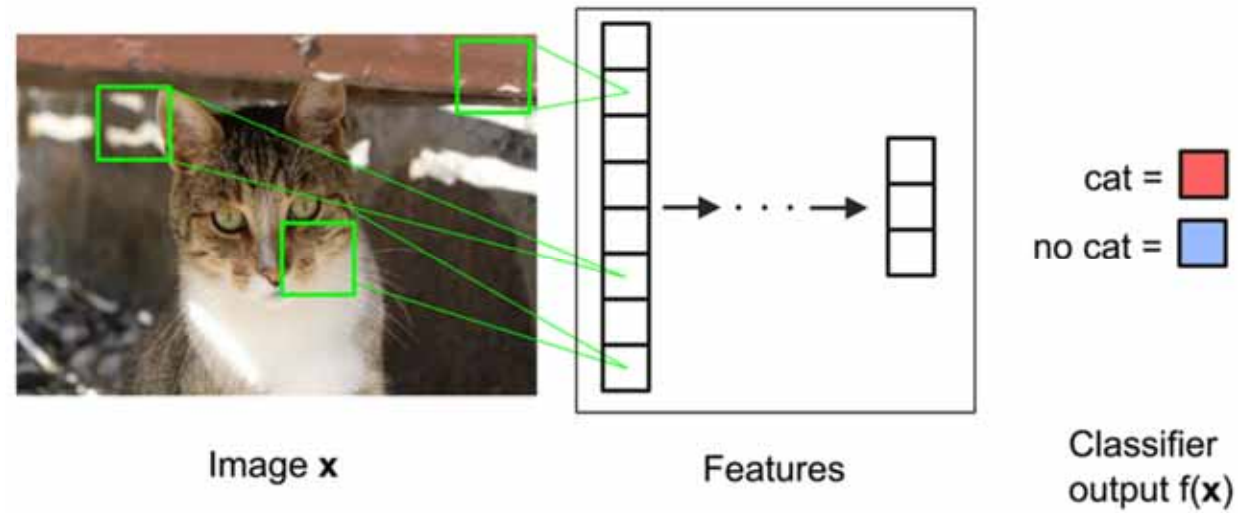
Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh & Dhruv Batra. Grad-CAM: Visual Explanations from Deep Networks via Gradient-Based Localization. ICCV, 2017. 618-626.

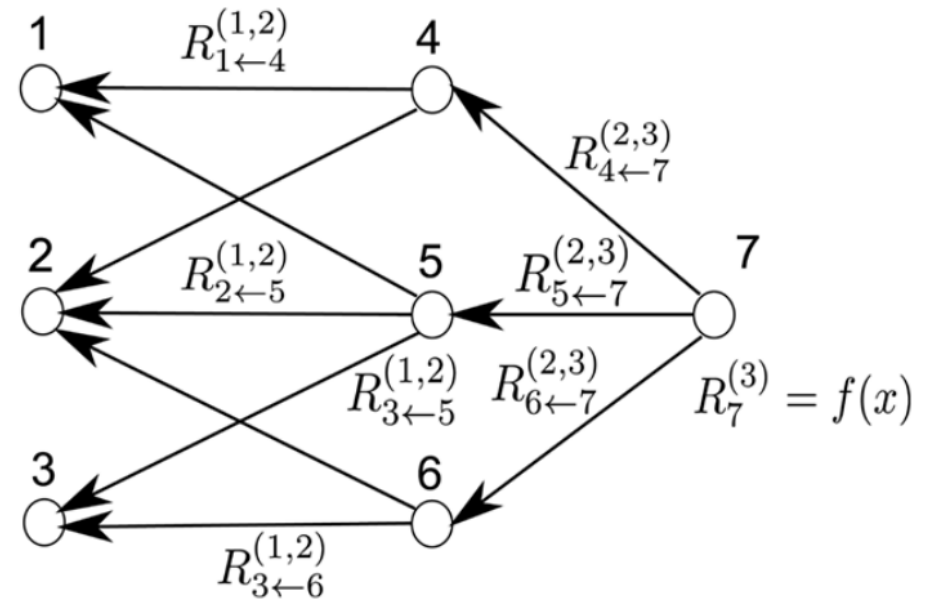
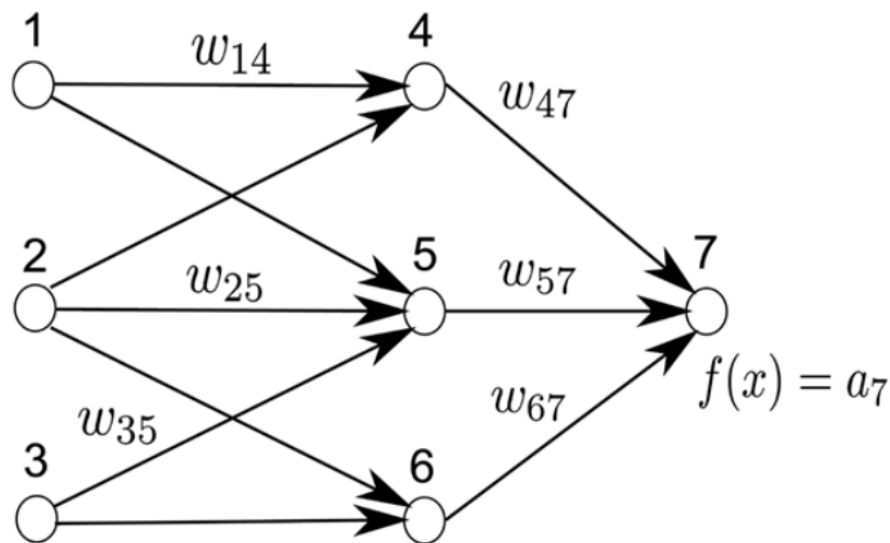


Karen Simonyan, Andrea Vedaldi & Andrew Zisserman 2013. Deep inside convolutional networks: Visualising image classification models and saliency maps. arXiv:1312.6034.

David Baehrens, Timon Schroeter, Stefan Harmeling, Motoaki Kawanabe, Katja Hansen & Klaus-Robert Mueller 2010. How to explain individual classification decisions. Journal of machine learning research (JMLR), 11, (6), 1803-1831.

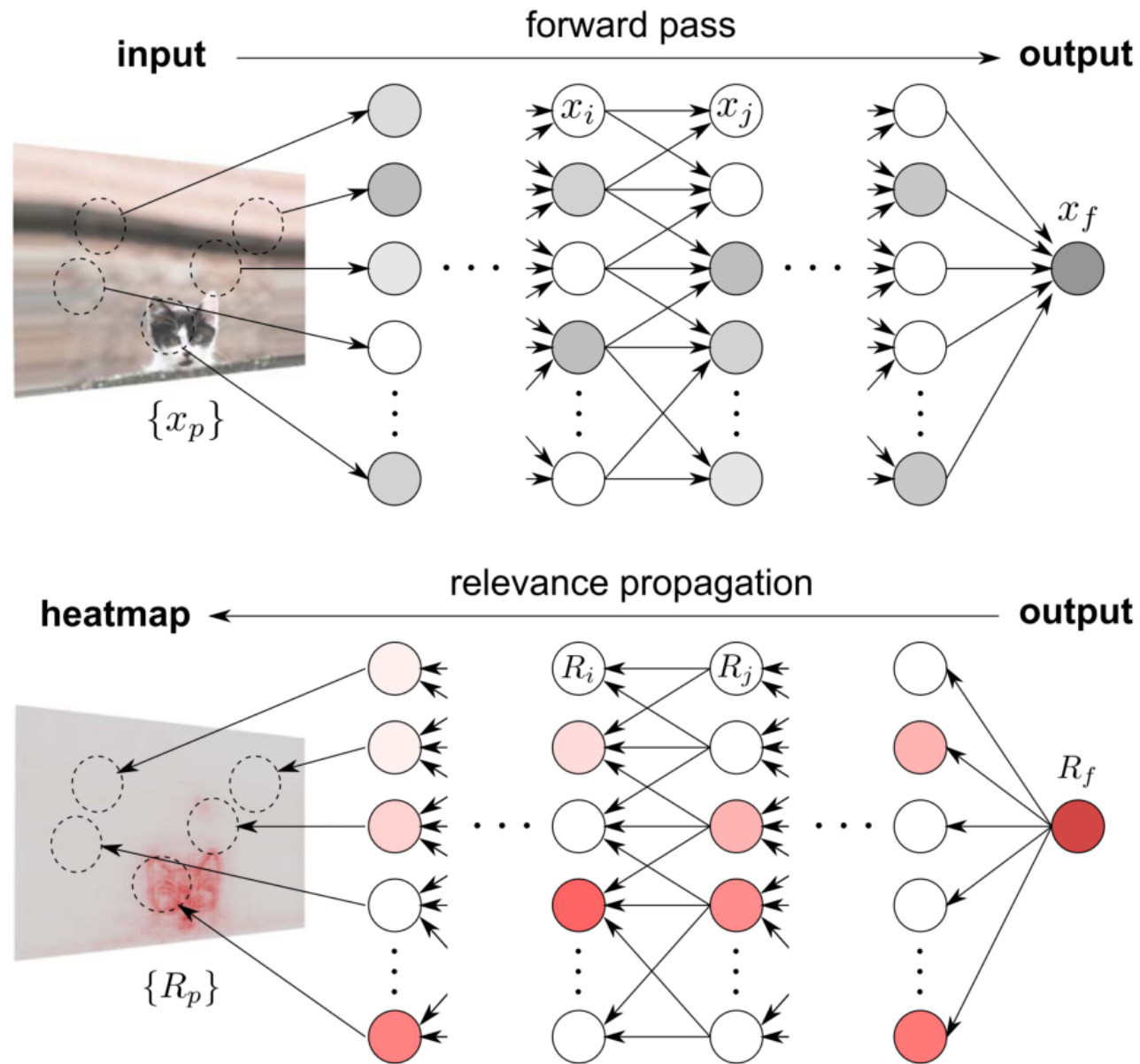
Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller & Wojciech Samek 2015. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. PLoS one, 10, (7), e0130140, doi:10.1371/journal.pone.0130140.

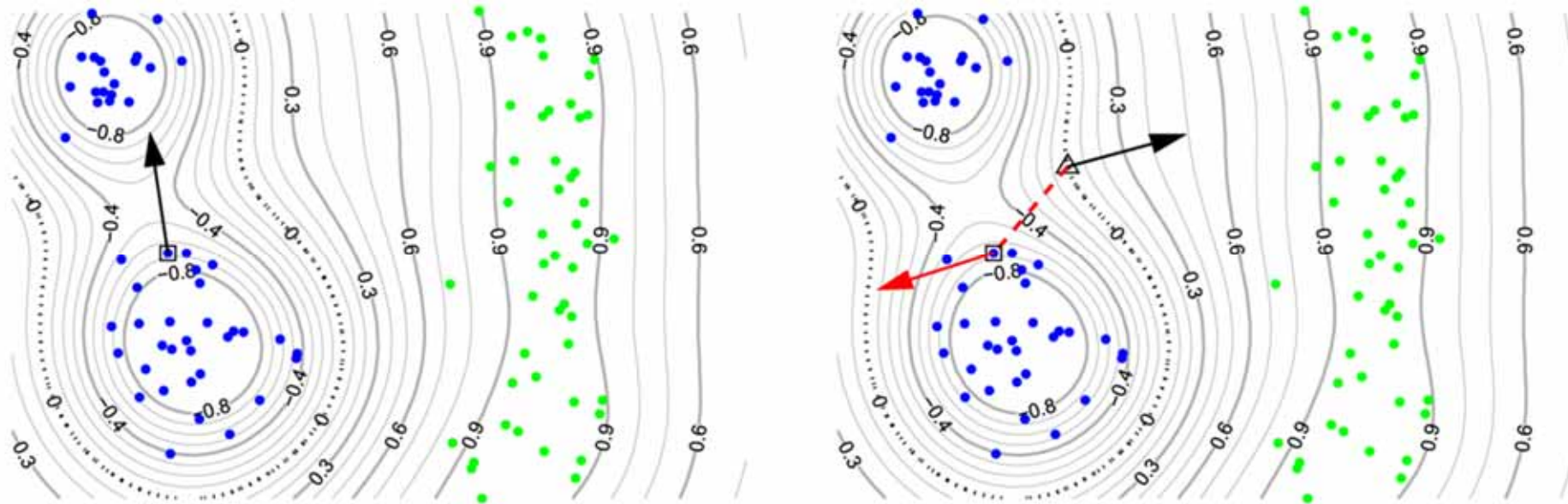




$$f(x) = \dots = \sum_{d \in l+1} R_d^{(l+1)} = \sum_{d \in l} R_d^{(l)} = \dots = \sum_d R_d^{(1)}$$

