

From Explainable AI to Human-Centered AI

Andreas Holzinger

Human-Centered AI Lab (Holzinger Group)
Institute for Medical Informatics/Statistics,
Medical University Graz, Austria
and TU Austria (TU Graz & TU Wien)



HCAI
HUMAN-CENTERED.AI

@aholzin #KandinskyPatterns



Biobank Graz
The largest biobank
in Europe

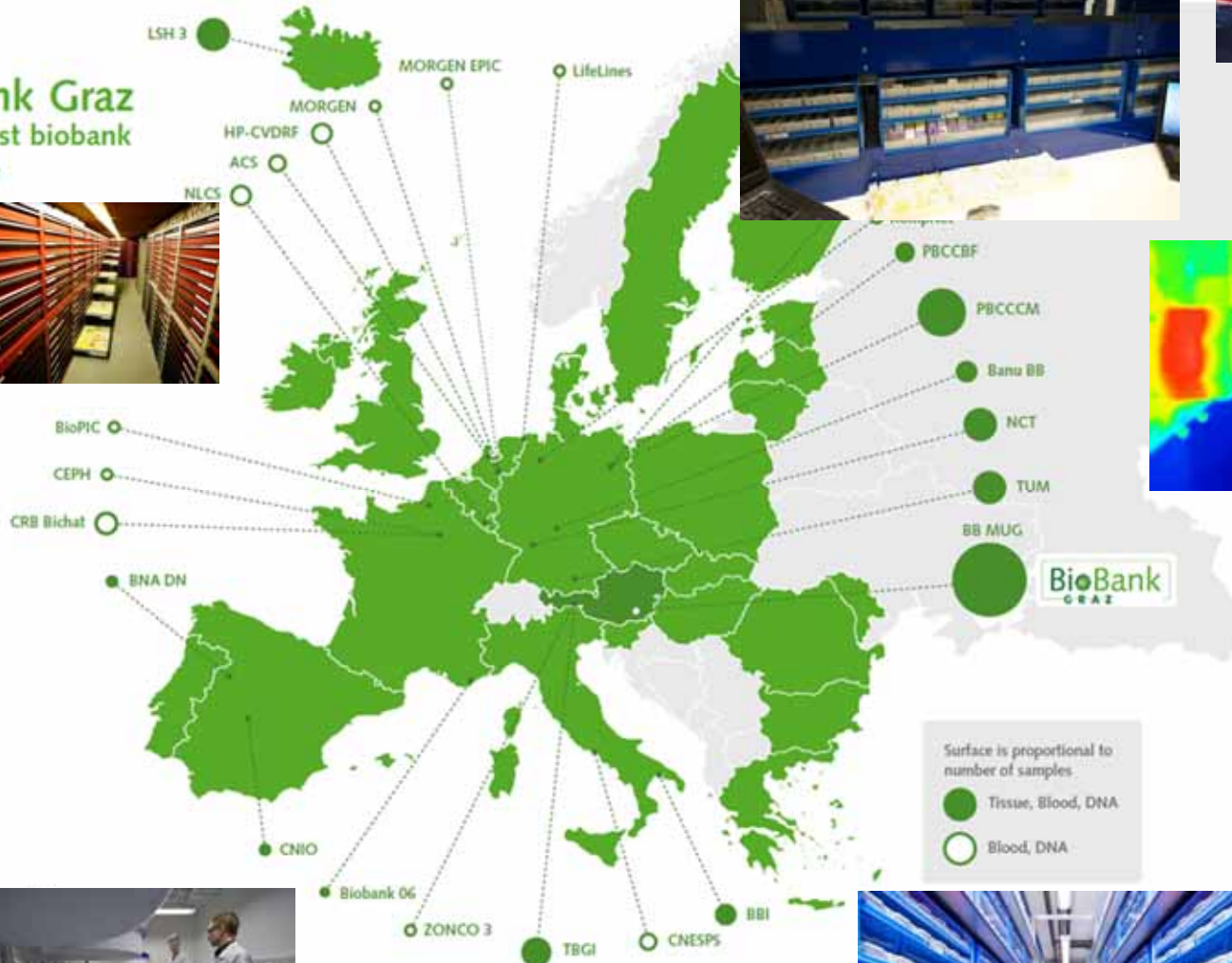
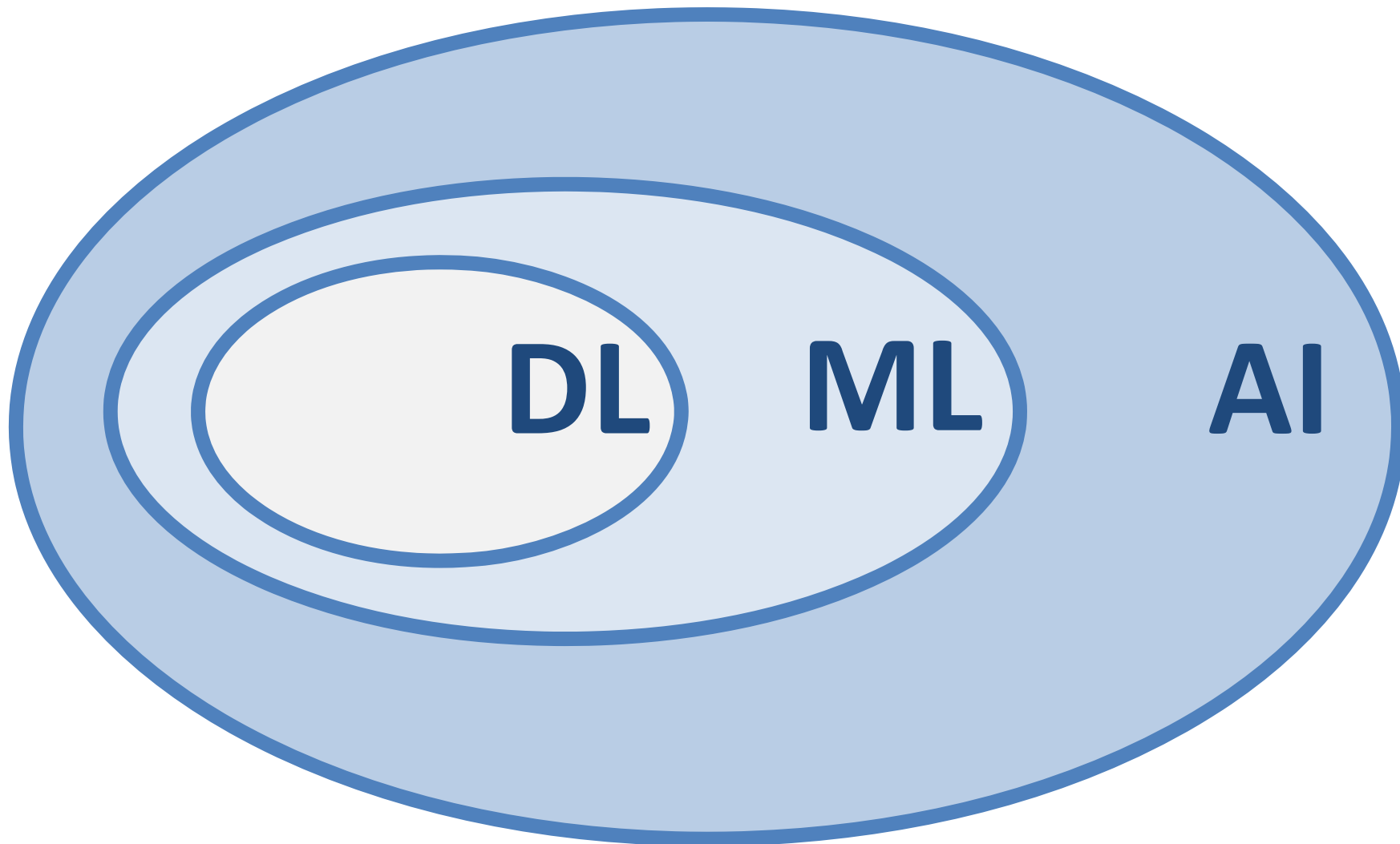




Image is in the public domain





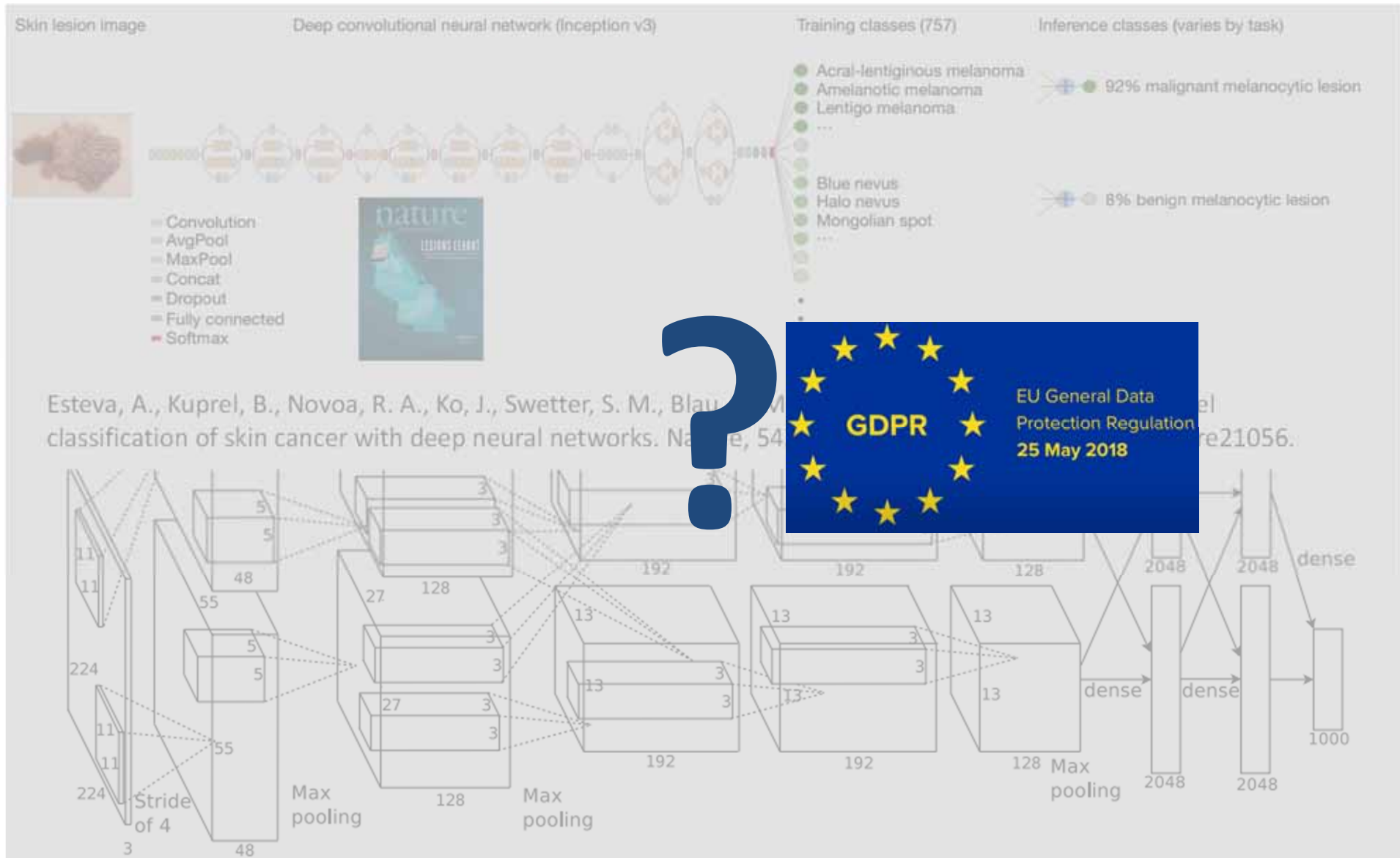


Andreas Holzinger, Peter Kieseberg, Edgar Weippl & A Min Tjoa 2018. Current Advances, Trends and Challenges of Machine Learning and Knowledge Extraction: From Machine Learning to Explainable AI. Springer Lecture Notes in Computer Science LNCS 11015. Cham: Springer, pp. 1-8, doi:[10.1007/978-3-319-99740-7_1](https://doi.org/10.1007/978-3-319-99740-7_1)

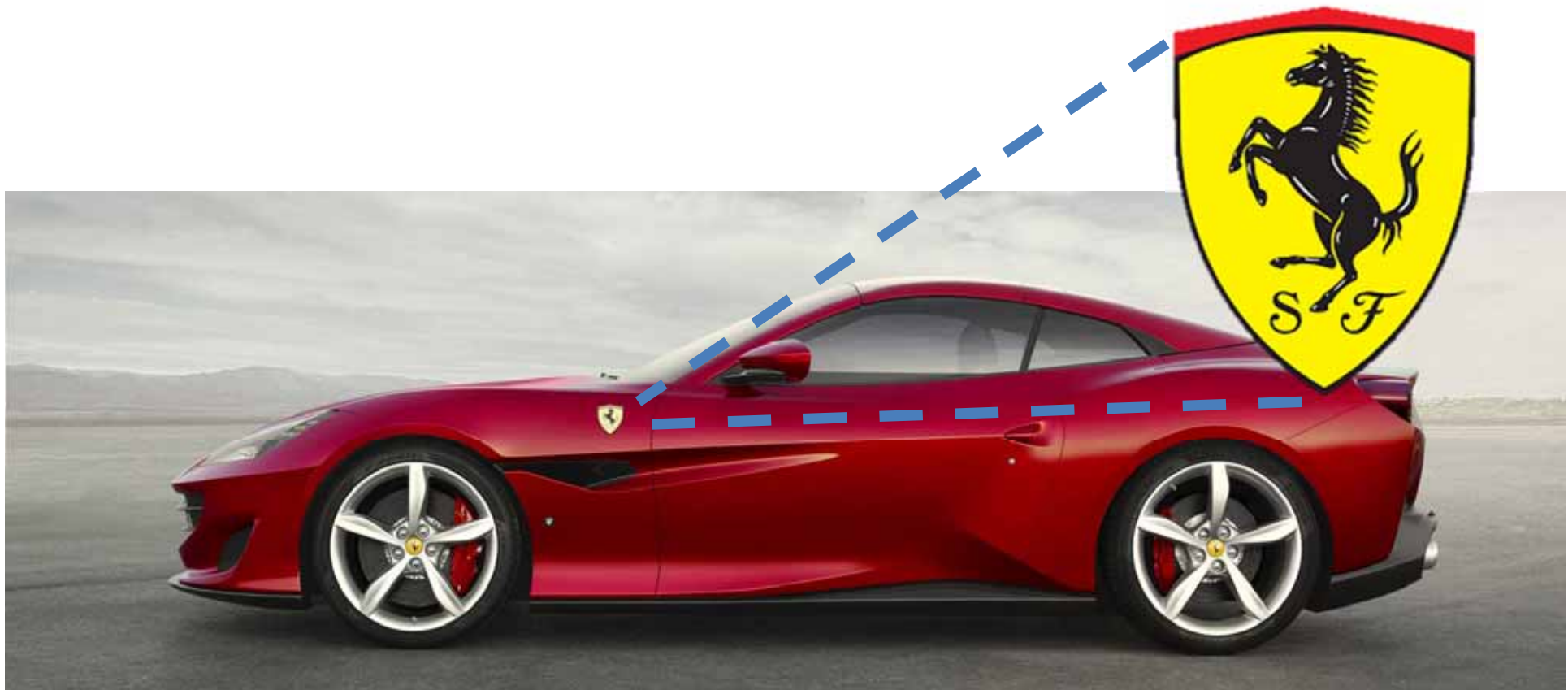


- **High-dimensional**
- **Non-convex**
- **Resource and data hungry**
- **Lacking interpretability ...**

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



- Result of the classifier: **This is a horse**
- **Why?**



Source: Image is in the public domain



+ .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 !