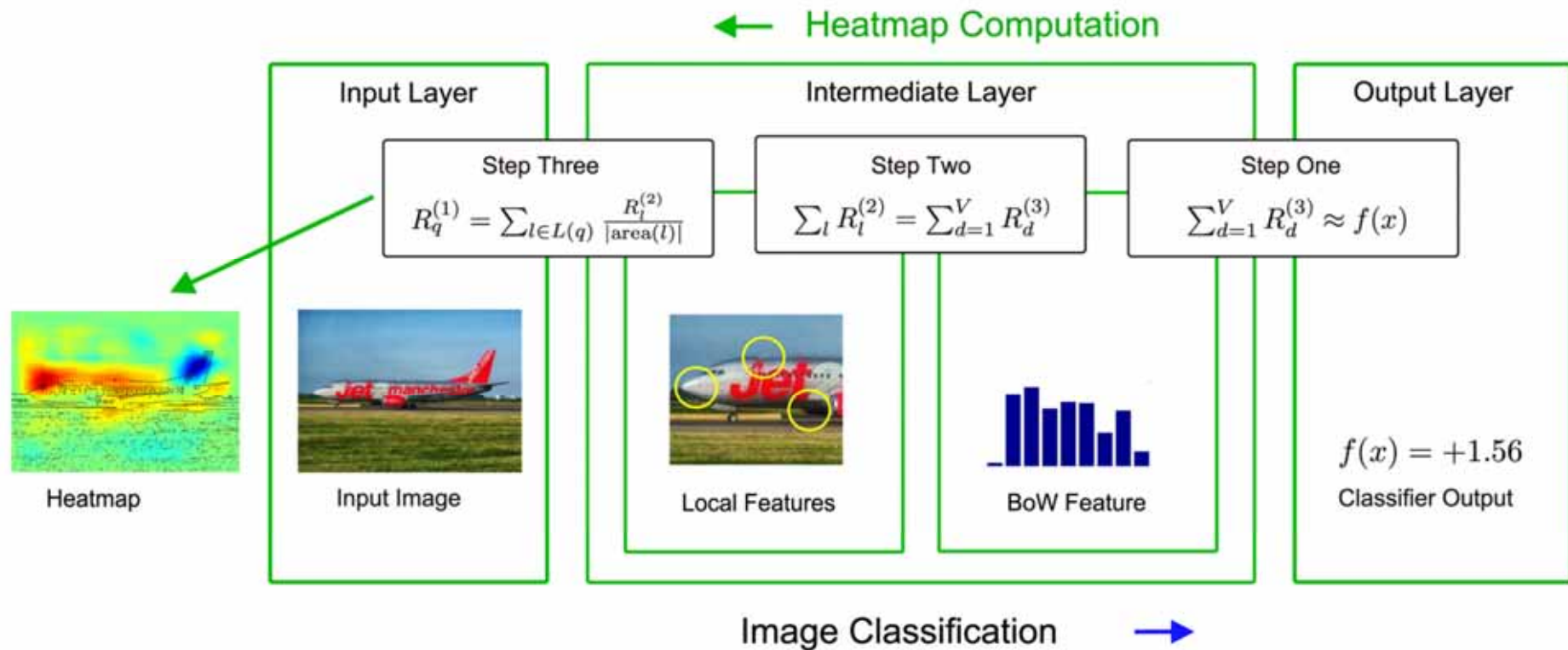
Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller & Wojciech Samek 2015. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. PLoS one, 10, (7), e0130140, doi:10.1371/journal.pone.0130140.



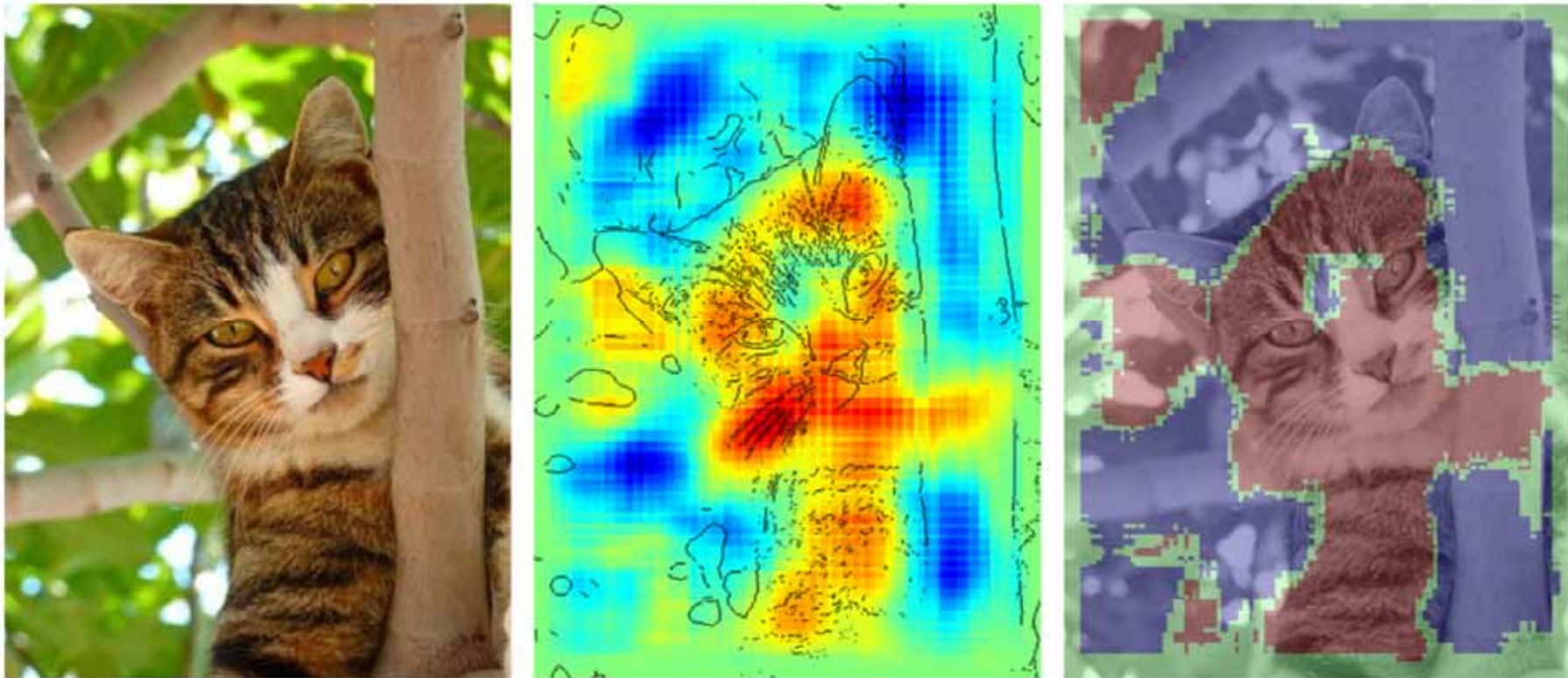


Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller & Wojciech Samek 2015. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. PloS one, 10, (7), e0130140, doi:10.1371/journal.pone.0130140.



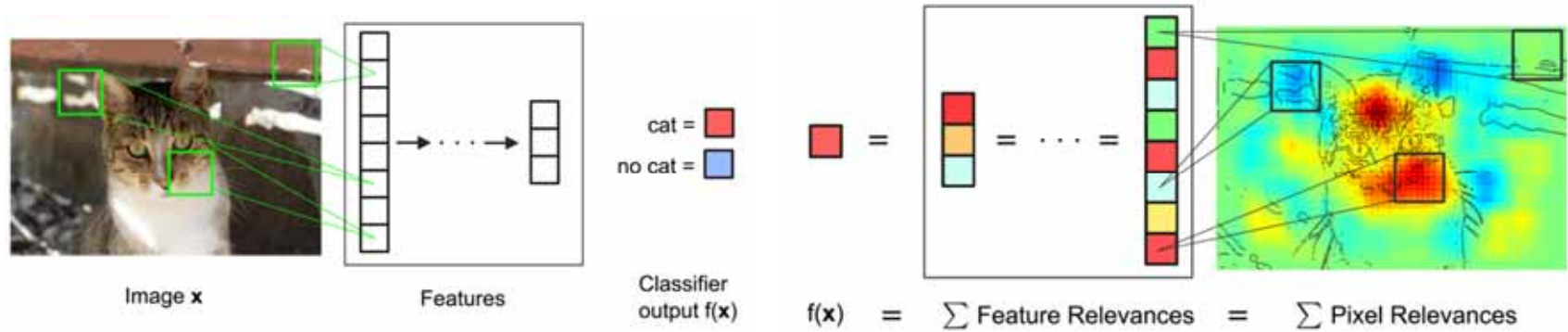


Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller & Wojciech Samek  
 2015. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. PLoS one, 10,  
 (7), e0130140, doi:10.1371/journal.pone.0130140.



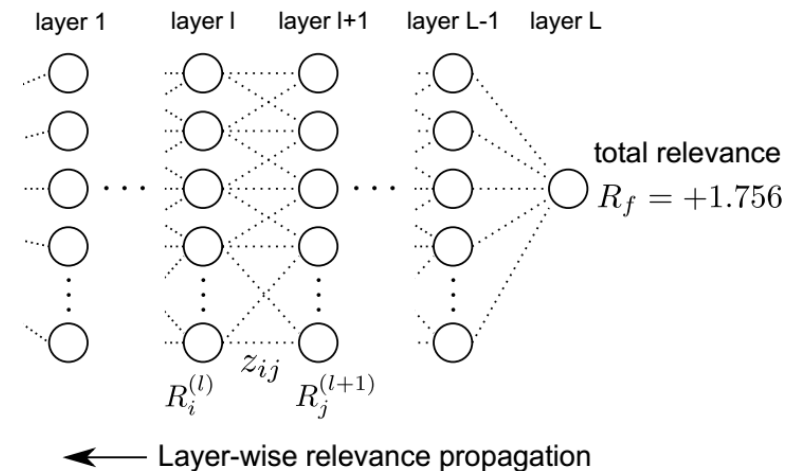
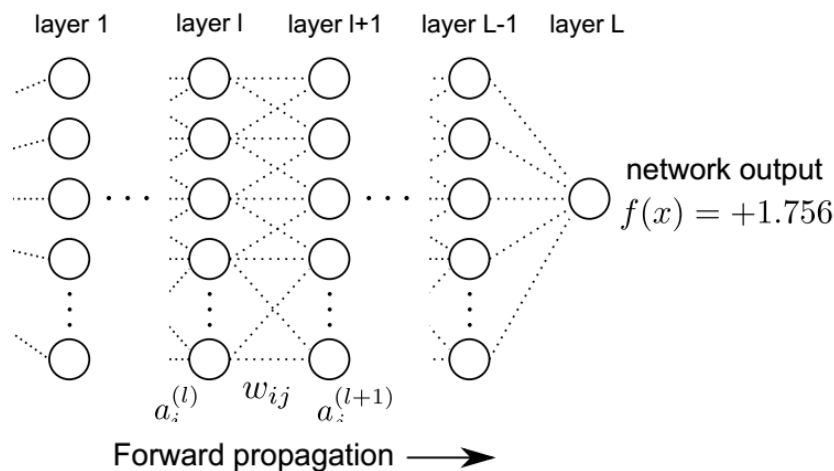
Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller & Wojciech Samek  
2015. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. PloS one, 10,  
(7), e0130140, doi:10.1371/journal.pone.0130140.

Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller & Wojciech Samek 2015. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. PloS one, 10, (7), e0130140, doi:10.1371/journal.pone.0130140.

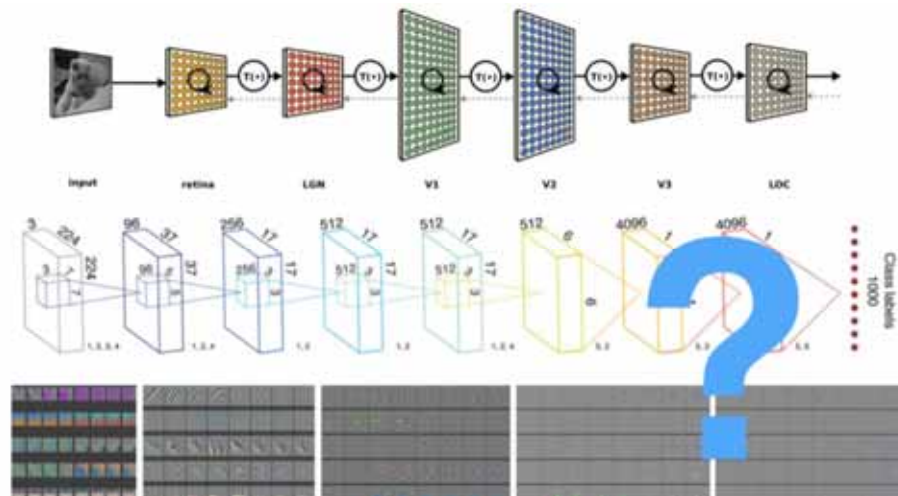


$$a_j^{(l+1)} = \sigma \left( \sum_i a_i^{(l)} w_{ij} + b_j^{(l+1)} \right)$$

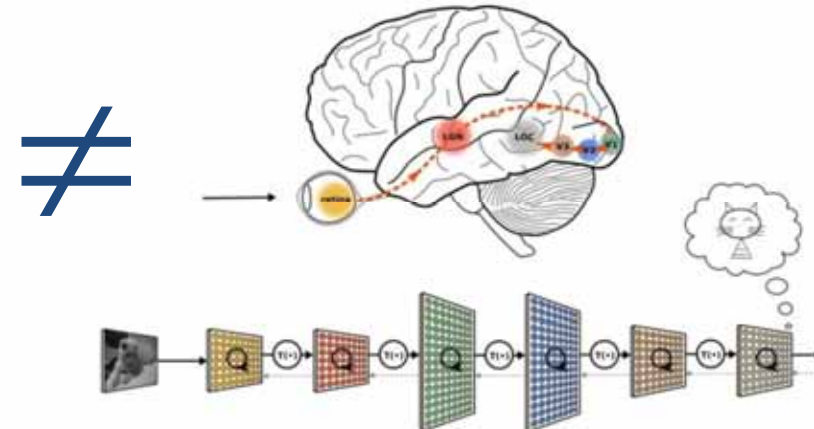
$$R_i^{(l)} = \sum_j \frac{z_{ij}}{\sum_{i'} z_{i'j}} R_j^{(l+1)}$$



$$R_i = \left\| \frac{\partial}{\partial x_i} f(\mathbf{x}) \right\| \quad \sum_i R_i = \dots = \sum_j R_j = \sum_k R_k = \dots = f(\mathbf{x})$$



Yann Lecun, Yoshua Bengio & Geoffrey Hinton 2015. Deep learning. Nature, 521, (7553), 436-444, doi:10.1038/nature14539.

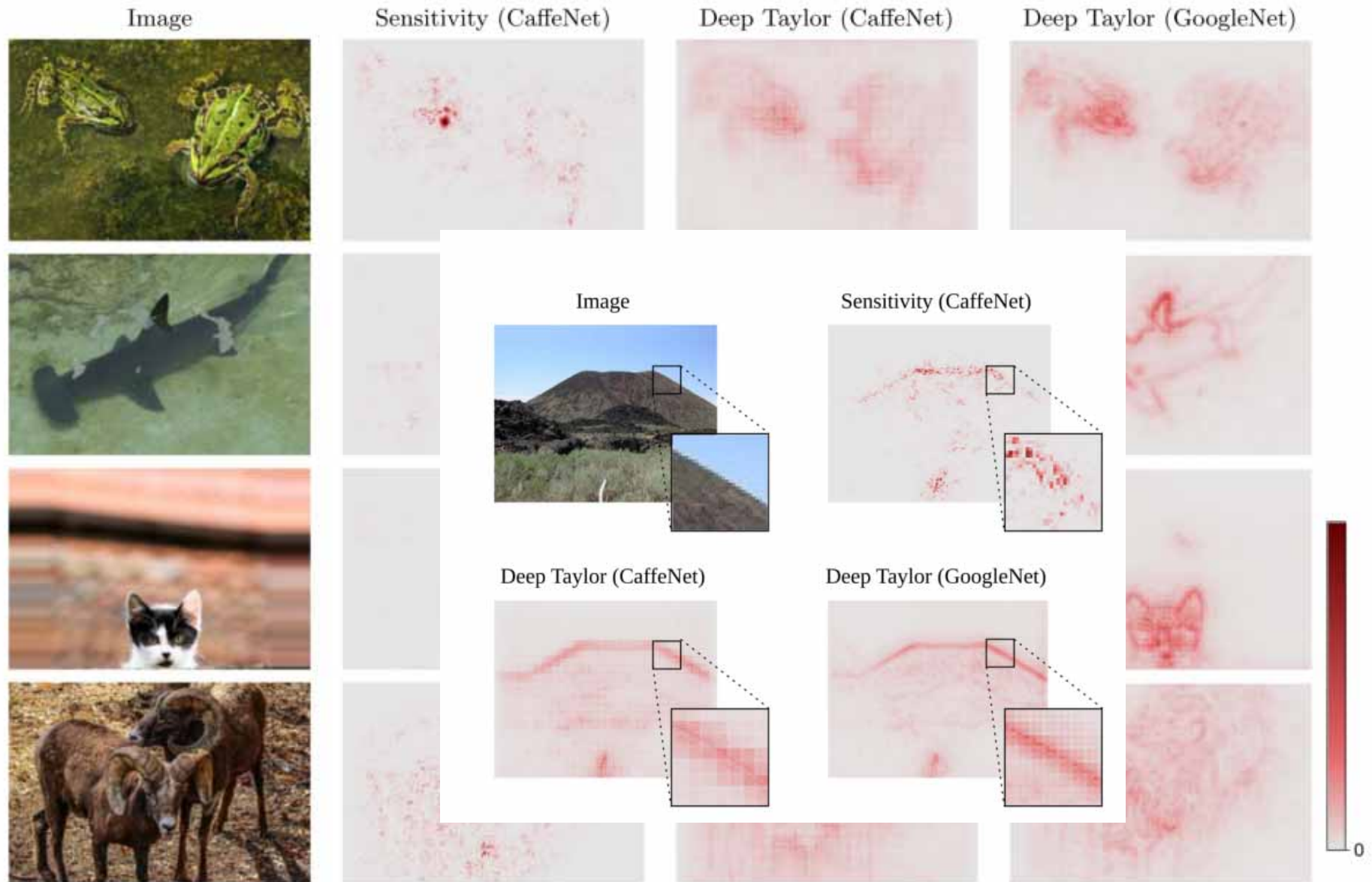


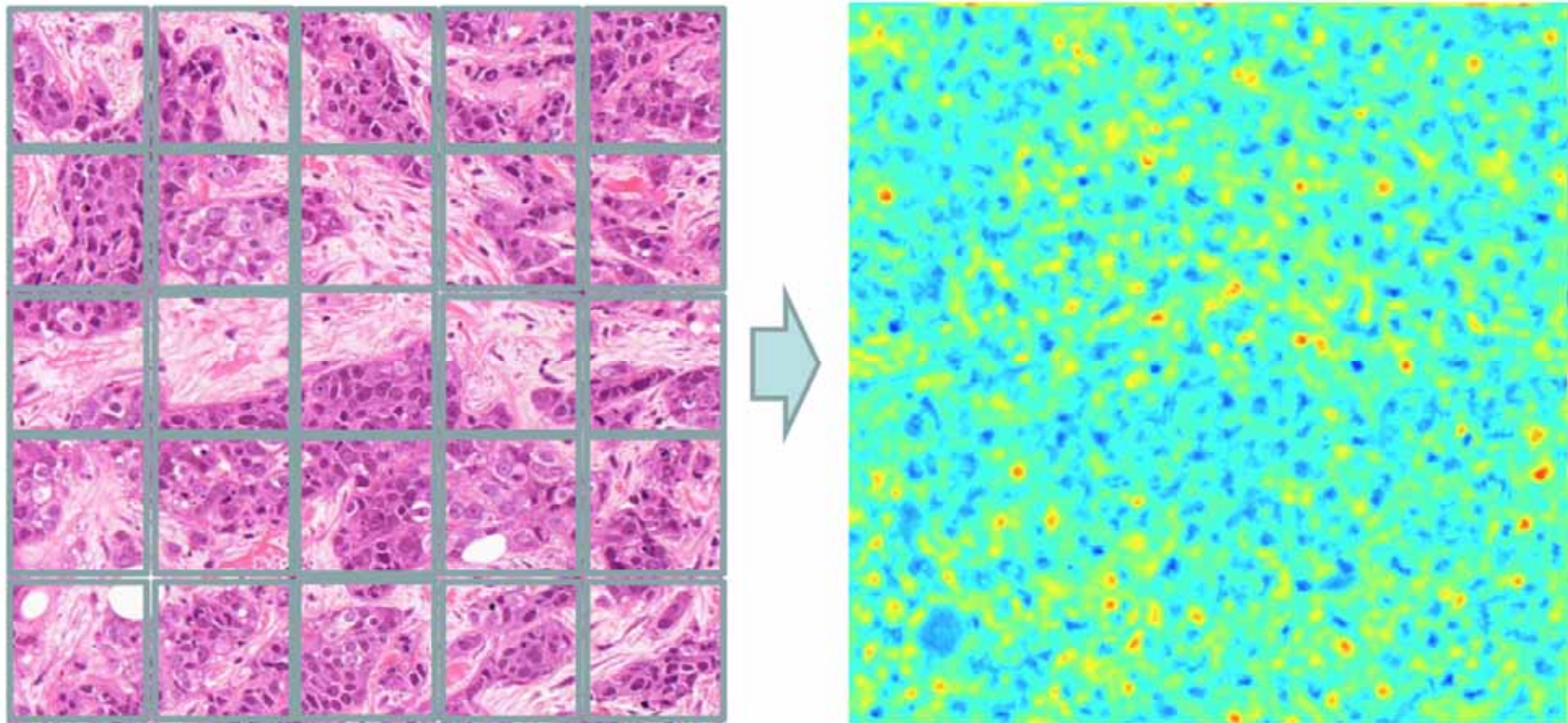
$$\frac{\partial h_k(x)}{\partial x_{a,b}}$$

Humans work in another vector space which is spanned by **implicit knowledge** vectors corresponding to an unknown set of human interpretable concepts.

$$S_{C,k,l}(x) = \lim_{\epsilon \rightarrow 0} \frac{h_{l,k}(f_l(x) + \epsilon v_C^l) - h_{l,k}(f_l(x))}{\epsilon} = \nabla h_{l,k}(f_l(x)) \cdot v_C^l$$

Been Kim, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler & Fernanda Viegas. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors, ICML, 2018. 2673-2682.



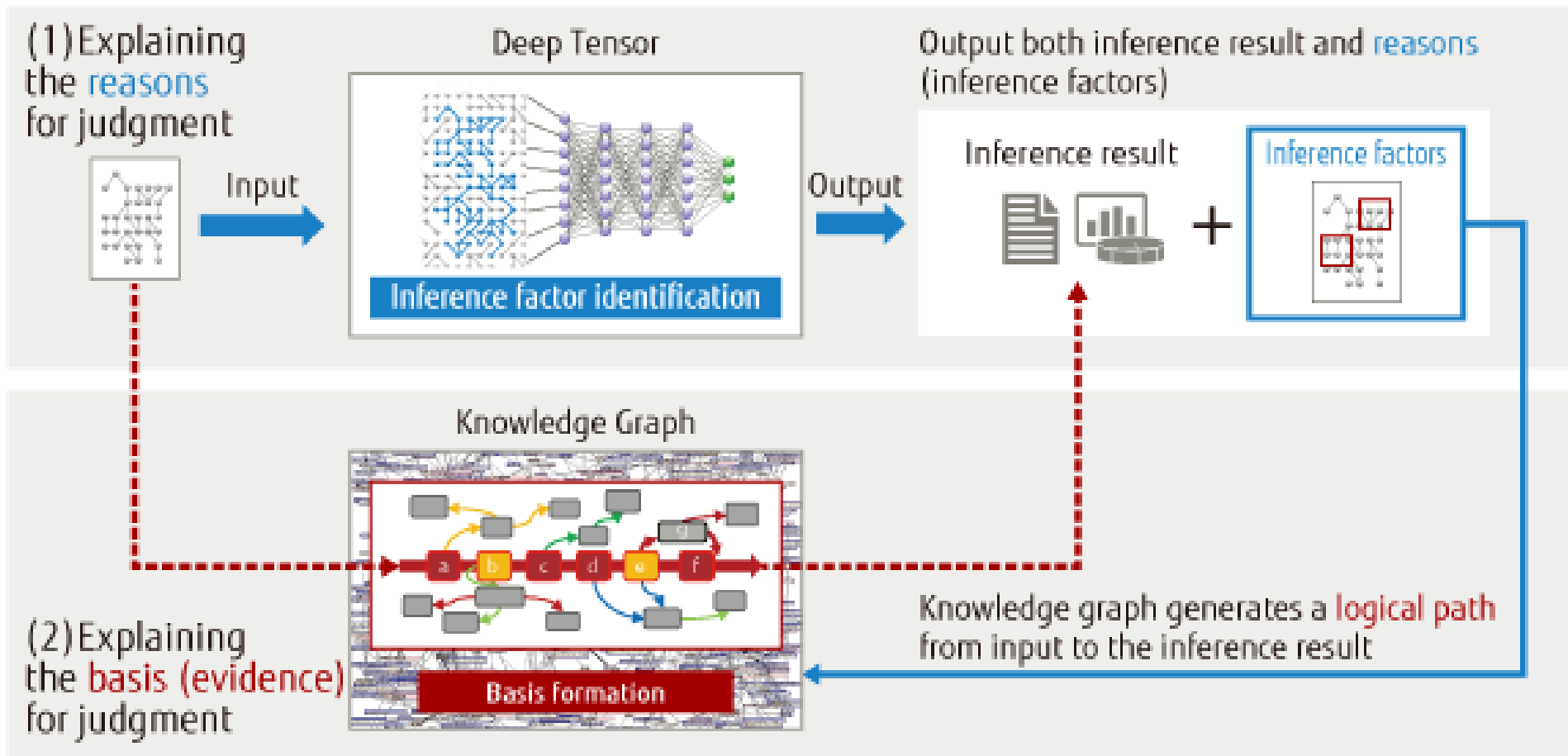


Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller & Wojciech Samek 2015. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. PloS one, 10, (7), e0130140, doi:10.1371/journal.pone.0130140.