- 00 Motivation
- 01 What is the HCAI approach?
- 02 Application Area: Health Informatics
- 03 Probabilistic Information
- 04 aML
- 05 iML
- 06 Explainable AI (state-of-the-art)
- 07 Towards interpretable ML
- 08 Measuring Machine Intelligence (Kandinsky)

01 What is the

Privacy 4 – Transparency, Accountability, Ethics



Space and Time 5 - Graphs, 6-Topology, 7-Entropy

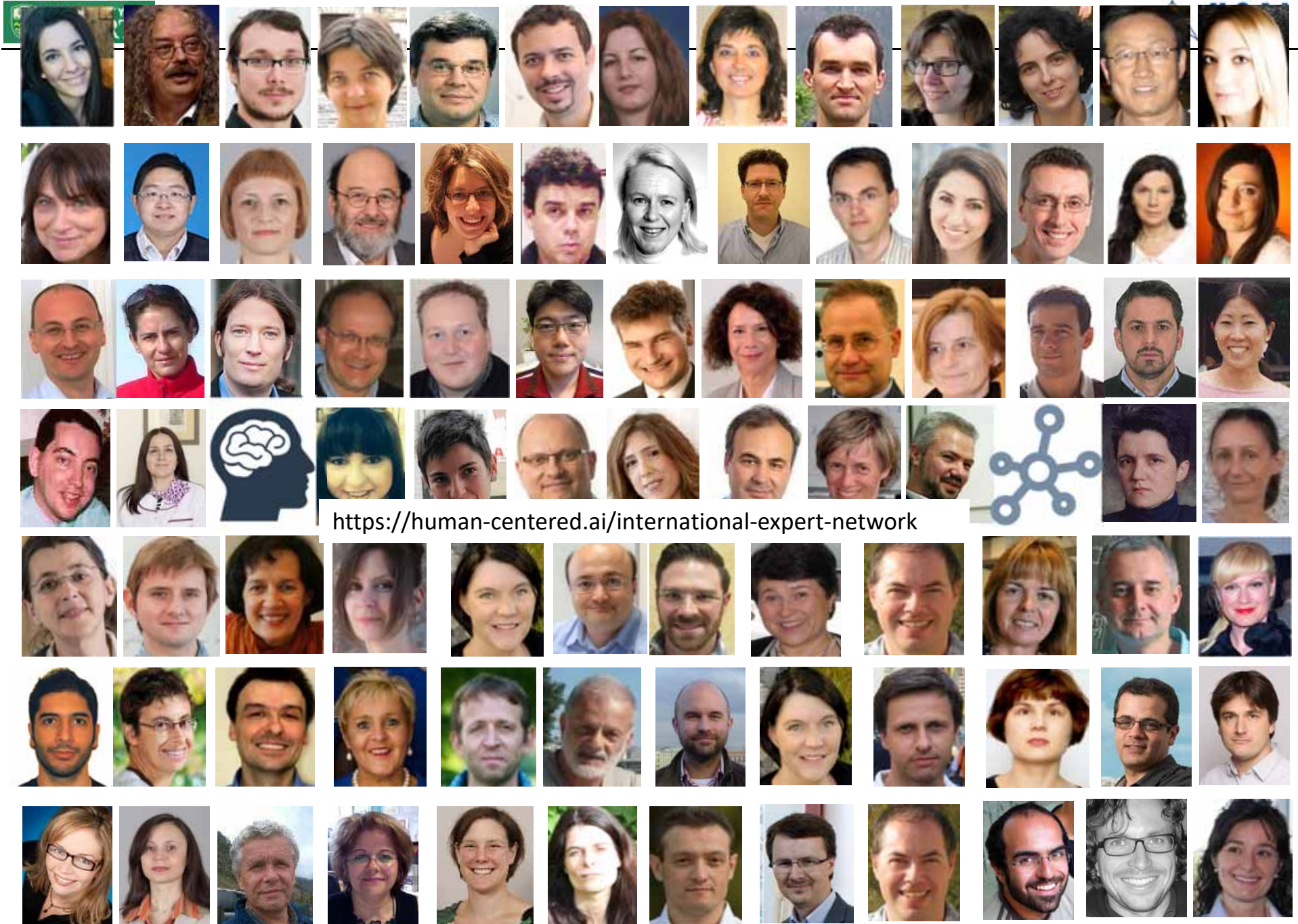
Andreas Holzinger 2017. Introduction to Machine Learning and Knowledge Extraction (MAKE). *Machine Learning and Knowledge Extraction*, 1, (1), 1-20, doi:10.3390/make1010001.

A large orchestra performing in a concert hall. The musicians are seated in rows, playing various instruments including violins, violas, cellos, double basses, trumpets, trombones, and percussion. A conductor is visible on the left side of the stage. The hall has a high ceiling with wooden paneling and a large organ in the background.

Needs a concerted effort

international
without boundaries ...

Image Source: <http://www.bach-cantatas.com>



<https://human-centered.ai/international-expert-network>



Lecture Notes in
Computer Science
LNCS LNAI LNBI

CD-MAKE 2019

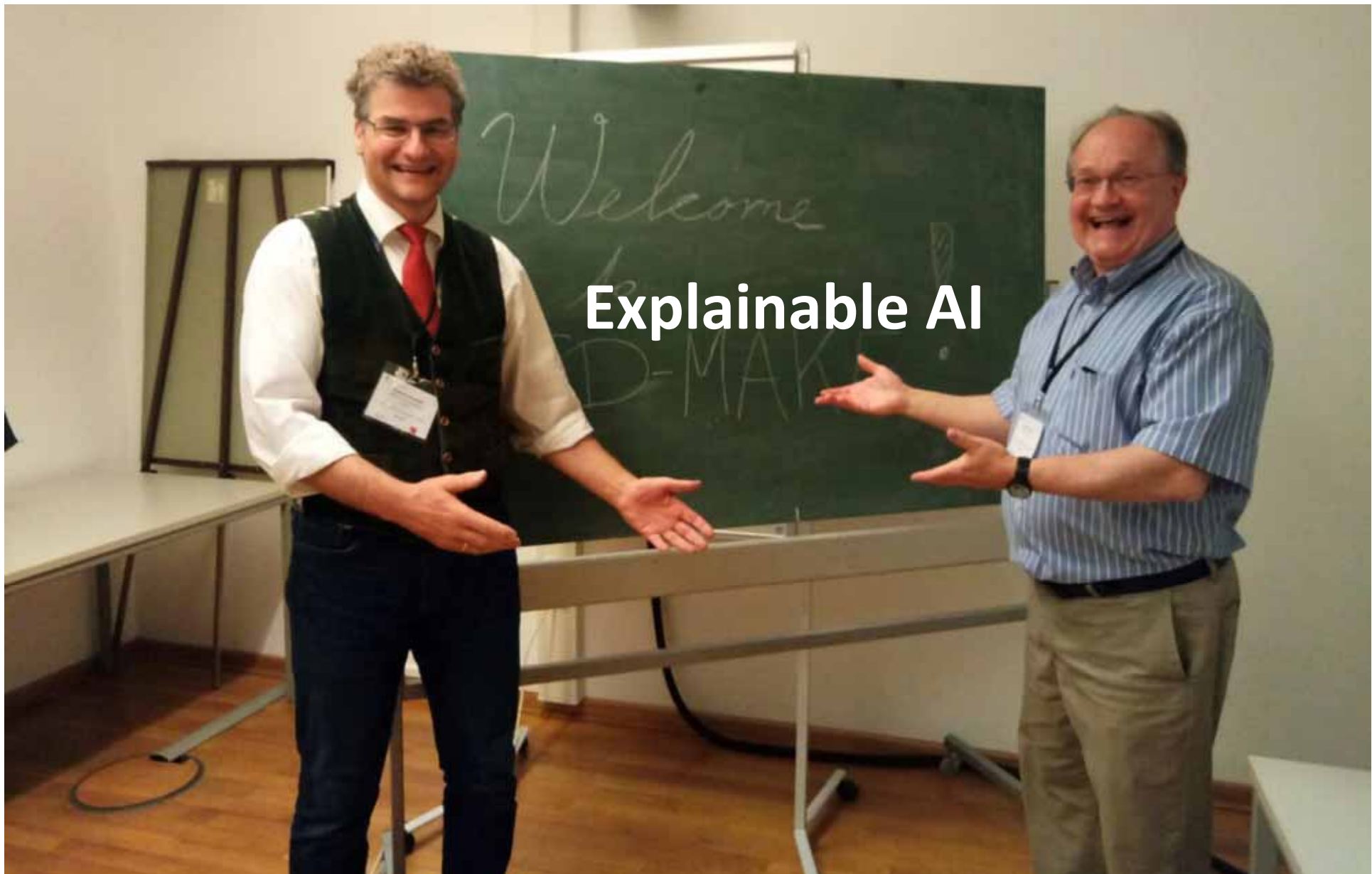


Cross Domain Conference for Machine Learning and Knowledge Extraction



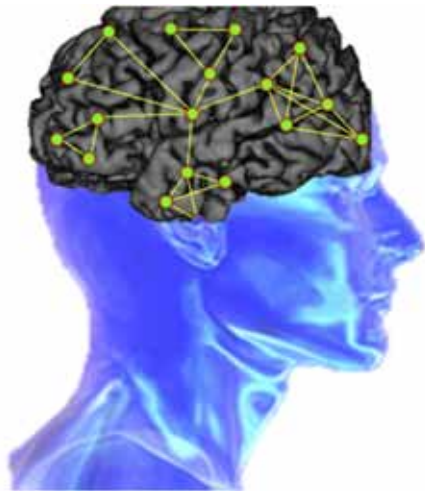
<https://cd-make.net>

Image with friendly permission of Michael D. Beckwith

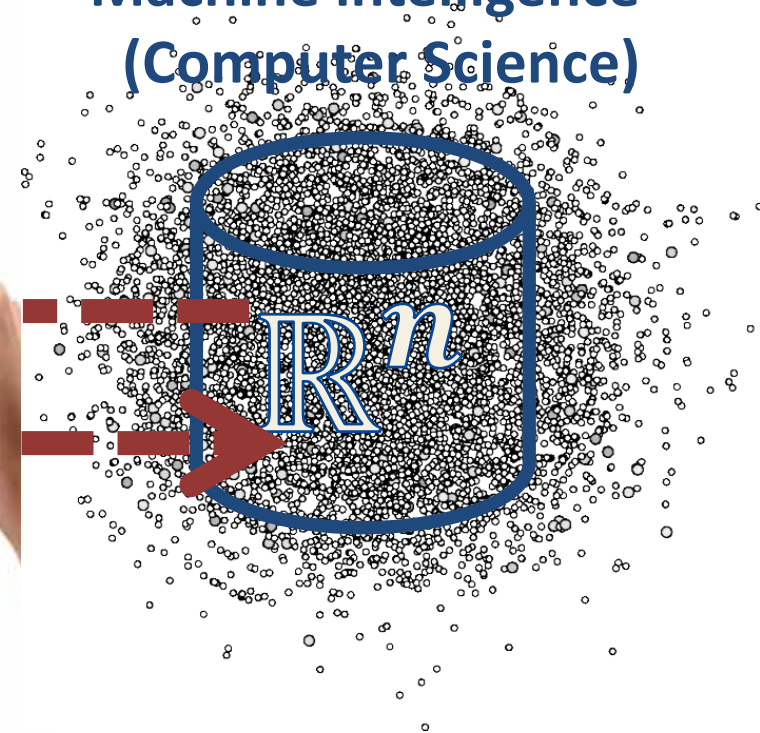


Our goal is that human values are aligned to ensure responsible machine learning

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 !



Humanoid AI

≠

Human-centered AI



02 Application Area Health Informatics

Why is this application area complex ?



Our central hypothesis: Information may bridge this gap

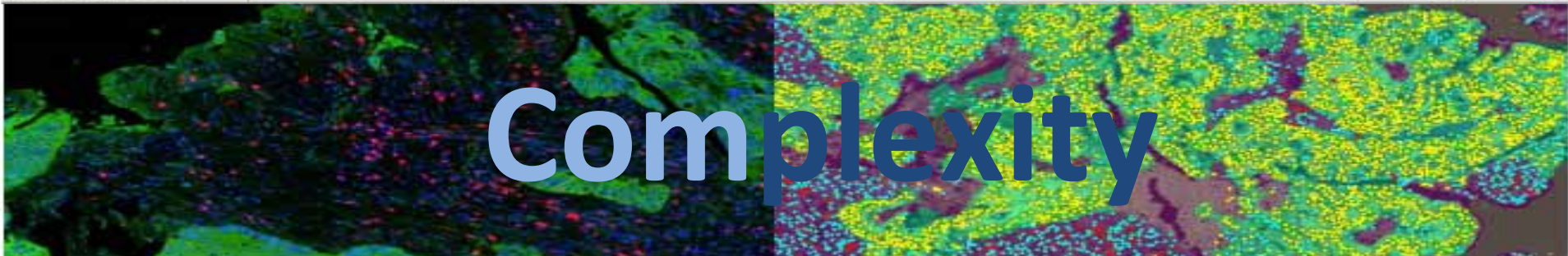
Andreas Holzinger & Klaus-Martin Simonic (eds.) 2011. Information Quality in e-Health. Lecture Notes in Computer Science LNCS 7058, Heidelberg, Berlin, New York: Springer, doi:10.1007/978-3-642-25364-5.



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	6	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	43	43	43	43	43	43	43	
Harnmenge/Zeit															
Harn/kg/Std															
Total Ma-Darm	6	6	6	6	6	6	6	6	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	29	118	118	118	22	140	140	
Nettobilanz 24h	+96	+105	+70	+70	+70	+70	+70	+95	+95	+114	+101	+106	+18	+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

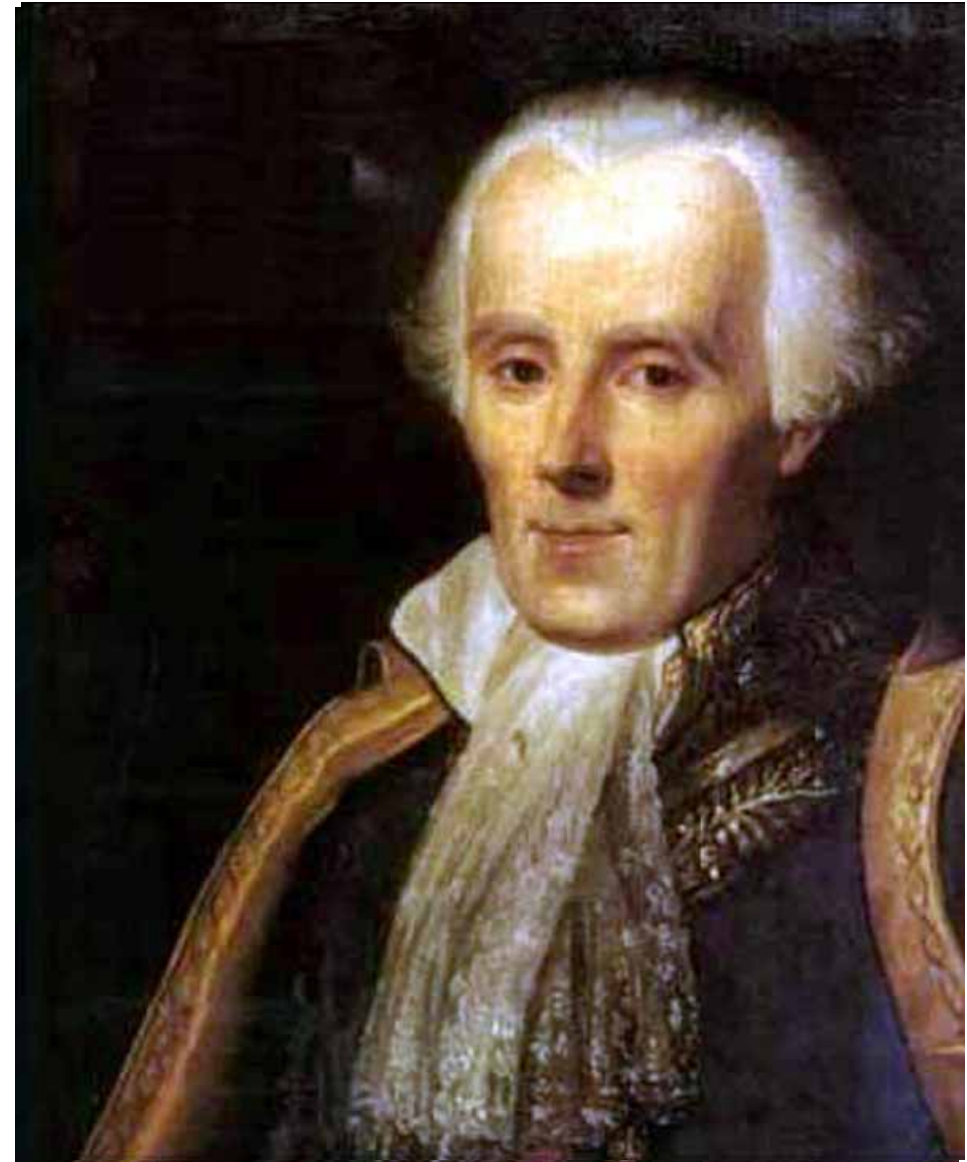
The true logic of this world is
in the calculus of
probabilities.

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



This image is in the Public Domain

Pierre Simon de Laplace (1749-1827)

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

Observed data:



\approx Training data: $\mathcal{D} = x_{1:n} = \{x_1, x_2, \dots, x_n\}$

Feature Parameter: θ or hypothesis h $h \in \mathcal{H}$

Prior belief \approx prior probability of hypothesis h : $p(\theta)$ $p(h)$

Likelihood $\approx p(x)$ of the data that h is true $p(\mathcal{D}|\theta)$ $p(d|h)$

Data evidence \approx marginal $p(x)$ that $h = \text{true}$ $p(\mathcal{D})$ $\sum_{h \in \mathcal{H}} p(d|h) * p(h)$

Posterior $\approx p(x)$ of h after seen ("learn") data d $p(\theta|\mathcal{D})$ $p(h|d)$

$$\text{posterior} = \frac{\text{likelihood} * \text{prior}}{\text{evidence}} \quad p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta) * p(\theta)}{p(\mathcal{D})}$$

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

d ... data

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

$\forall h, d \dots$

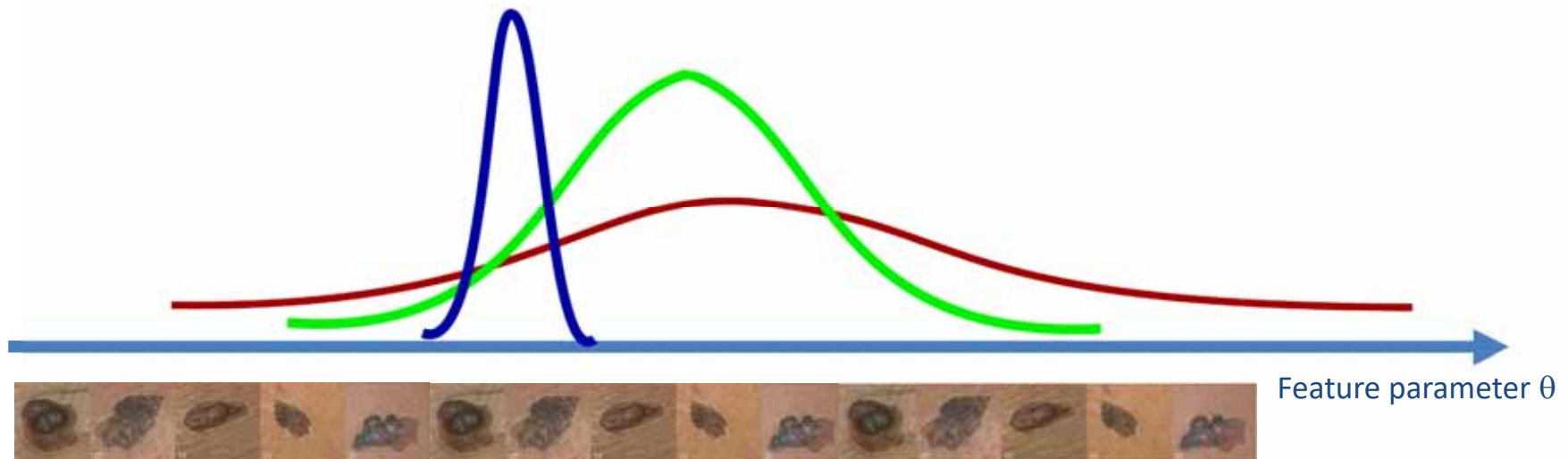
h ... hypotheses

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

Likelihood
Prior Probability

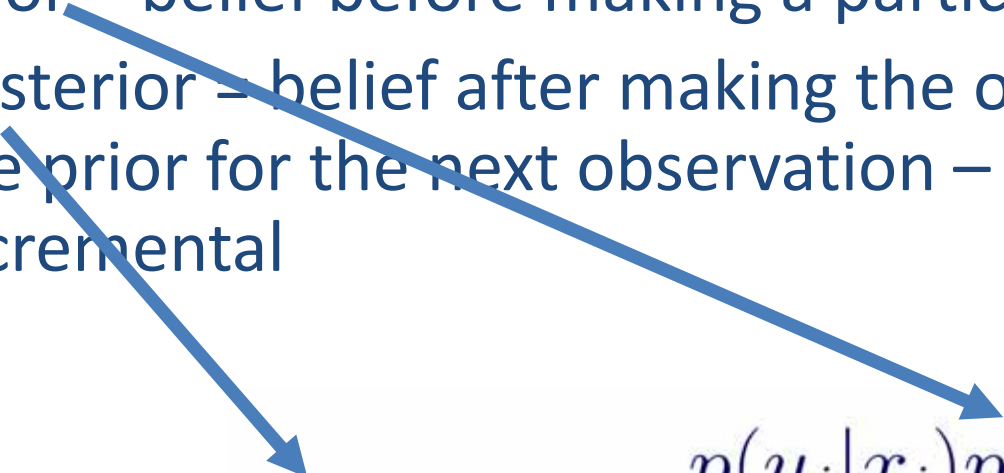
Posterior Probability

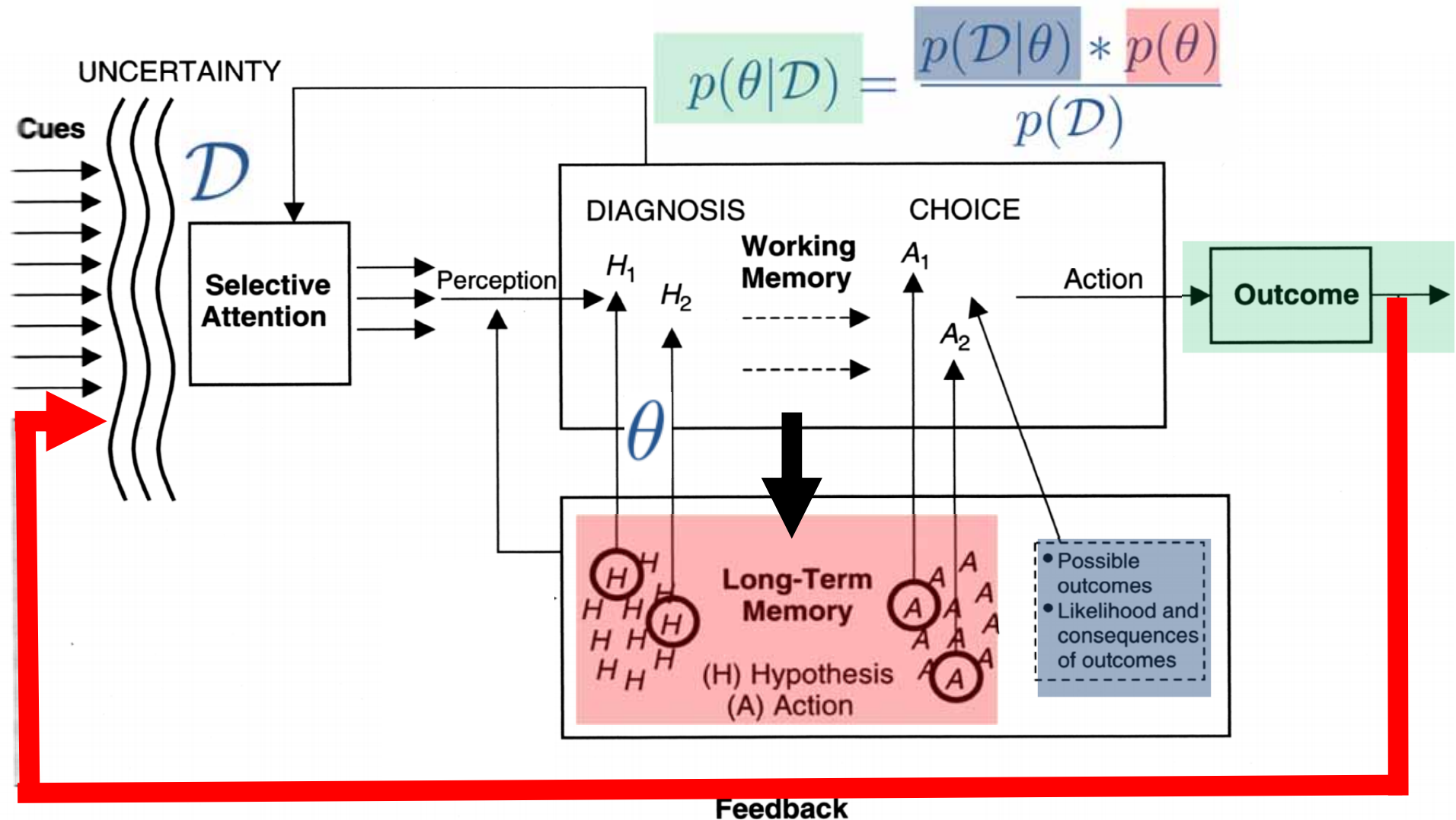
Problem in $\mathbb{R}^n \rightarrow$ complex



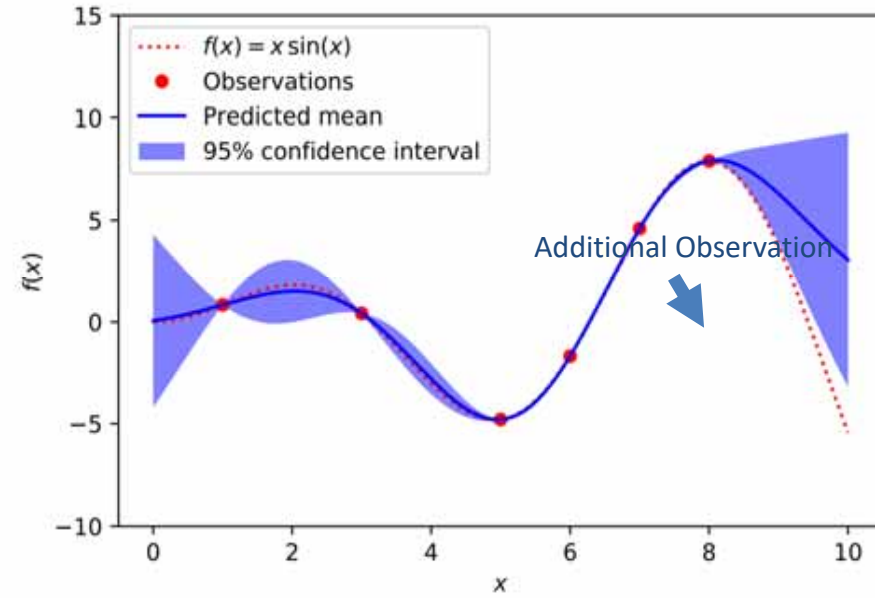
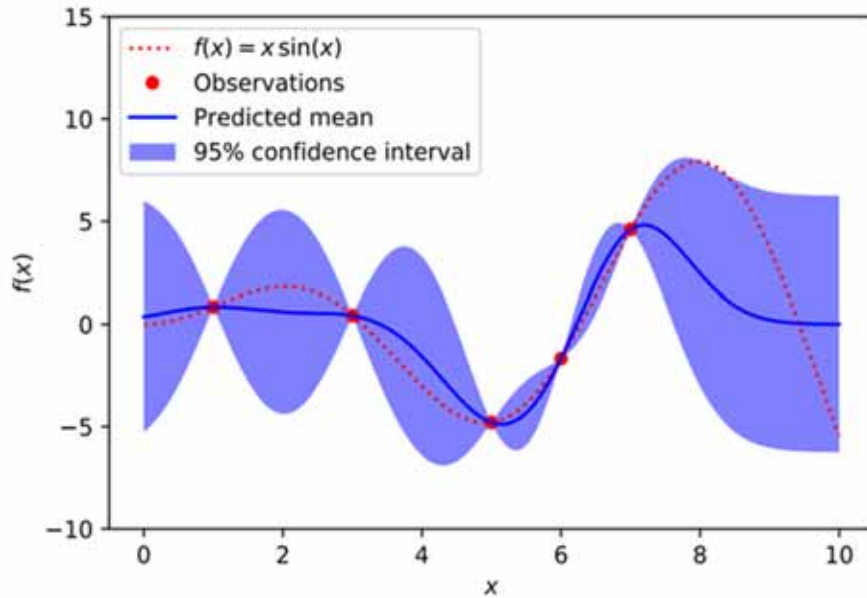
Why is this relevant for medicine?