# 07 Towards human interpretable models

**End-users shall be able to retrace  
the results on demand  
and we engineers need to  
understand our own  
machine learning models!**



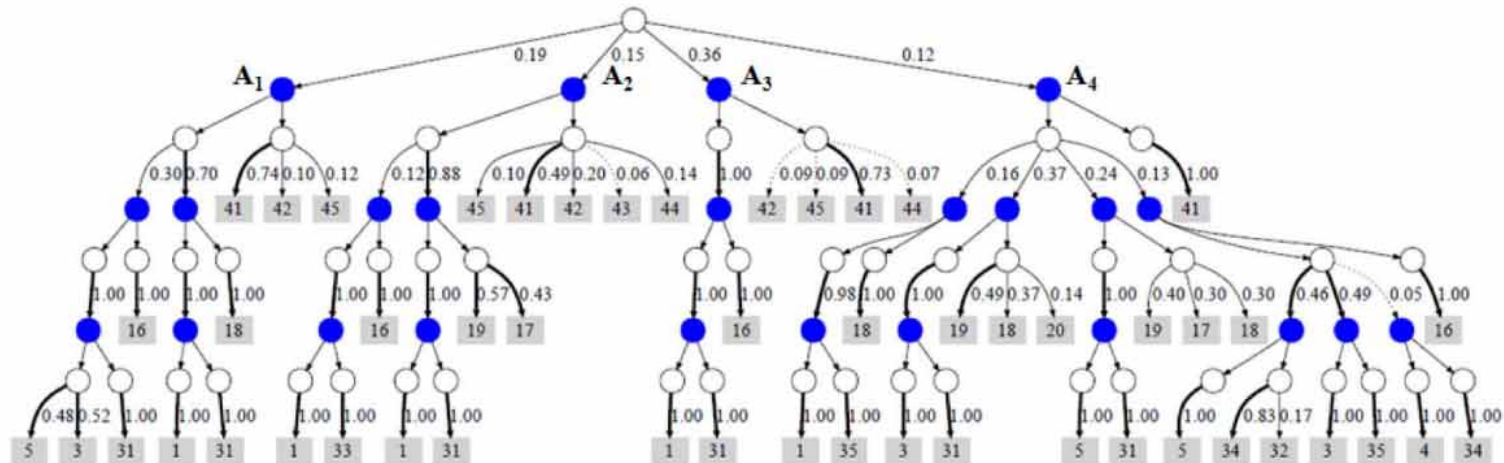


Randy Goebel, Ajay Chander, Katharina Holzinger, Freddy Lecue, Zeynep Akata, Simone Stumpf, Peter Kieseberg & Andreas Holzinger 2018. Explainable AI: the new 42? Springer Lecture Notes in Computer Science LNCS 11015

Input images



Stochastic AOT



Part dictionary  
(terminal nodes)

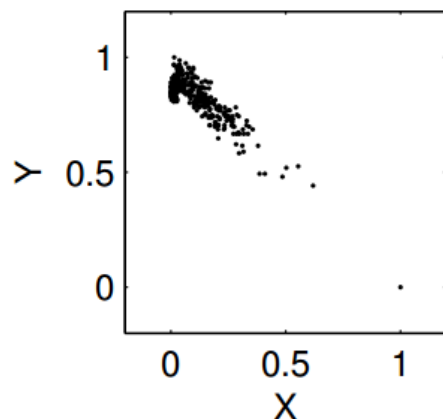
|          | 1 | 2 | 3 | 4 | 5 | 16 | 17 | 18 | 19 | 20 | 31 | 32 | 33 | 34 | 35 | 41 | 42 | 43 | 44 | 45 |
|----------|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| sketch   |   |   |   |   |   |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |
| texture  |   |   |   |   |   |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |
| flatness |   |   |   |   |   |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |

Valid configurations

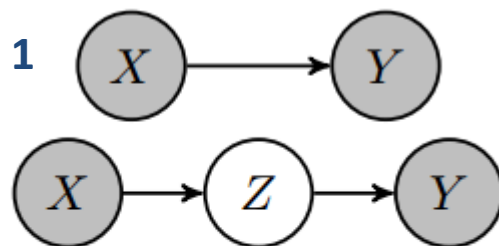


Zhangzhang Si & Song-Chun Zhu 2013. Learning and-or templates for object recognition and detection. IEEE transactions on pattern analysis and machine intelligence, 35, (9), 2189-2205, doi:10.1109/TPAMI.2013.35.

# 08 Digression: Causal Reasoning (but we need it)

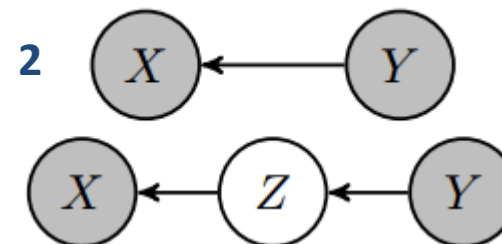


Joris M. Mooij, Jonas Peters, Dominik Janzing, Jakob Zscheischler & Bernhard Schölkopf 2016. Distinguishing cause from effect using observational data: methods and benchmarks. The Journal of Machine Learning Research, 17, (1), 1103-1204.



$$\mathbb{P}_Y \neq \mathbb{P}_{Y|\text{do}(x)} = \mathbb{P}_{Y|x}$$

$$\mathbb{P}_X = \mathbb{P}_{X|\text{do}(y)} \neq \mathbb{P}_{X|y}$$



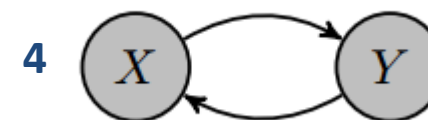
$$\mathbb{P}_Y = \mathbb{P}_{Y|\text{do}(x)} \neq \mathbb{P}_{Y|x}$$

$$\mathbb{P}_X \neq \mathbb{P}_{X|\text{do}(y)} = \mathbb{P}_{X|y}$$



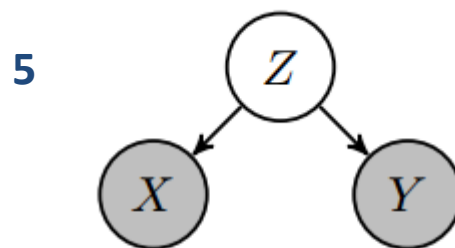
$$\mathbb{P}_Y = \mathbb{P}_{Y|\text{do}(x)} = \mathbb{P}_{Y|x}$$

$$\mathbb{P}_X = \mathbb{P}_{X|\text{do}(y)} = \mathbb{P}_{X|y}$$



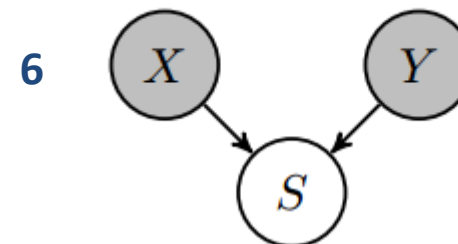
$$\mathbb{P}_Y \neq \mathbb{P}_{Y|\text{do}(x)} \neq \mathbb{P}_{Y|x}$$

$$\mathbb{P}_X \neq \mathbb{P}_{X|\text{do}(y)} \neq \mathbb{P}_{X|y}$$



$$\mathbb{P}_Y = \mathbb{P}_{Y|\text{do}(x)} \neq \mathbb{P}_{Y|x}$$

$$\mathbb{P}_X = \mathbb{P}_{X|\text{do}(y)} \neq \mathbb{P}_{X|y}$$

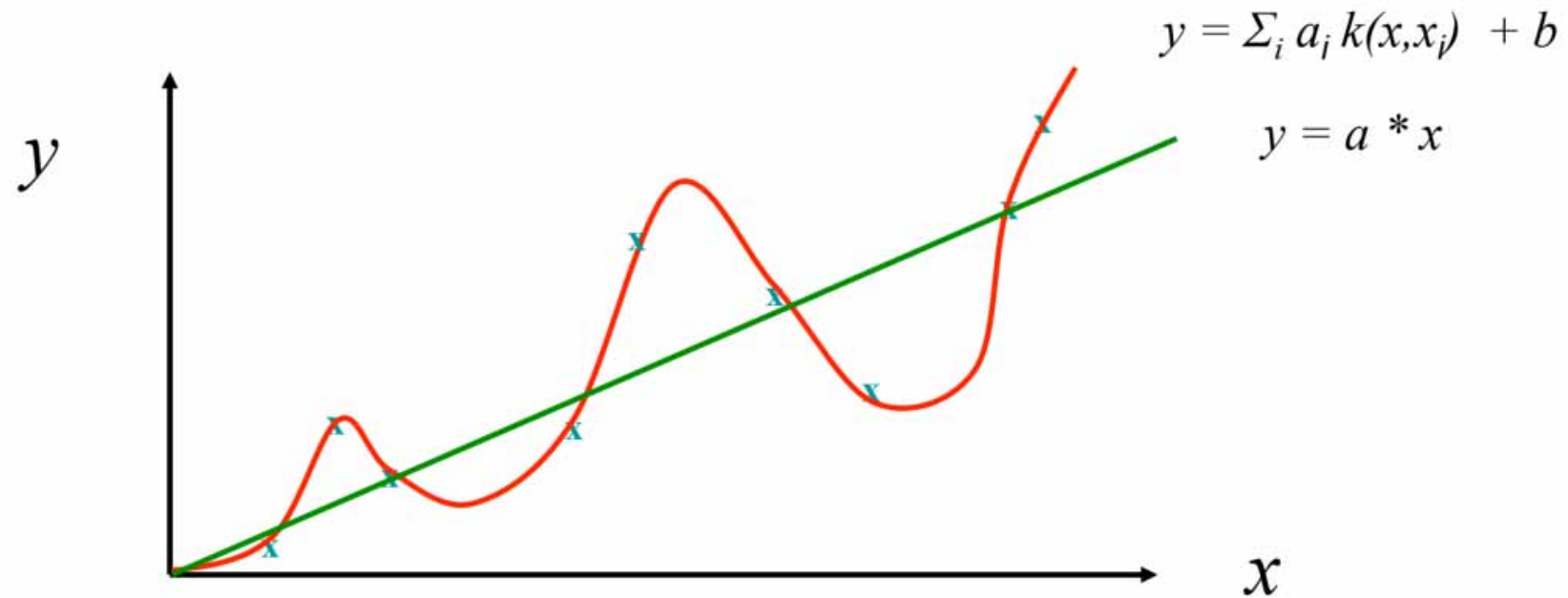


$$\mathbb{P}_{Y|s} \neq \mathbb{P}_{Y|\text{do}(x),s} = \mathbb{P}_{Y|x,s}$$

$$\mathbb{P}_{X|s} \neq \mathbb{P}_{X|\text{do}(y),s} = \mathbb{P}_{X|y,s}$$

- **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.