- 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


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



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

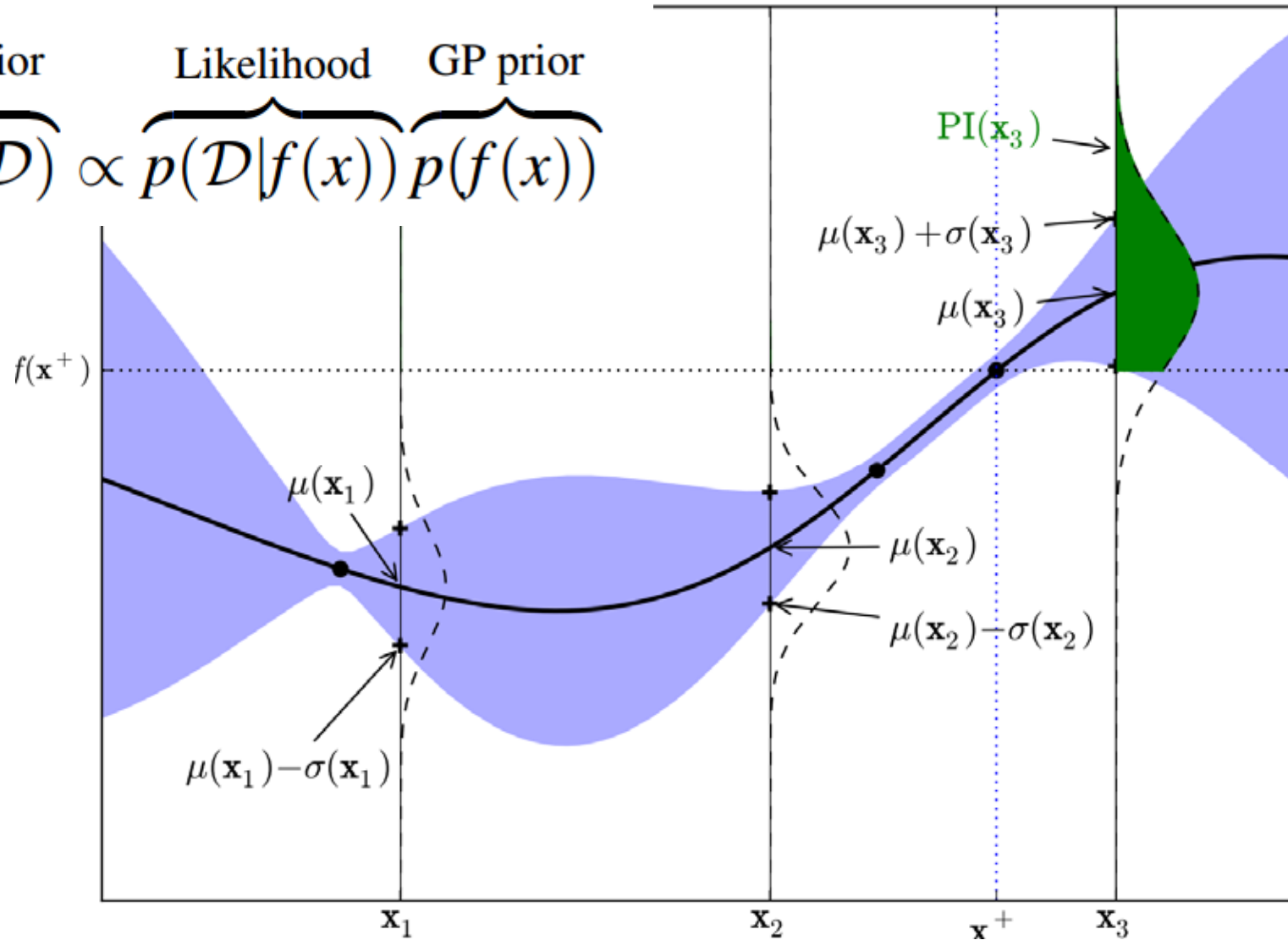


$$\mathbb{E}[f] = \int p(x) f(x) dx$$

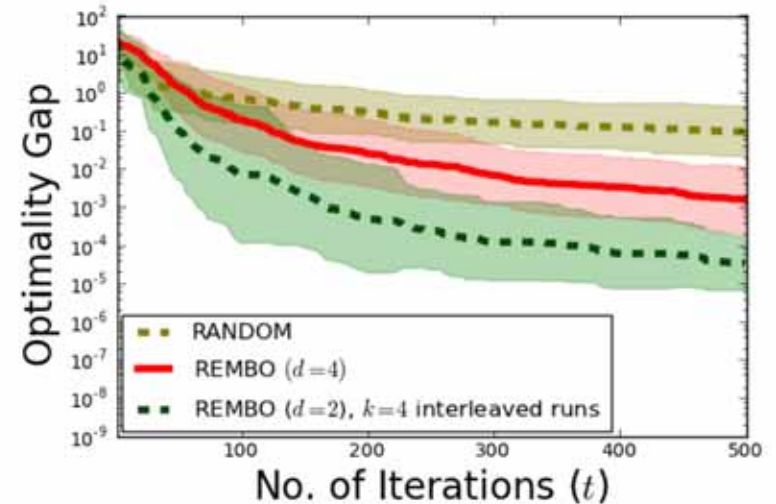
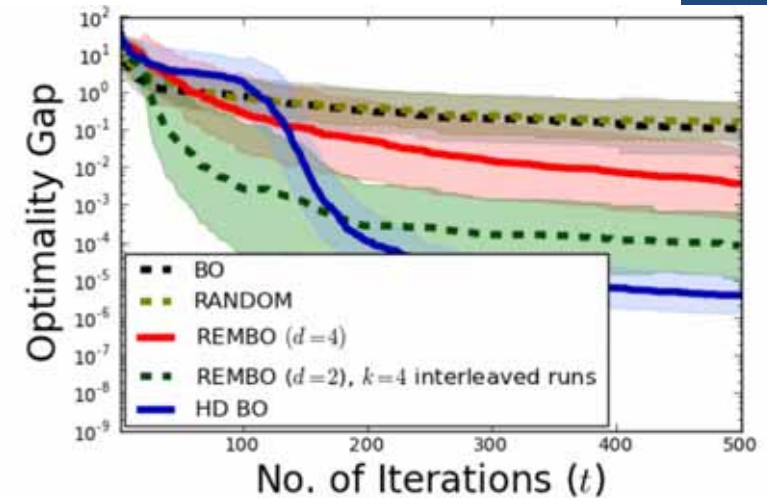
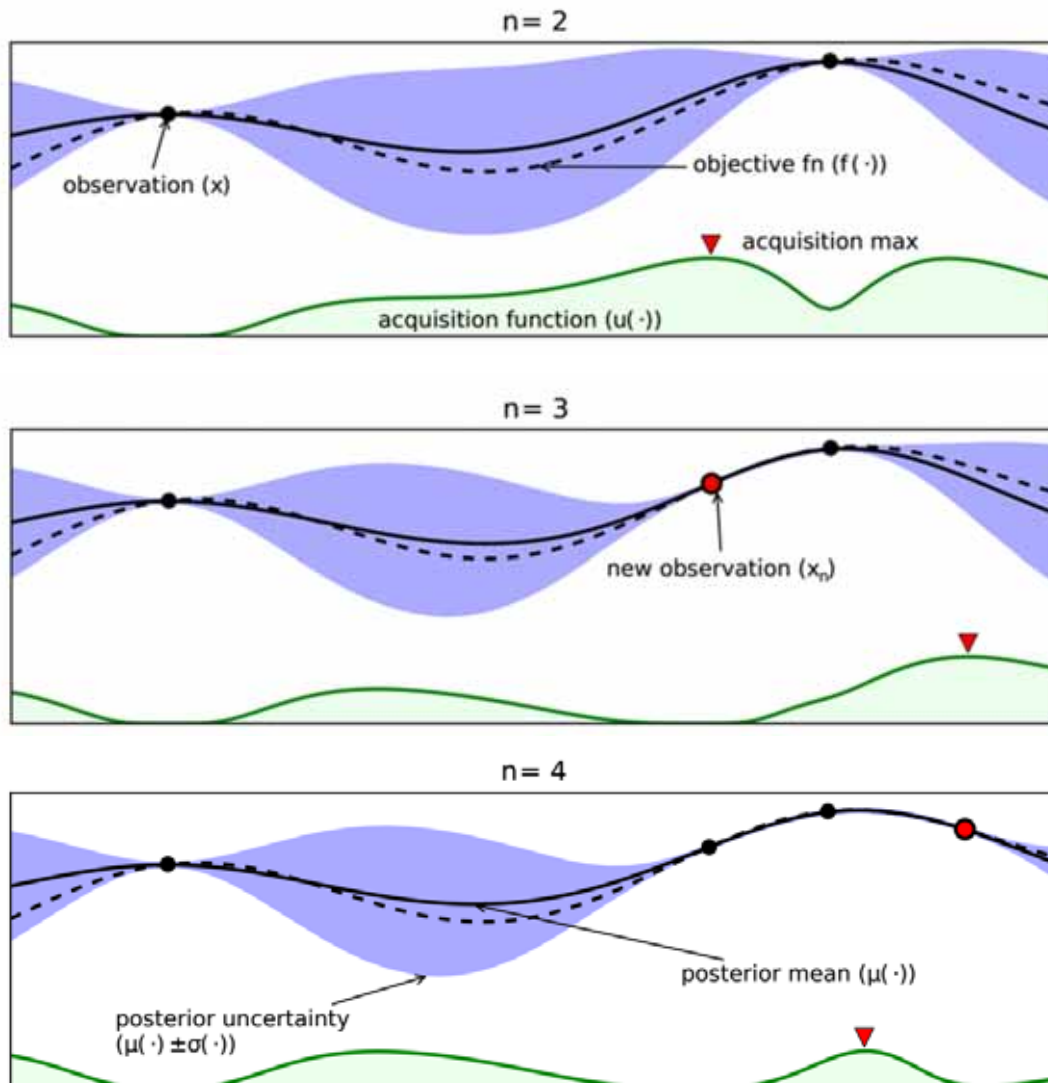
$$\mathbb{E}[f] \simeq \frac{1}{N} \sum_{n=1}^N f(x_n)$$

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

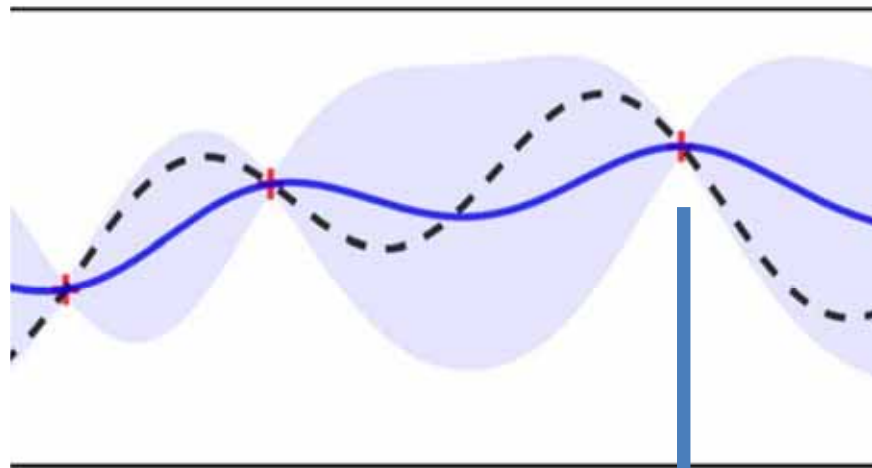
$$\overbrace{p(f(x)|\mathcal{D})}^{\text{GP posterior}} \propto \overbrace{p(\mathcal{D}|f(x))}^{\text{Likelihood}} \overbrace{p(f(x))}^{\text{GP prior}}$$



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.

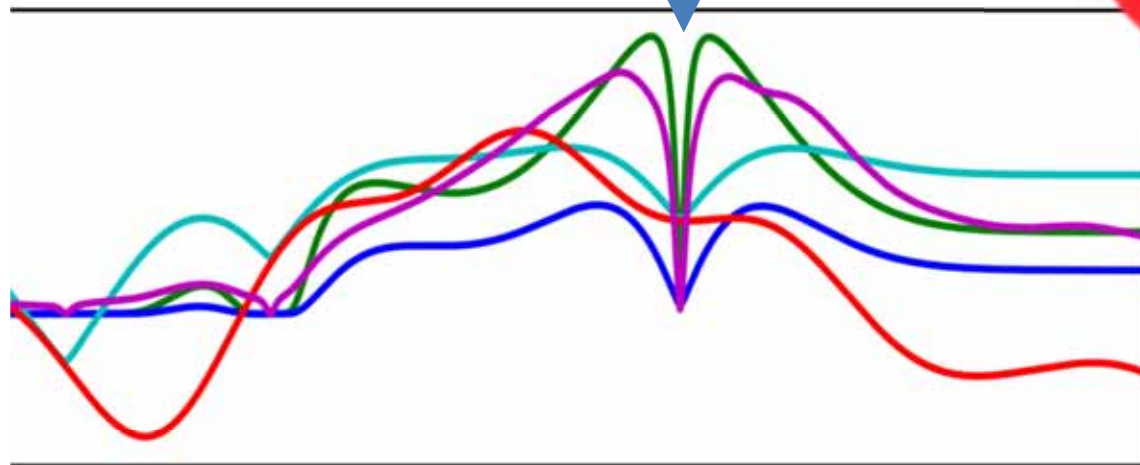
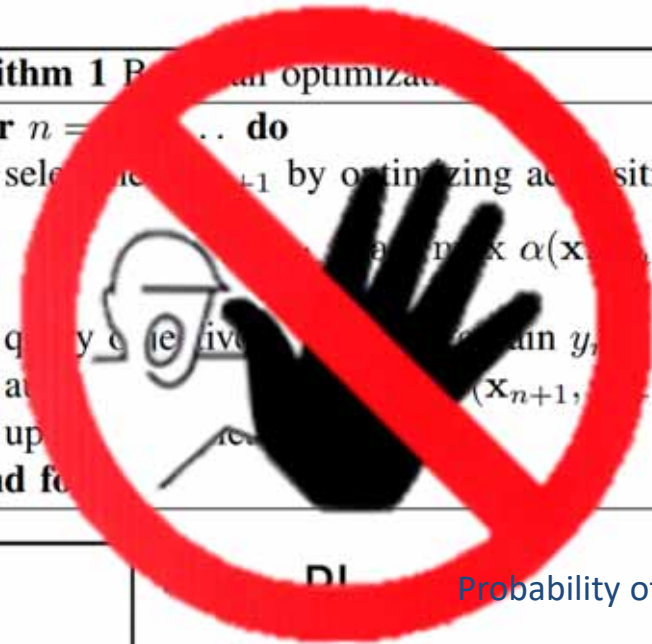


Wang, Z., Hutter, F., Zoghi, M., Matheson, D. & De Feitas, N. 2016. Bayesian optimization in a billion dimensions via random embeddings. Journal of Artificial Intelligence Research, 55, 361-387, doi:10.1613/jair.4806.



```

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

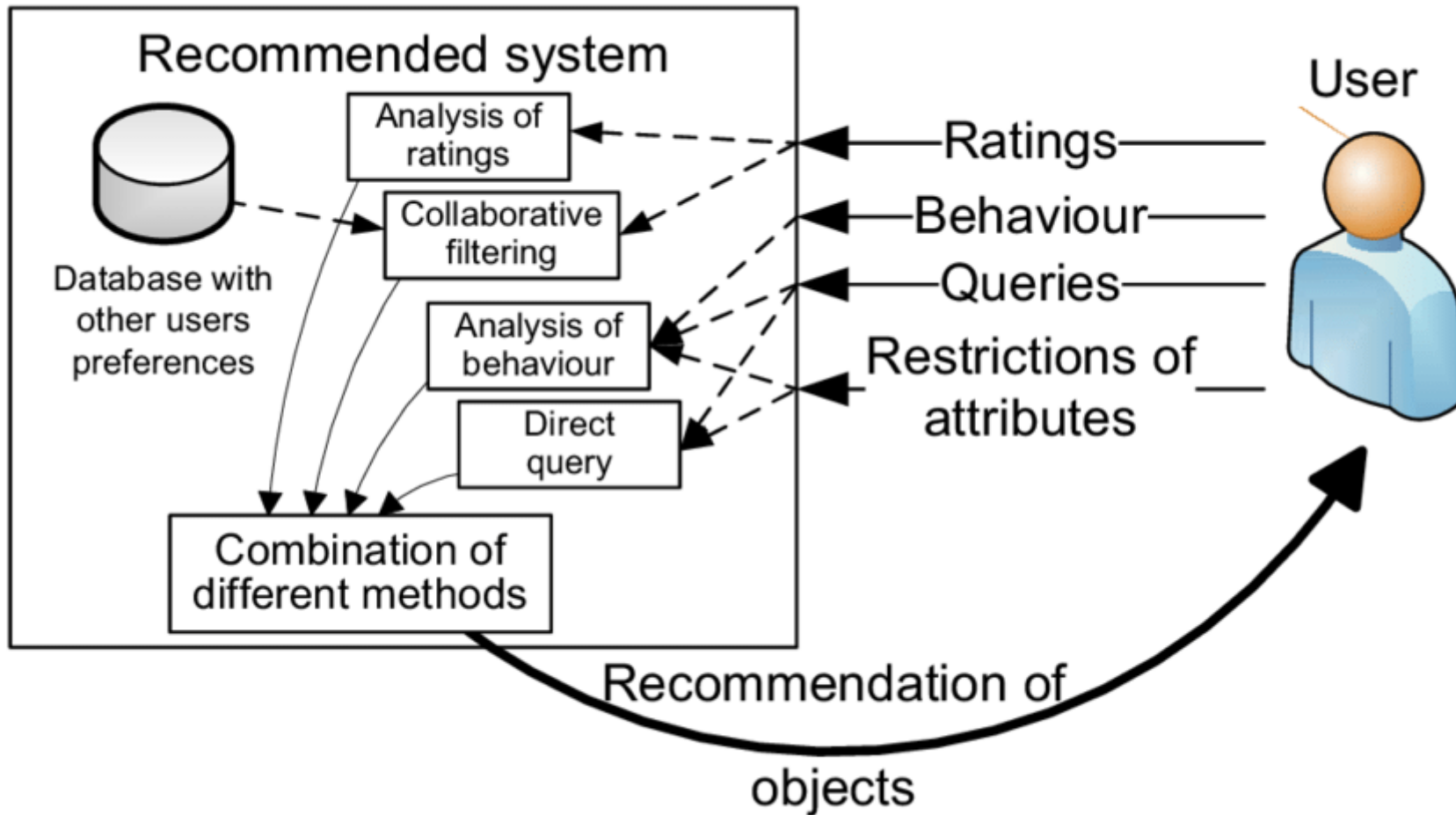


- EI Expected Improvement
- UCB Upper Confidence Bound
- TS Thompson Sampling
- PES Predictive Entropy Search

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.



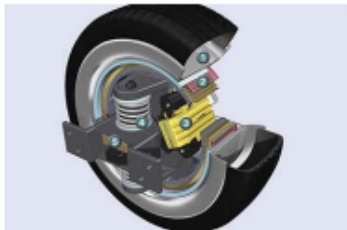
Human			Machine		
LEVEL 0	LEVEL 1	LEVEL 2	LEVEL 3	LEVEL 4	LEVEL 5
No Active Assistance System	Longitudinal or Transverse Guide	Traffic Control	Awareness for Take Over	No Driver Intervention	No Driver
	Longitudinal or Transverse Guide	Longitudinal and Transverse Guide	Take Over Request	No Take Over Request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



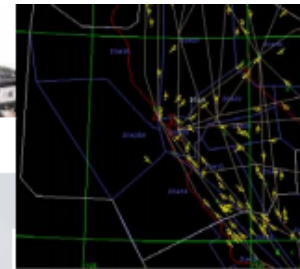
E-Corner, Siemens

Building Systems

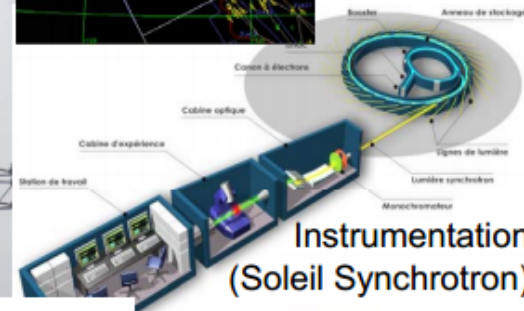


Avionics

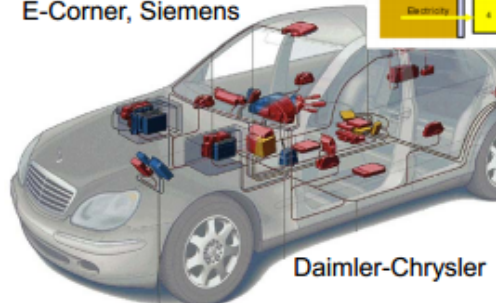
Telecommunications



Transportation
(Air traffic control at SFO)



Instrumentation
(Soleil Synchrotron)



Daimler-Chrysler

Power generation and distribution



Courtesy of General Electric

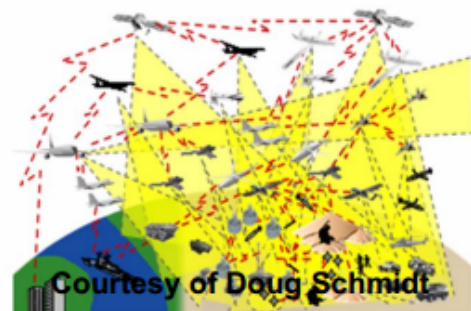
Factory automation



Courtesy of Kuka Robotics Corp.



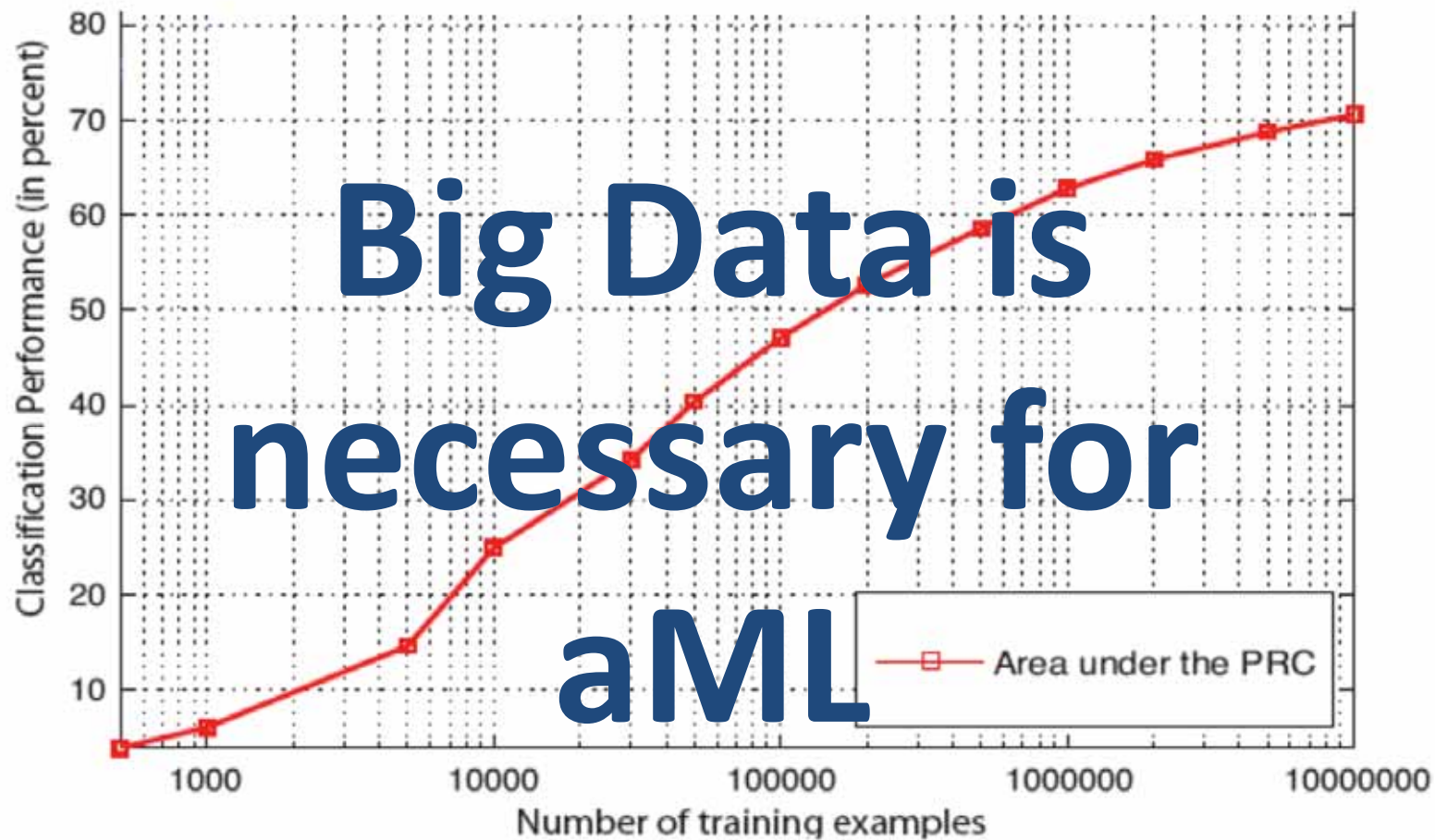
Military systems:



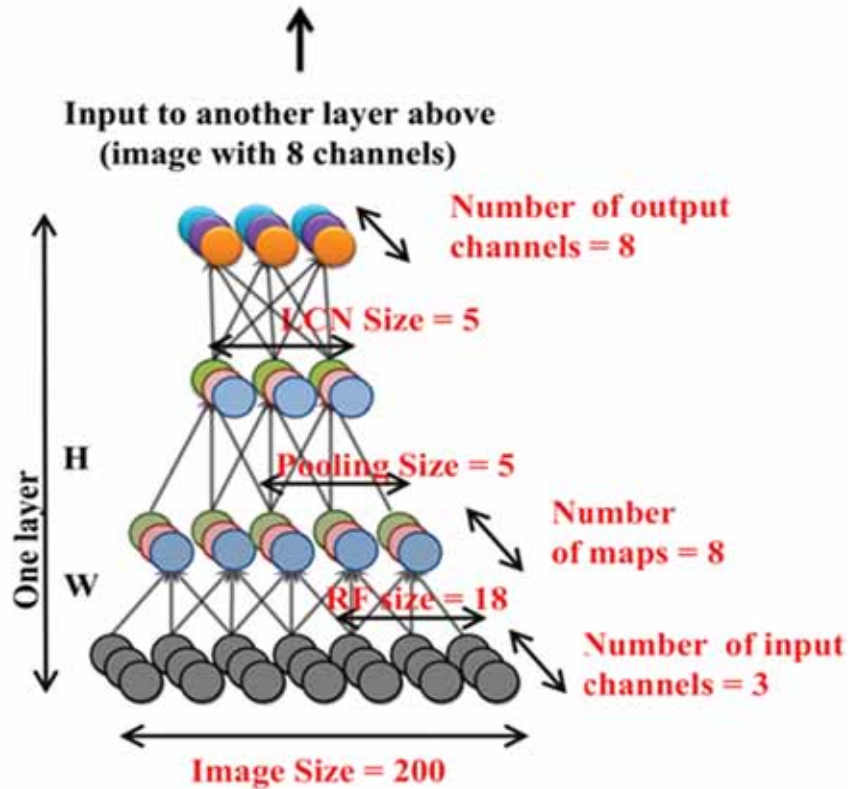
Courtesy of Doug Schmidt