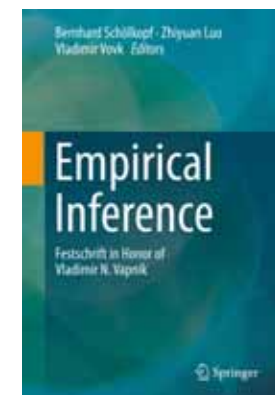
Gottfried W. Leibniz (1646-1716)

Hermann Weyl (1885-1955)

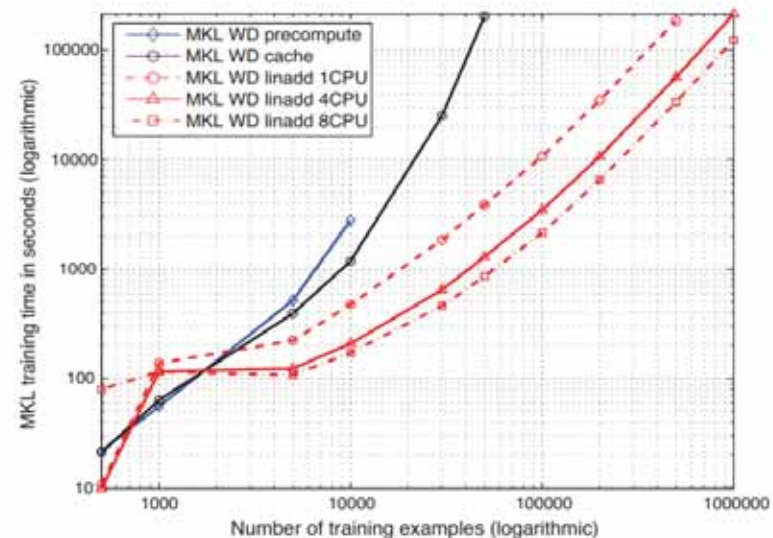
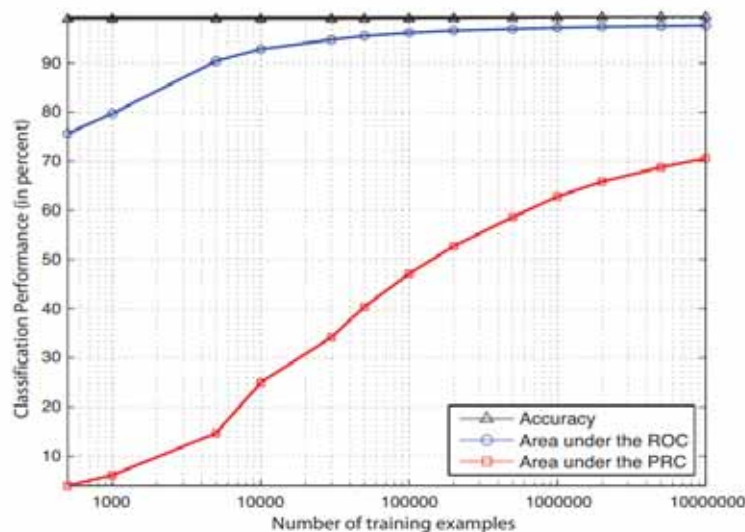
Vladimir Vapnik (1936-)

Alexey Chervonenkis (1938-2014)

Gregory Chaitin (1947-)



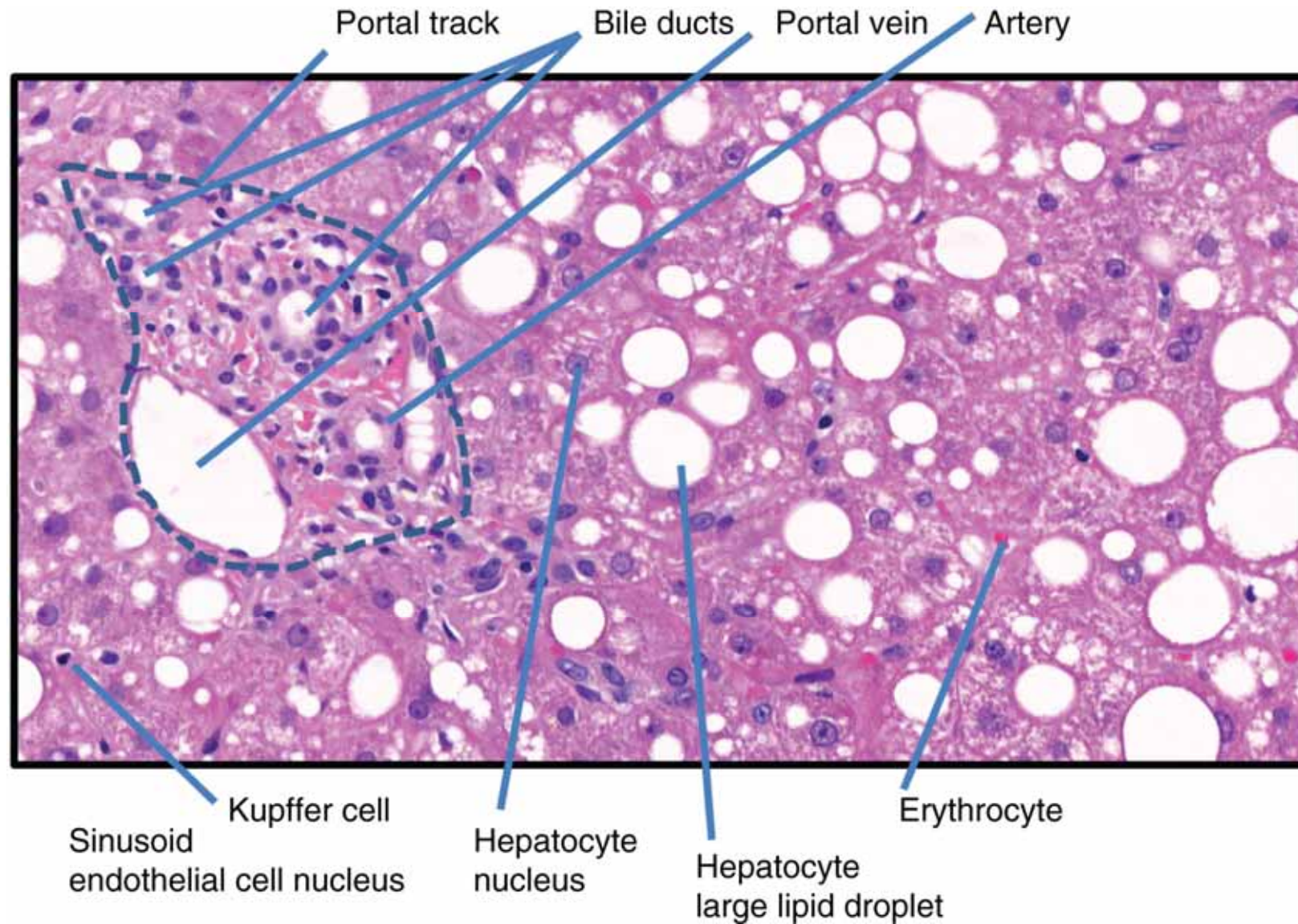
- High dimensionality (curse of dim., many factors contribute)
- Complexity (real-world is non-linear, non-stationary, non-IID \*)
- Need of large top-quality data sets
- Little prior data (no mechanistic models of the data)
  - \*) = Def.: a sequence or collection of random variables is independent and identically distributed if each random variable has the same probability distribution as the others and all are mutually independent



Sören Sonnenburg, Gunnar Rätsch, Christin Schaefer & Bernhard Schölkopf 2006. Large scale multiple kernel learning. Journal of Machine Learning Research, 7, (7), 1531-1565.



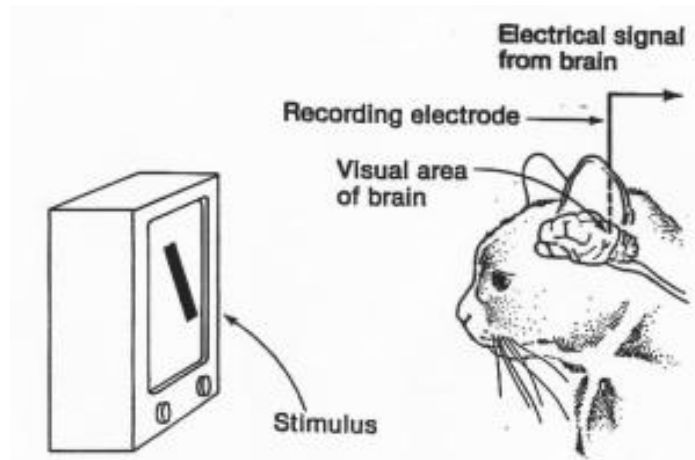
# 09 Measuring Machine Intelligence #KANDINSKYpatterns



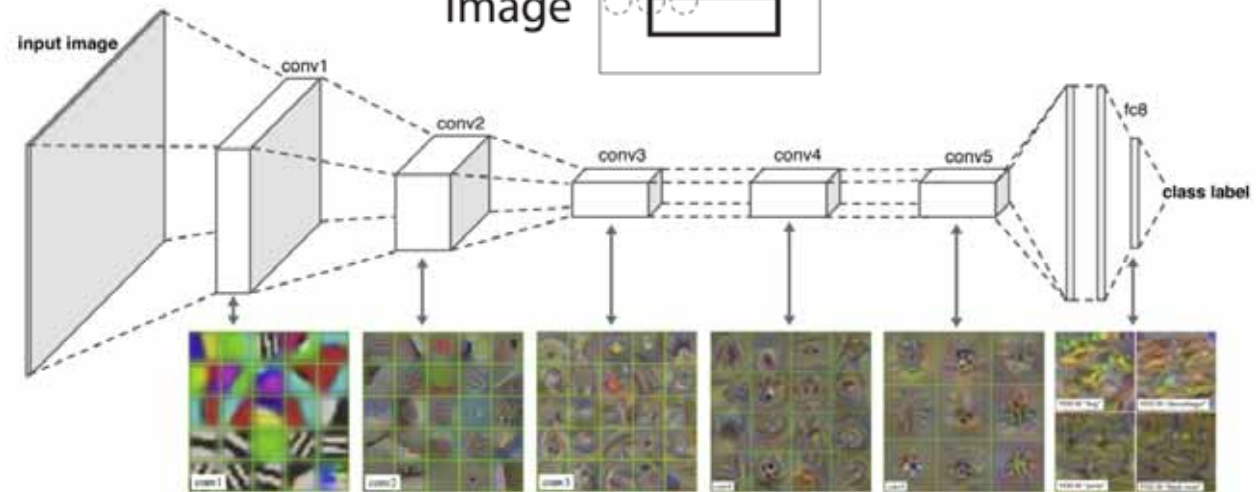
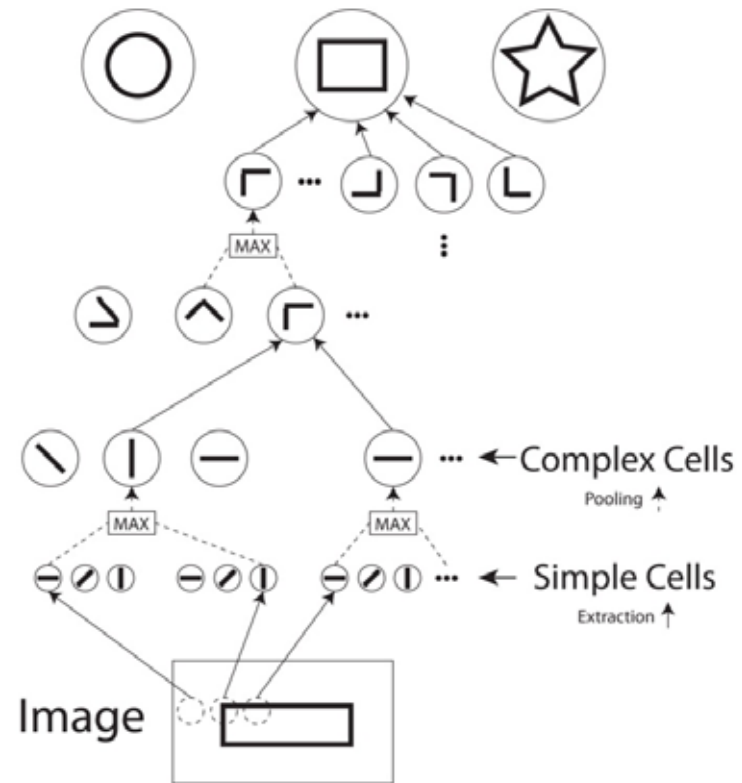
Andreas Holzinger, Georg Langs, Helmut Denk, Kurt Zatloukal & Heimo Mueller 2019. Causability and Explainability of AI in Medicine. Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery, doi:10.1002/widm.1312.