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, ...



05 iML

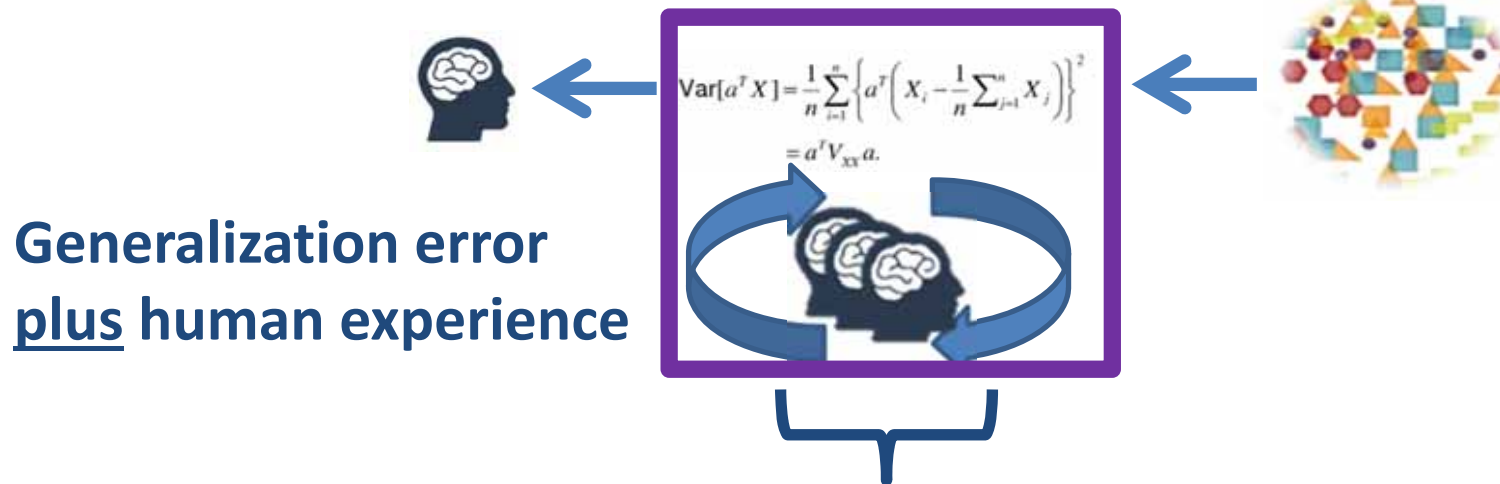
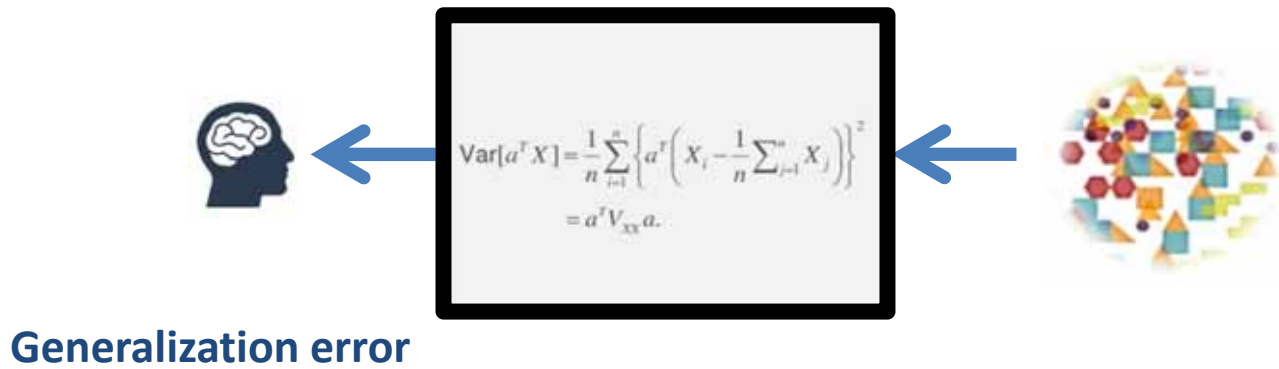
Holzinger, A. 2016. Interactive Machine Learning (iML). Informatik Spektrum, 39, (1), 64-68, doi:10.1007/s00287-015-0941-6.

- **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.



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 intuition?

Humans can generalize from few examples, and ...

- understand relevant representations,
- find concepts between $P(x)$ and $P(Y|X)$,
- with a causal link between $Y \rightarrow X$

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)



See also: Ian J. Goodfellow, Jonathon Shlens & Christian Szegedy 2014. Explaining and harnessing adversarial examples. arXiv:1412.6572. Images see: <https://imgur.com/a/K4RWn>

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.



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.

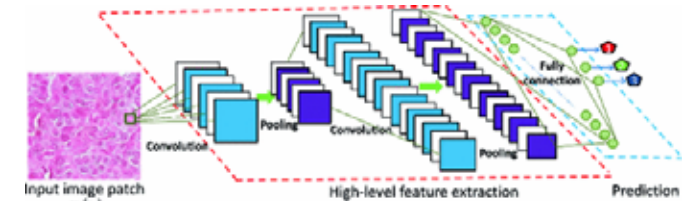
Image Captions by dee learning : github.com/karpathy/neuraltalk2

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

06 Methods of Explainable AI

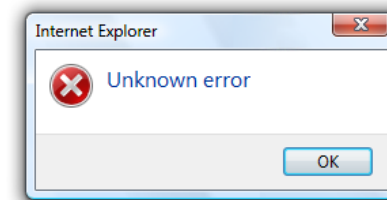
Verify that algorithms/classifiers work as expected ...

Wrong decisions can be costly and dangerous ...



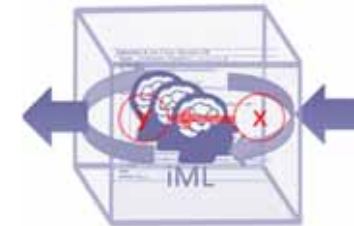
Understanding the errors ...

Detection of bias, weaknesses, unknowns, ...



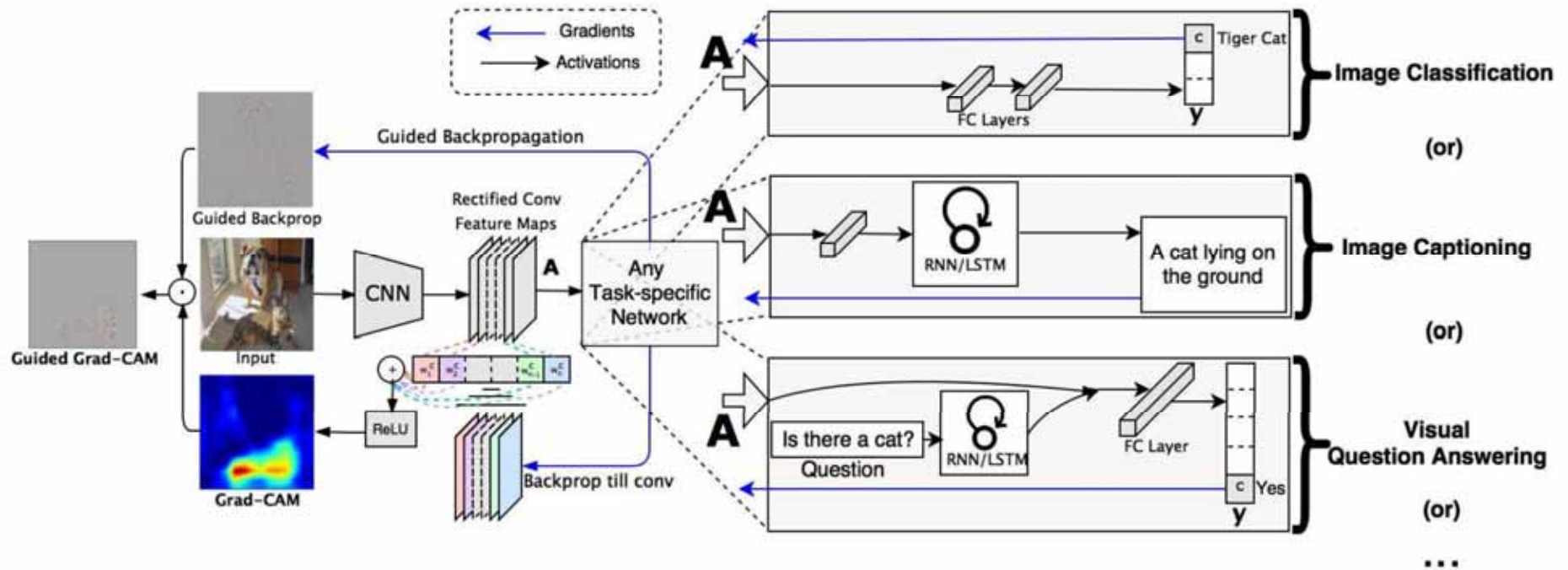
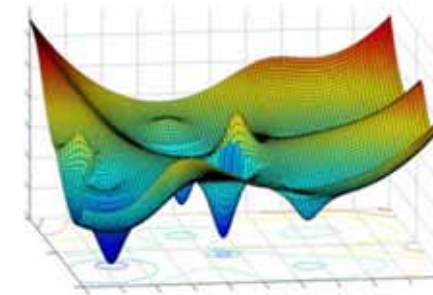