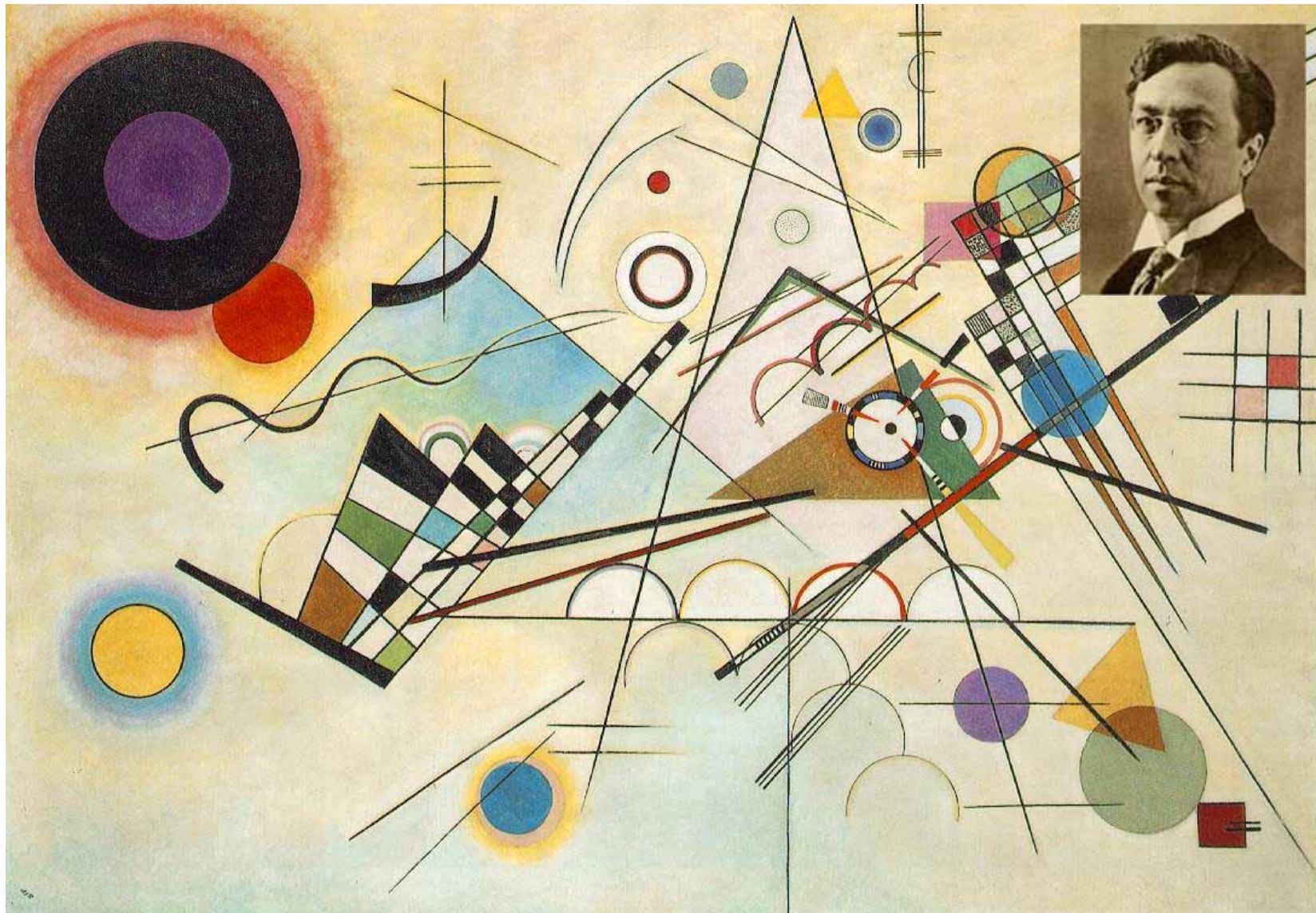
|  |                                      |   |
|--|--------------------------------------|---|
| <b>Radiologischer Befund</b>   |                                      | angelegt am 06.05.2006/20:26<br>geschr. von [REDACTED]<br>gedruckt am 17.11.2006/08:24<br>Anfo: NCHIN |
| <b>Kurzanamnese:</b>   | St.p. SHT                            |   |
| <b>Fragestellung:</b>  | -                                    |   |
| <b>Untersuchung:</b>   | Thorax eine Ebene liegend [REDACTED] |   |
| SB   |                                      |   |
| Bewegungsartefakte. Zustand nach Schädelhirntrauma.  |                                      |   |
| Das Cor in der Größennorm, keine akuten Stauungszeichen.<br>Fragliches Infiltrat parahilär li. im UF, RW-Erguss li.  |                                      |   |
| Zustand nach Anlage eines ET, die Spitze ca. 5cm cranial der Bifurkation, lieg. MS, orthotop positioniert. ZVK über re., die Spitze in Proj. auf die VCS. Kein Hinweis auf Pneumothorax.<br>Der re. Rezessus frei. |                                      |   |
| Mit kollegialen Grüßen   |                                      |   |
| [REDACTED]   |                                      |   |
| *** Elektronische Freigabe durch [REDACTED] am 09.05.2006 ***  |                                      |   |

Holzinger, A., Geierhofer, R. & Errath, M. 2007. Semantische Informationsextraktion in medizinischen Informationssystemen. *Informatik Spektrum*, 30, (2), 69-78.

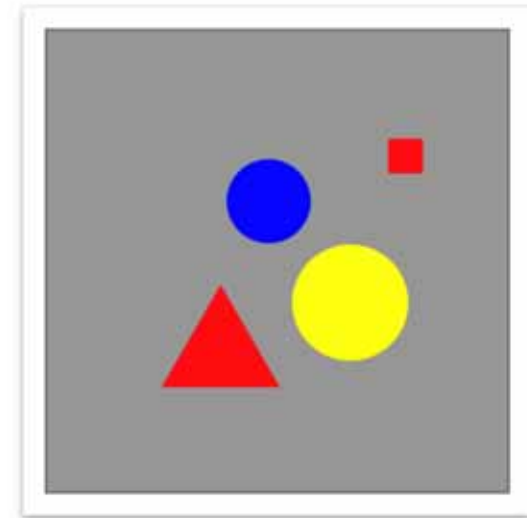


David H. Hubel & Torsten N. Wiesel  
1962. Receptive fields, binocular  
interaction and functional  
architecture in the cat's visual cortex.  
*The Journal of Physiology*, 160, (1),  
106-154,  
doi:10.1113/jphysiol.1962.sp006837.





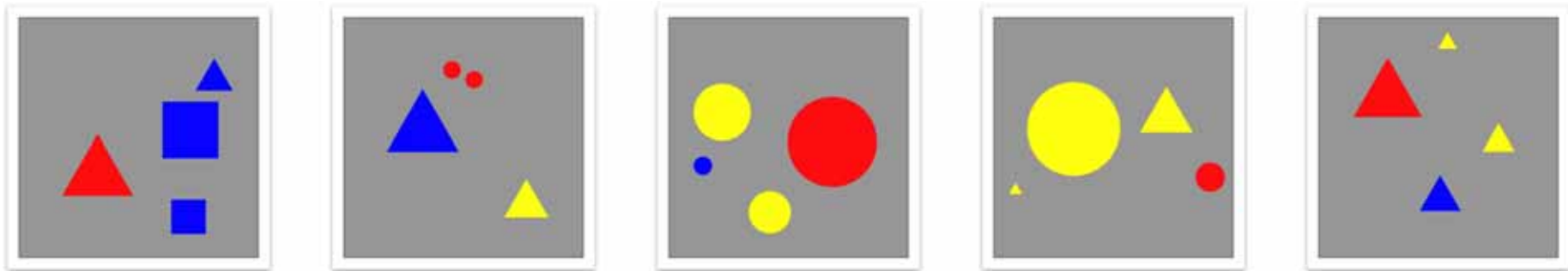
Komposition VIII, 1923, Solomon R. Guggenheim Museum, New York. Source: [https://de.wikipedia.org/wiki/Wassily\\_Kandinsky](https://de.wikipedia.org/wiki/Wassily_Kandinsky)  
This images are in the public domain.



- ... 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  $B \in \{0,1\}$
  - 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 following 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.

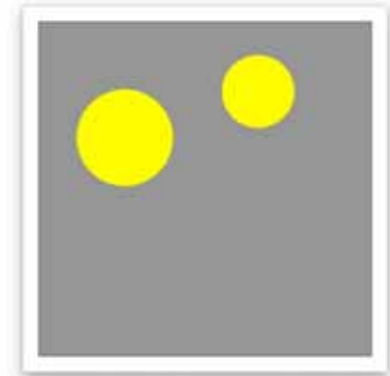
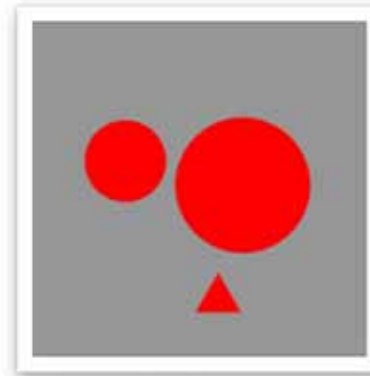
- ... is defined as the subset of all possible Kandinsky Figures  $k$  with  $s(k) \rightarrow 1$  or the natural language statement is true.
- $s(k)$  and a natural language statement are equivalent, if and only if the resulting Kandinsky Patterns contains the same Kandinsky Figures.
- $s(k)$  and the natural language statement are defined as the **Ground Truth** of a Kandinsky Pattern



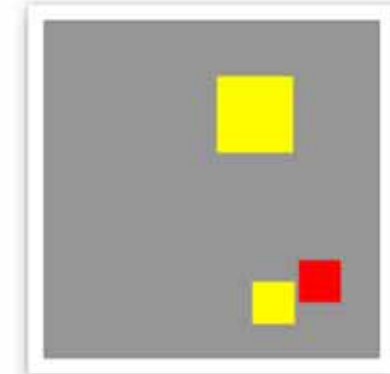
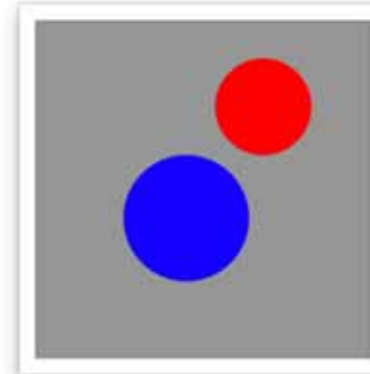
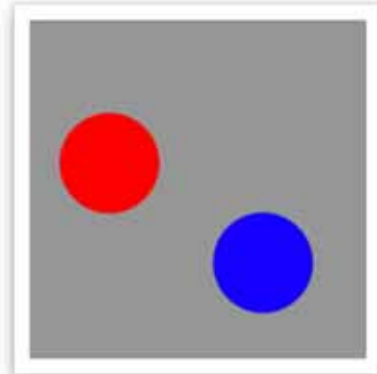
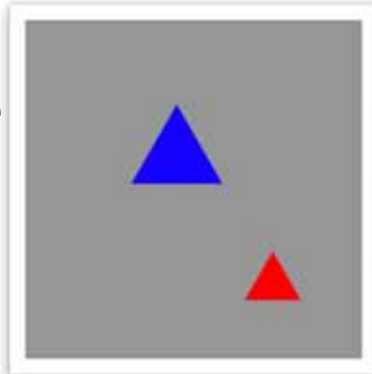
*"... the Kandinsky Figure has two pairs of objects with the same shape, in one pair the objects have the same color, in the other pair different colors, two pairs are always disjunct, i.e. they don't share a object ...".*



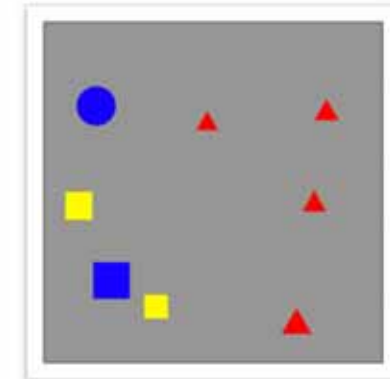
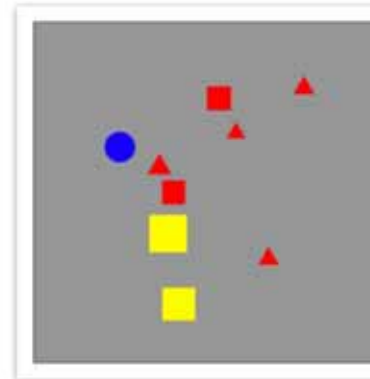
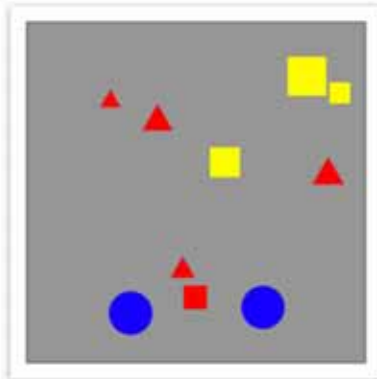
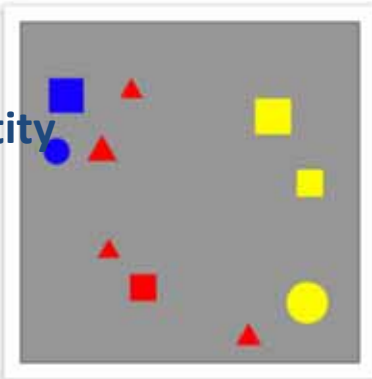
A  
Colour