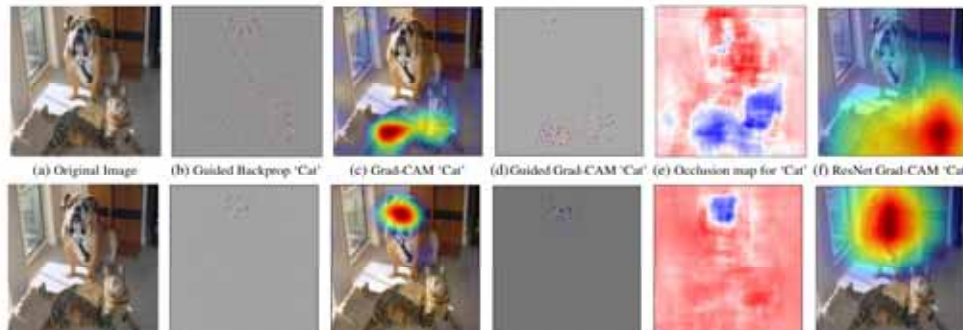
Scientific replicability and causality ...

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

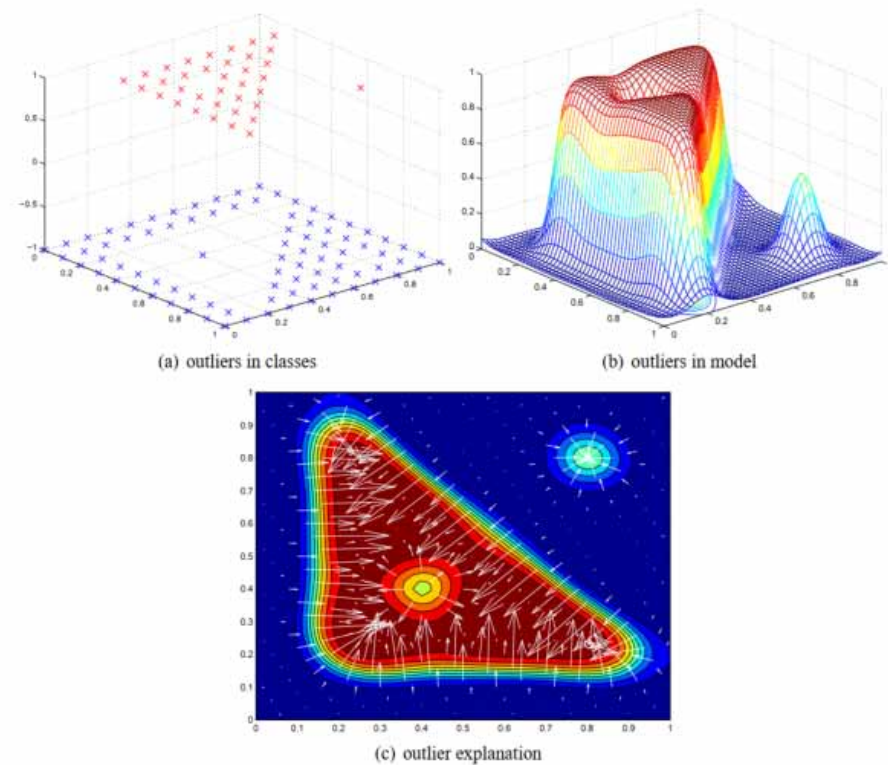
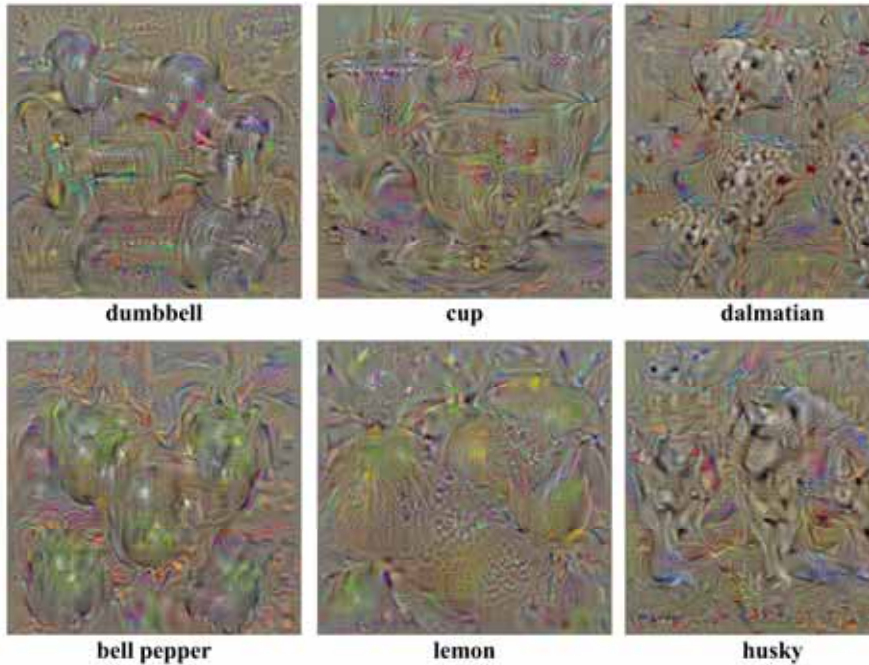


- 1) Gradients
- 2) Sensitivity Analysis
- 3) Decomposition Relevance Propagation (Pixel-RP, Layer-RP, Deep Taylor Decomposition, ...)
- 4) Optimization (Local-IME – model agnostic, BETA transparent approximation, ...)
- 5) Deconvolution and Guided Backpropagation
- 6) Model Understanding
 - Feature 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)



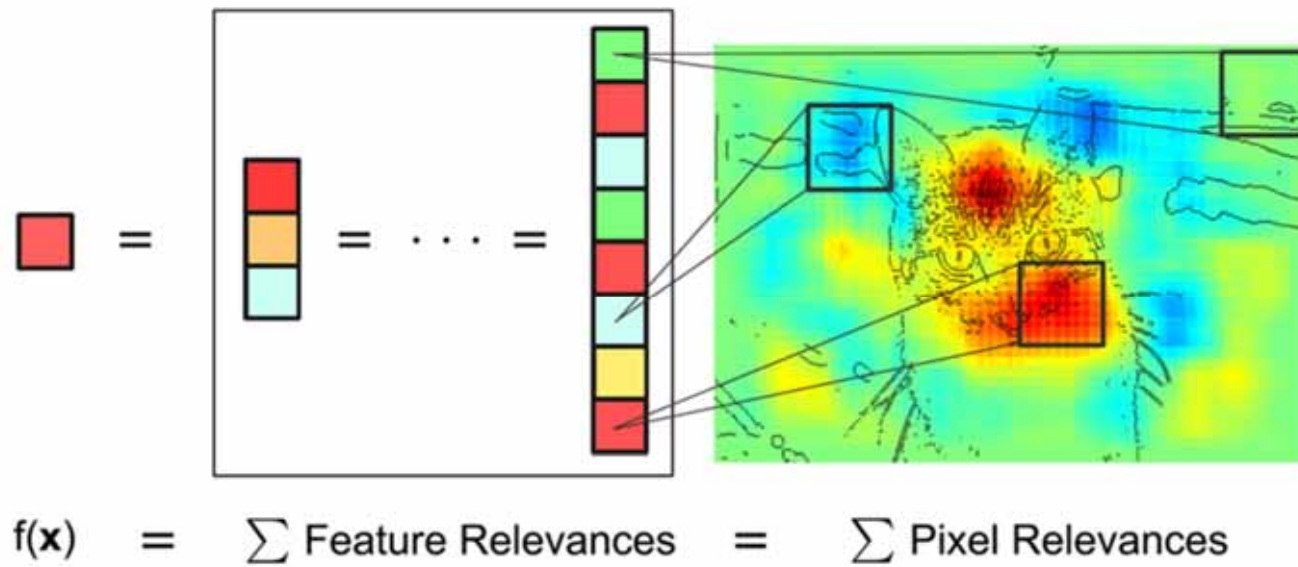
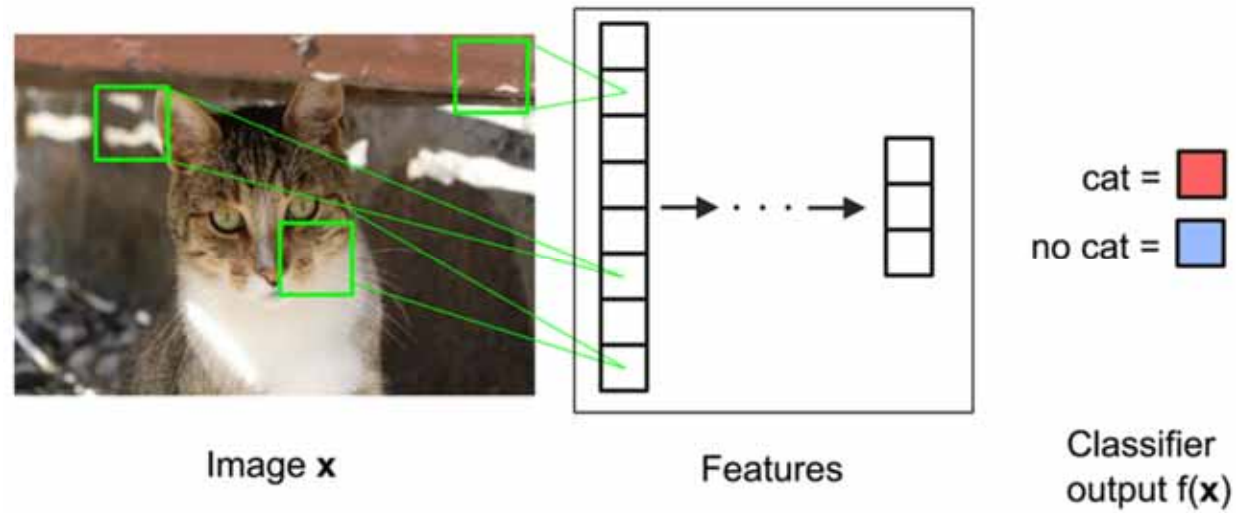
Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh & Dhruv Batra. Grad-CAM: Visual Explanations from Deep Networks via Gradient-Based Localization. ICCV, 2017. 618-626.

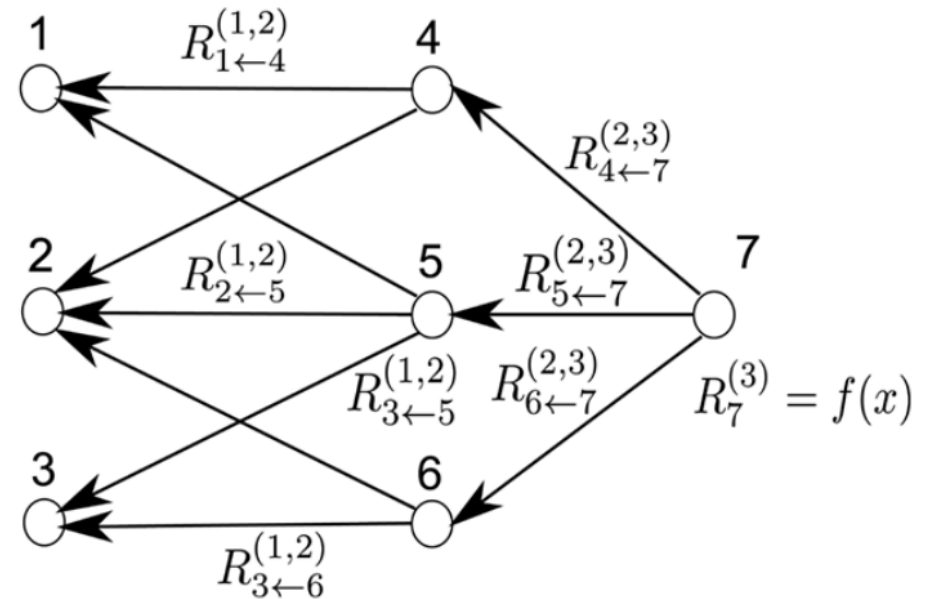
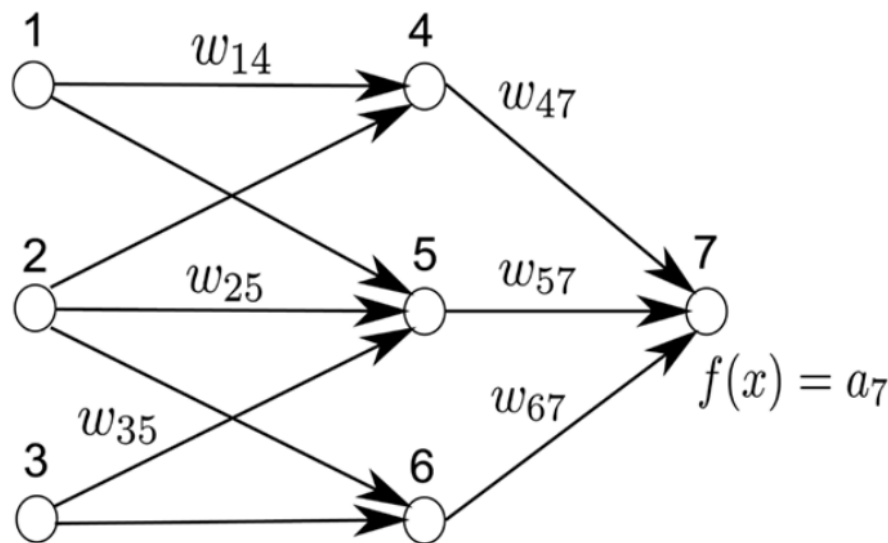


Karen Simonyan, Andrea Vedaldi & Andrew Zisserman 2013. Deep inside convolutional networks: Visualising image classification models and saliency maps. arXiv:1312.6034.

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.

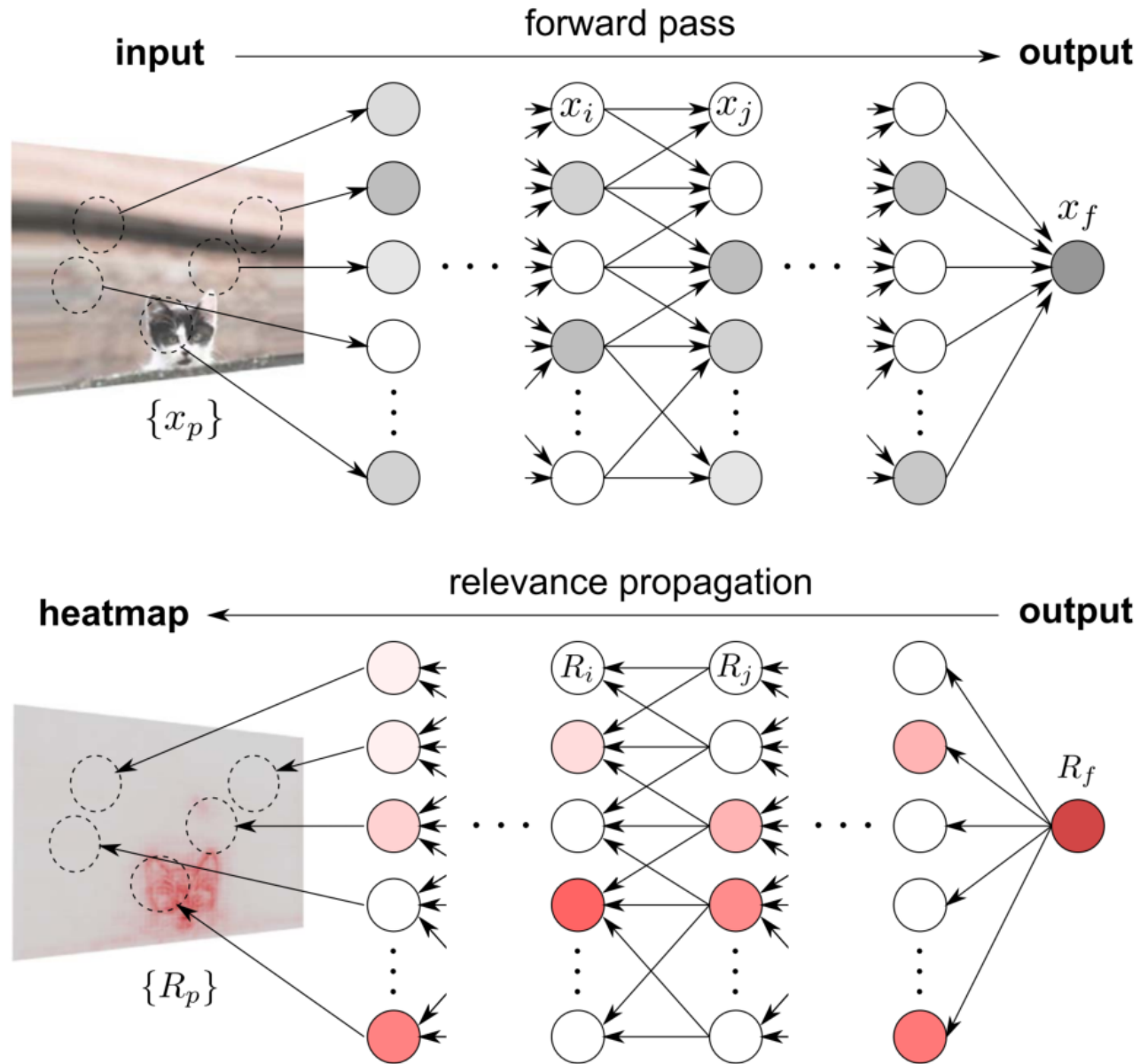
Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller & Wojciech Samek 2015. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. PLoS one, 10, (7), e0130140, doi:10.1371/journal.pone.0130140.

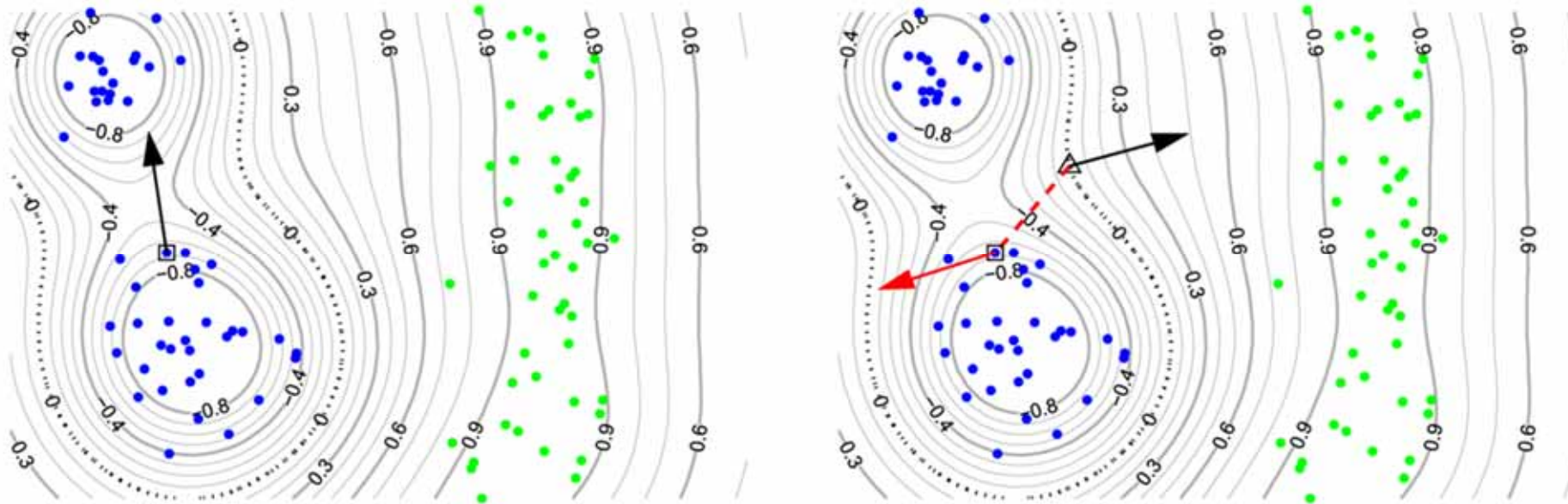




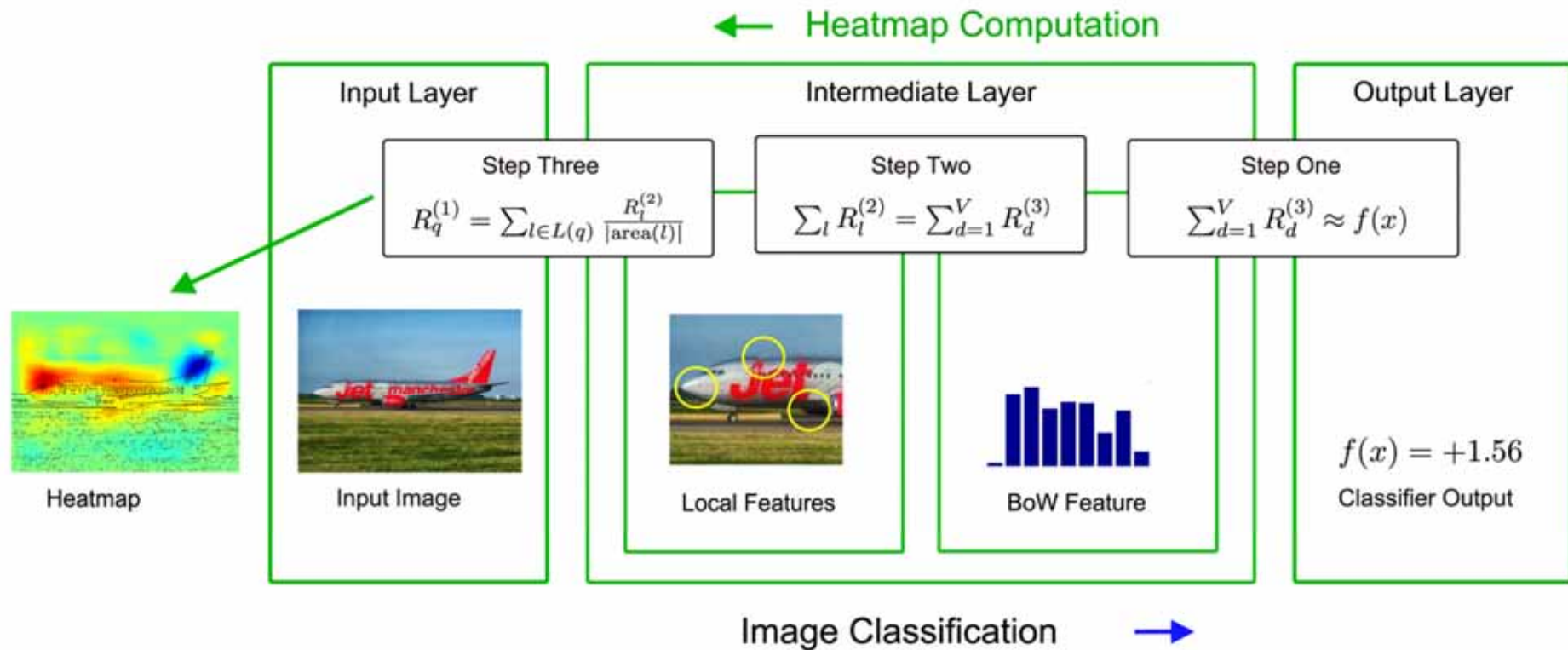