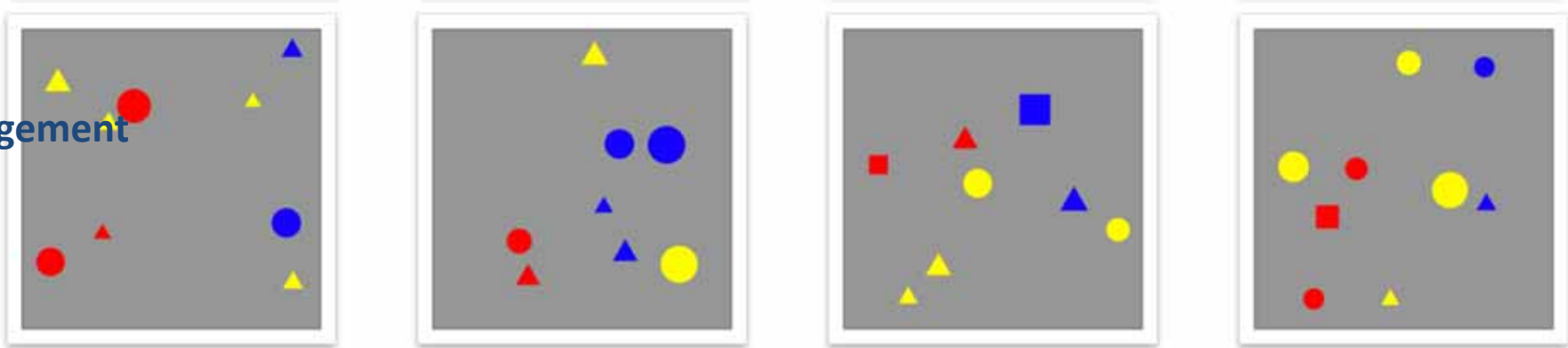
B  
Shape



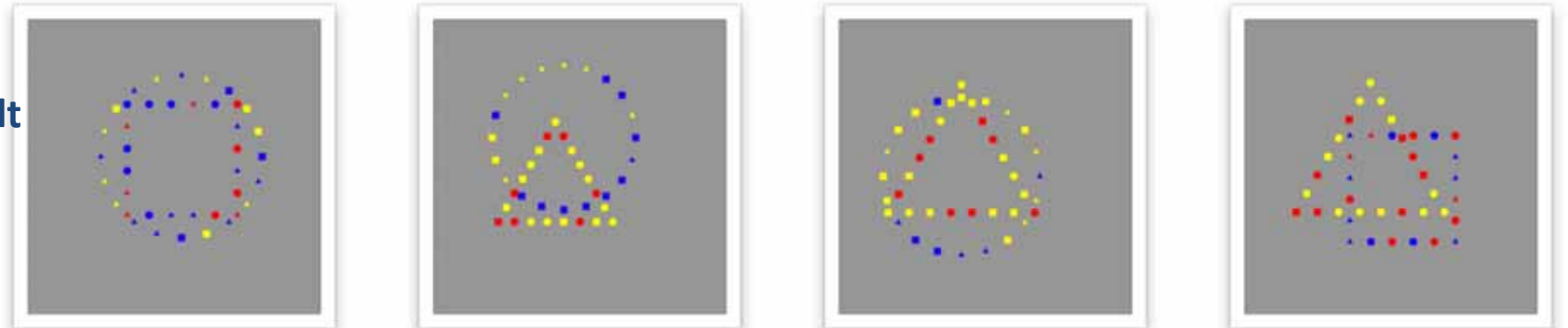
C  
Quantity



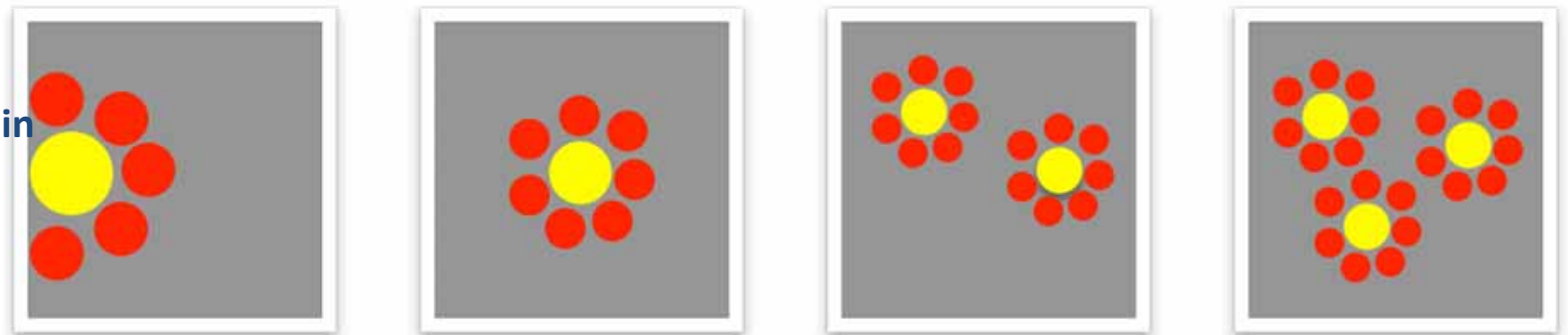
D  
Arrangement



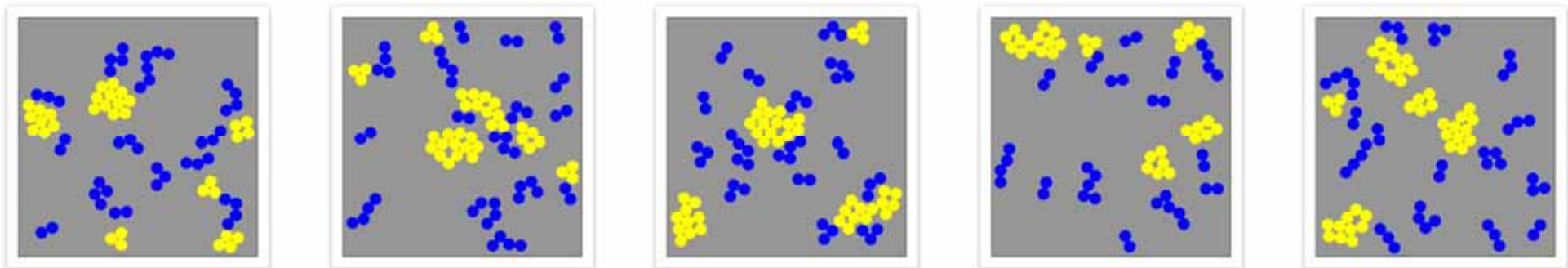
E  
Gestalt



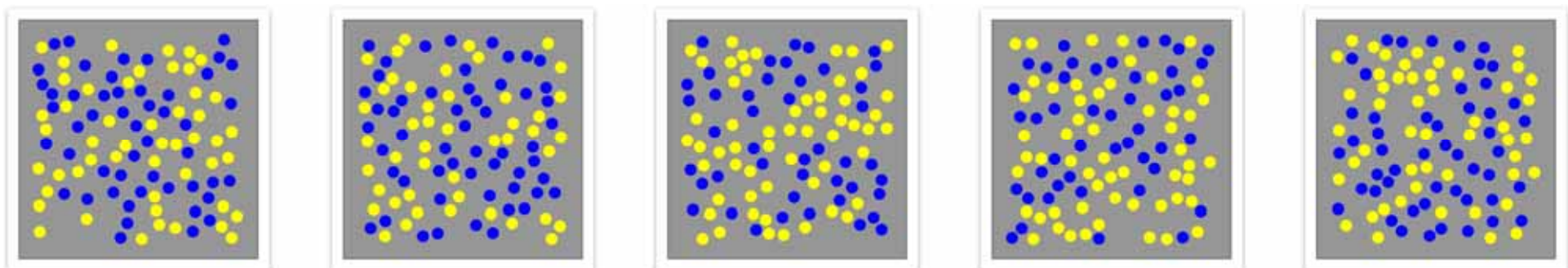
F  
Domain



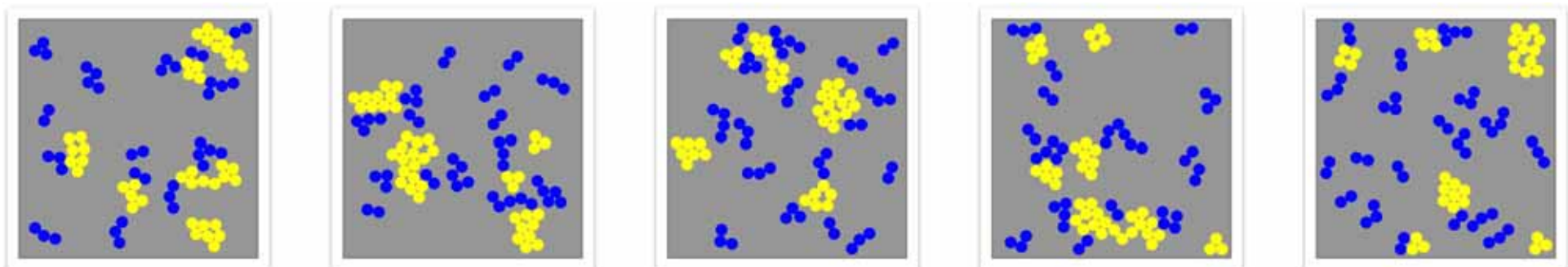
A) True



B) False

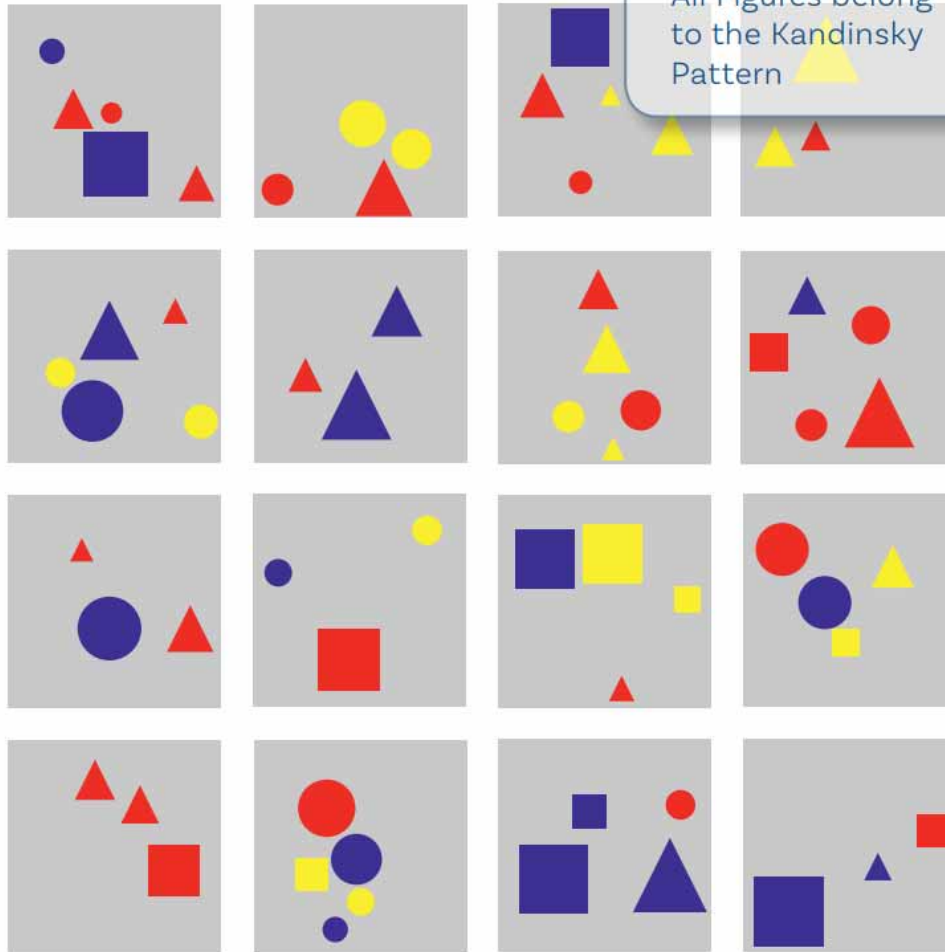


C) Counterfactual





☰ Part of the pattern



All Figures belong to the Kandinsky Pattern

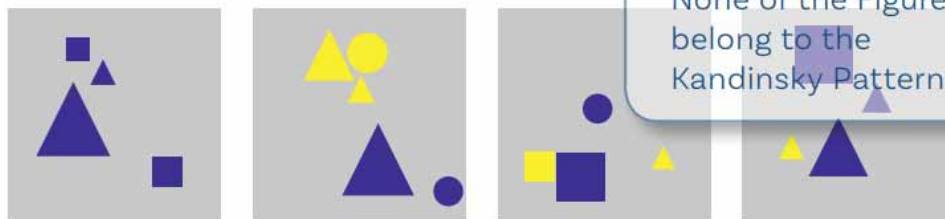
Hypothesis 1

*It only contains circles and triangles.*

Hypothesis 2

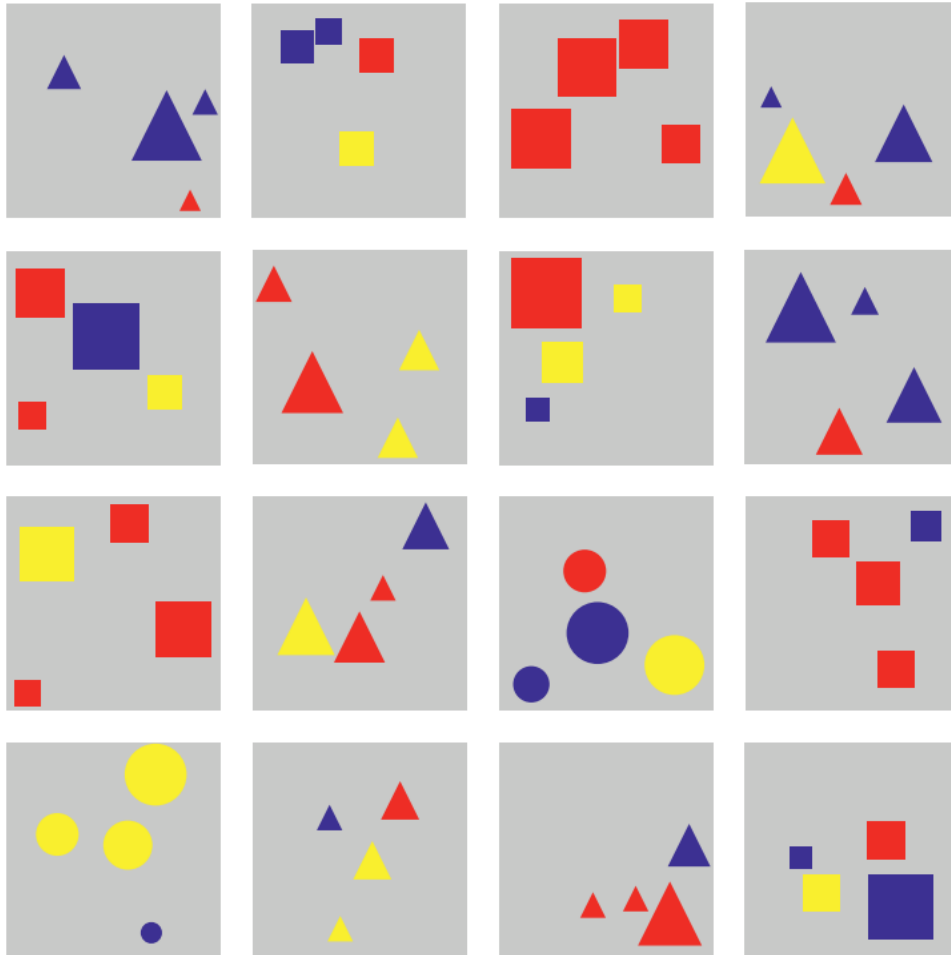
*It contains at least a red object.* ✓

≠ Not part of the pattern



None of the Figures belong to the Kandinsky Pattern

☰ Part of the pattern

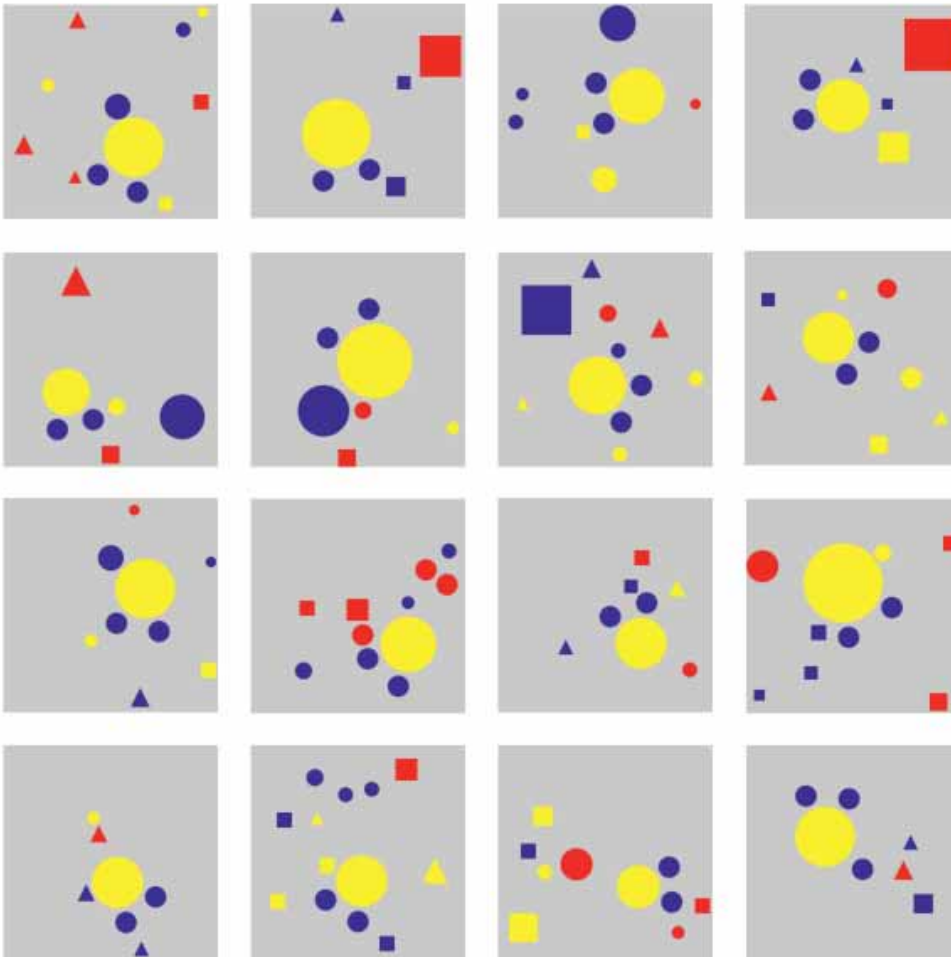


⚡ Not part of the pattern

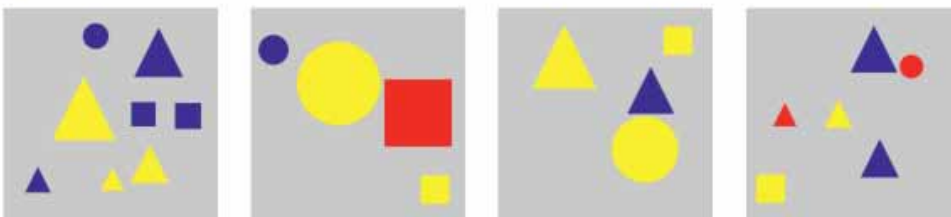


S2 Basic Pattern 2  
 Title: **All of Same Shape** ->  
 All objects have the same shape.  
 Hint: Don't be distracted by the colors

⊃ Part of the pattern



⊄ Not 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.

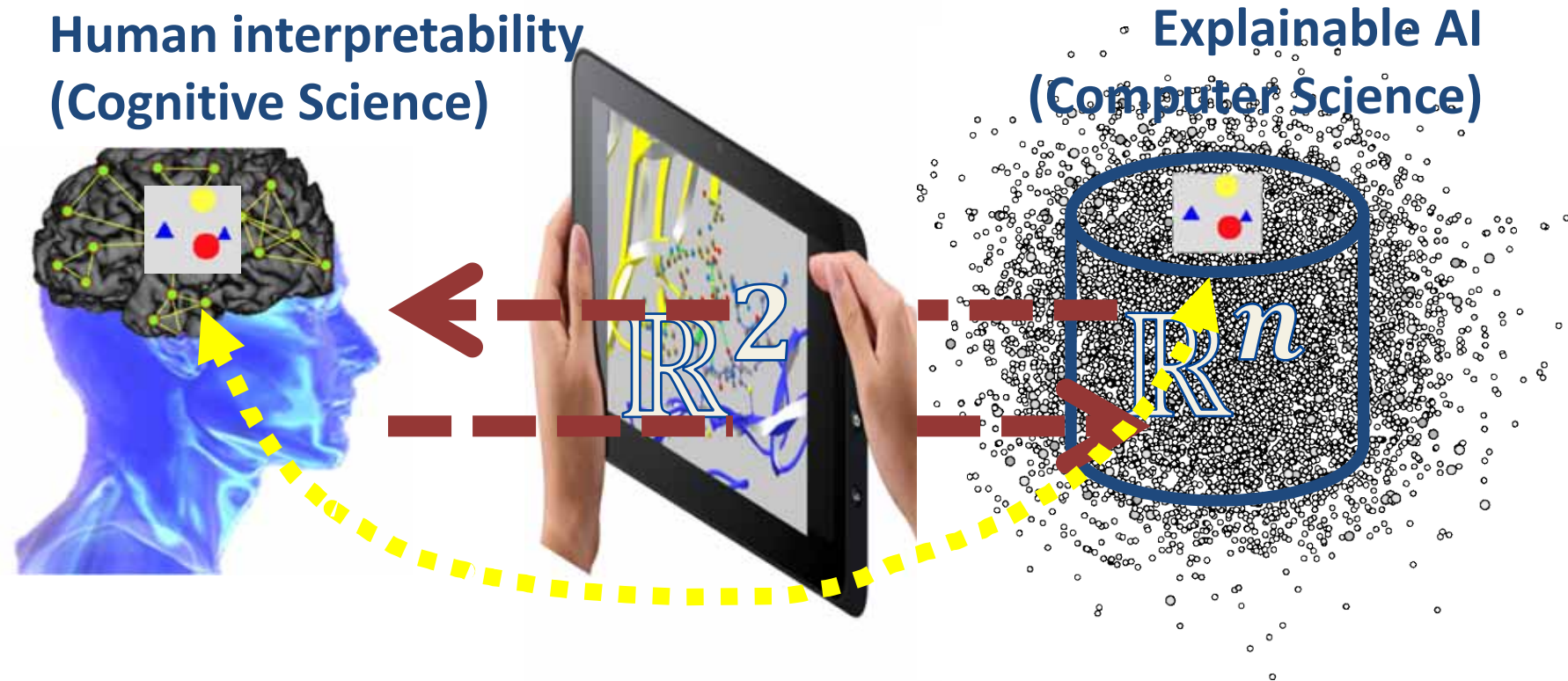
# Conclusion



- The results show that the majority of explanations was made based on the properties of individual elements in an image
- i.e., shape, color, size, ...
- and the appearance of individual objects (number)
- Comparisons of elements (e.g., more, less, bigger, smaller, etc.) were significantly less likely, and
- the location of objects, interestingly, played almost no role in the explanation of the images.

- Although humans tend to make more errors, human intelligence is more reliable and robust against catastrophic errors, whereas
- AI is vulnerable against software, hardware and energy failures.
- Human intelligence develops based on infinite interactions with an infinite environment, while AI is limited to the small world of a particular task.
- The development of intelligence, therefore, is the result of the incremental interplay between challenge/task, a conceptual change (physiological as well as mentally) of the system, and the assessment of the effects of the conceptual change.
- To advance AI, specifically in the direction of explainable AI, we suggest bridging the human strength and the human assessment methods with those of AI.

- Causability := a property of a person (Human)
- Explainability := a property of a system (Computer)



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.

**Thank you!**  
**Now, take part in the**  
**explainable-AI challenge:**

**<https://human-centered.ai/kandinsky-challenge>**



**August 26-29, 2020  
in Dublin**

**<https://cd-make.net>  
#cdmake**