$$f(x) = \dots = \sum_{d \in l+1} R_d^{(l+1)} = \sum_{d \in l} R_d^{(l)} = \dots = \sum_d R_d^{(1)}$$

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.

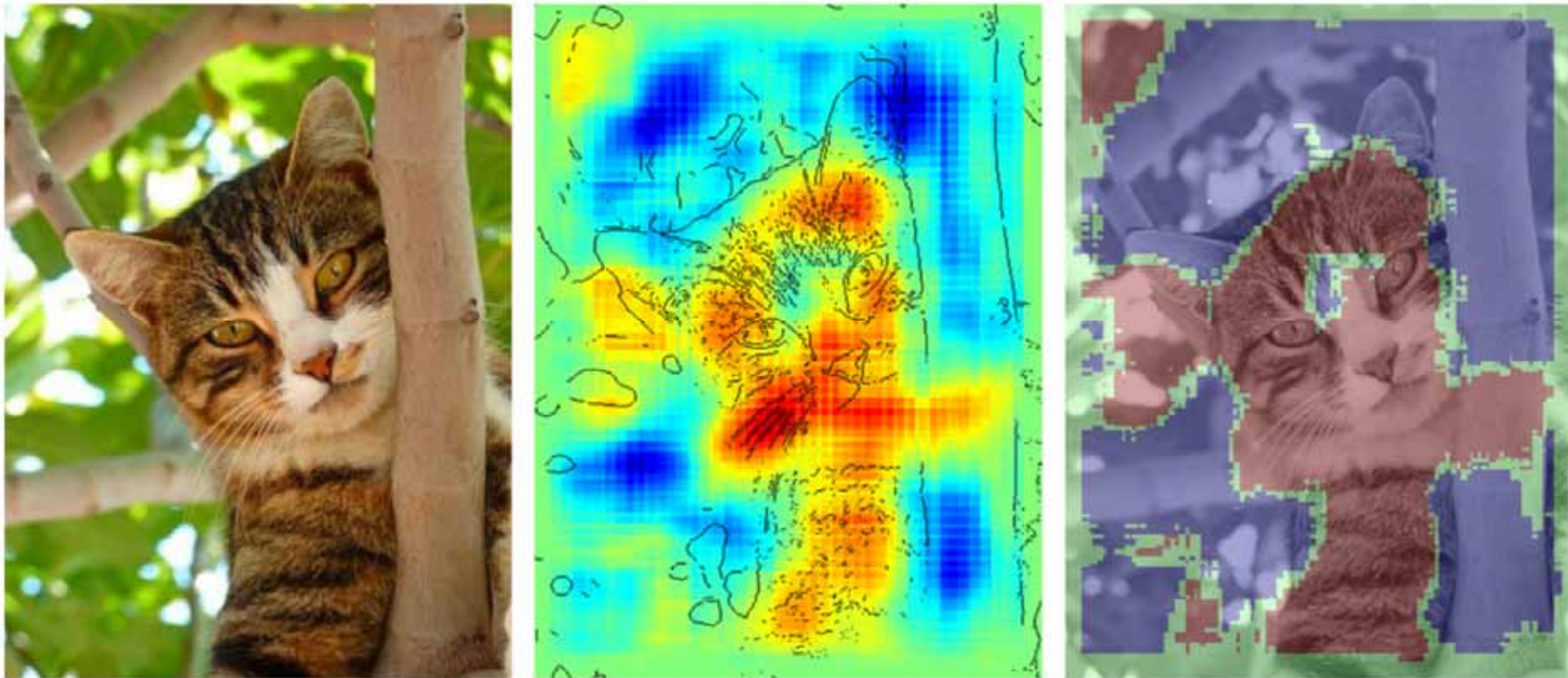




Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller & Wojciech Samek
2015. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. PloS one, 10,
(7), e0130140, doi:10.1371/journal.pone.0130140.

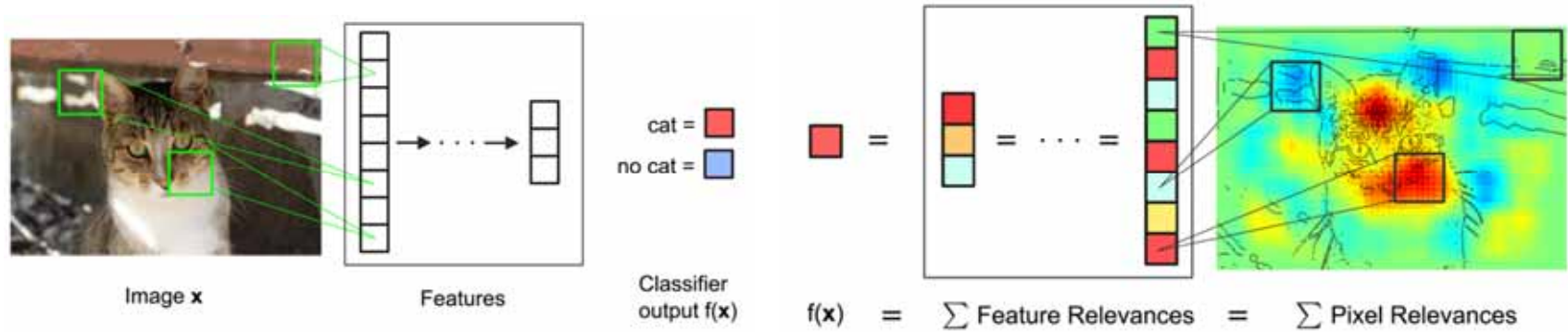


Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller & Wojciech Samek
 2015. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. PLoS one, 10,
 (7), e0130140, doi:10.1371/journal.pone.0130140.



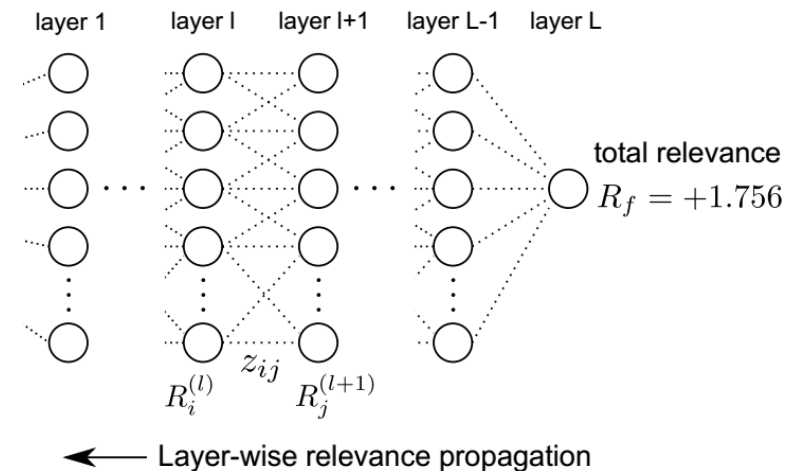
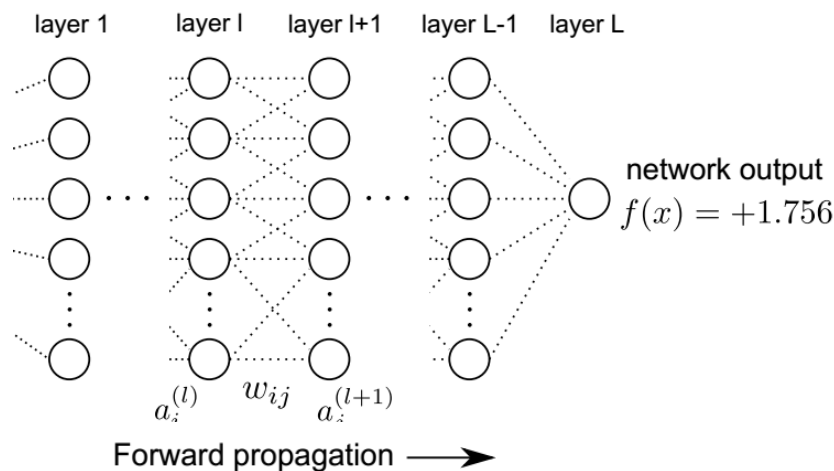
Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller & Wojciech Samek
2015. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. PloS one, 10,
(7), e0130140, doi:10.1371/journal.pone.0130140.

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.

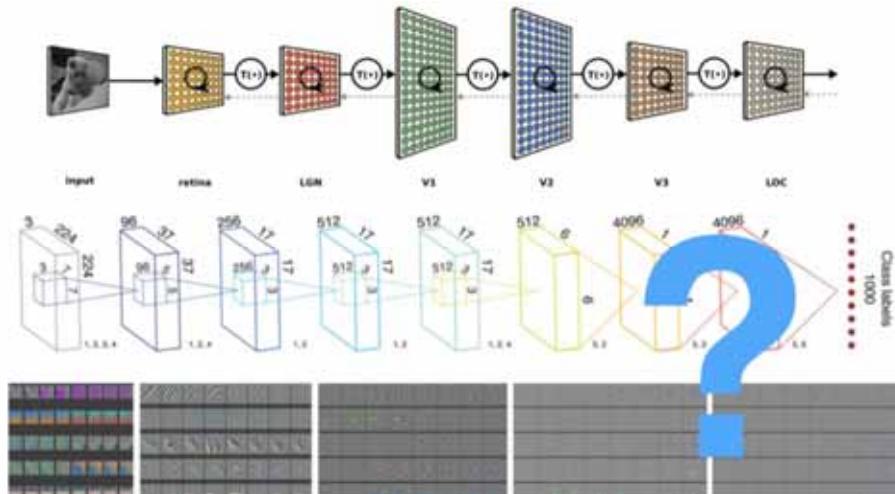


$$a_j^{(l+1)} = \sigma \left(\sum_i a_i^{(l)} w_{ij} + b_j^{(l+1)} \right)$$

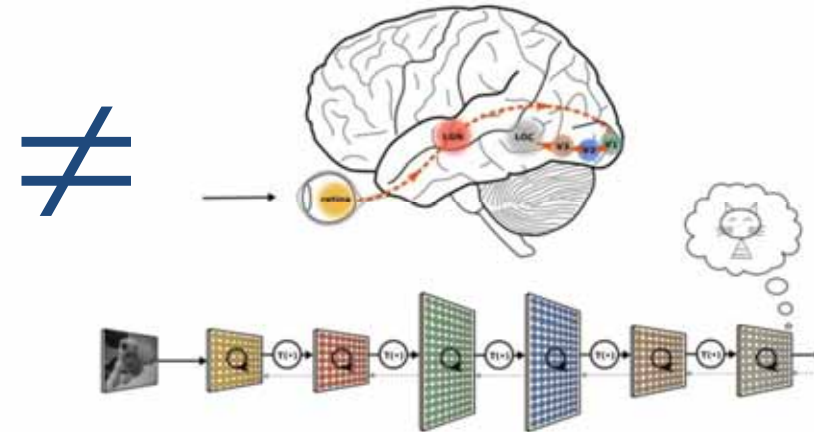
$$R_i^{(l)} = \sum_j \frac{z_{ij}}{\sum_{i'} z_{i'j}} R_j^{(l+1)}$$



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



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

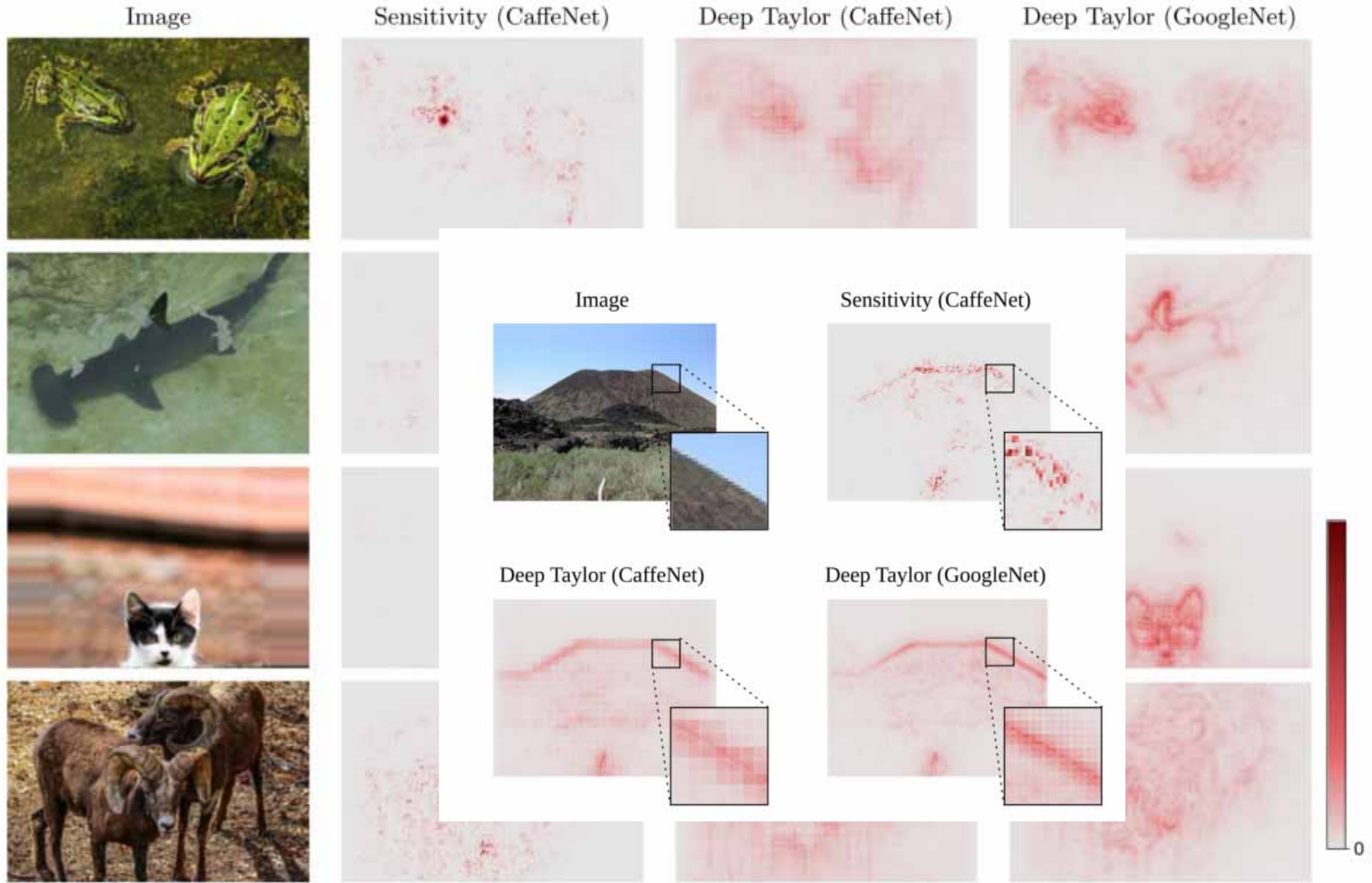


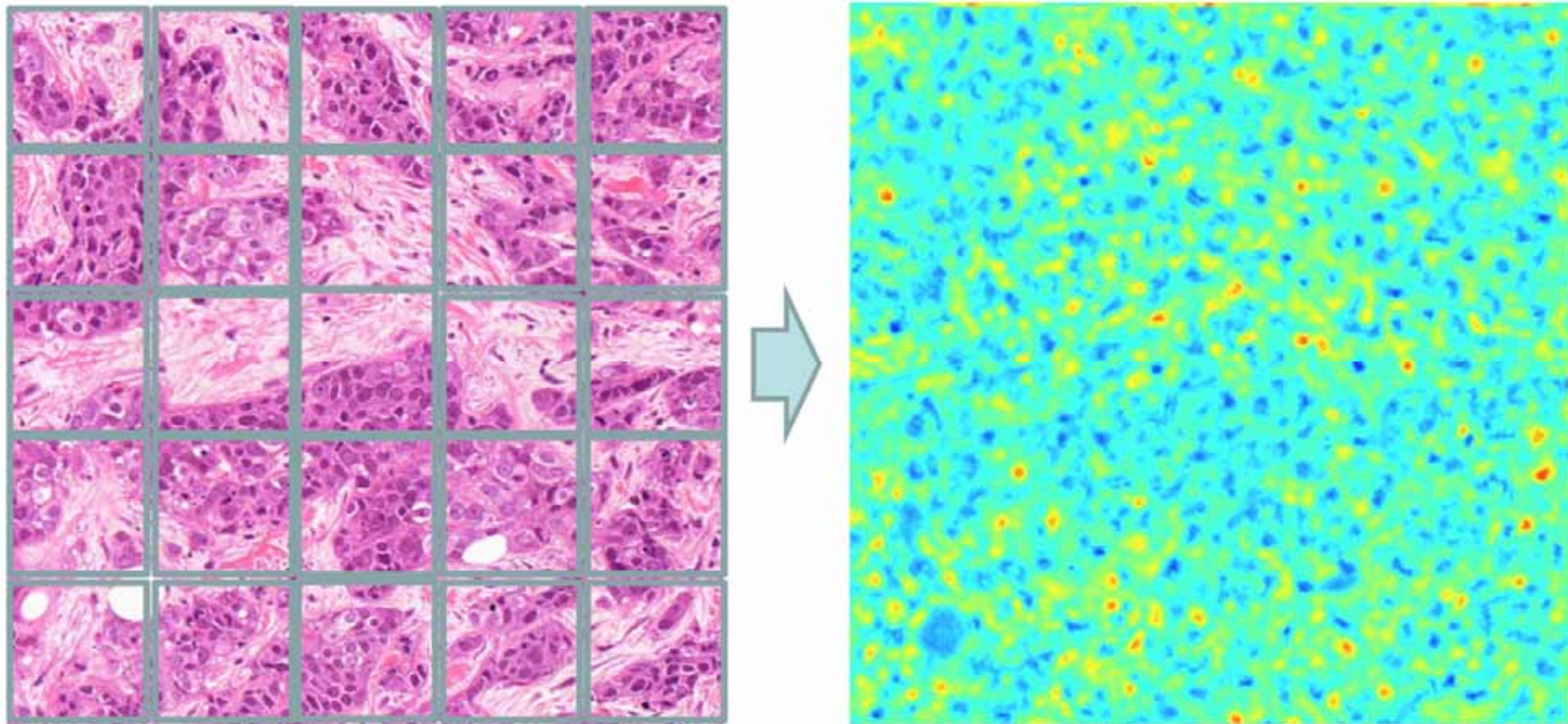
$$\frac{\partial h_k(x)}{\partial x_{a,b}}$$

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.

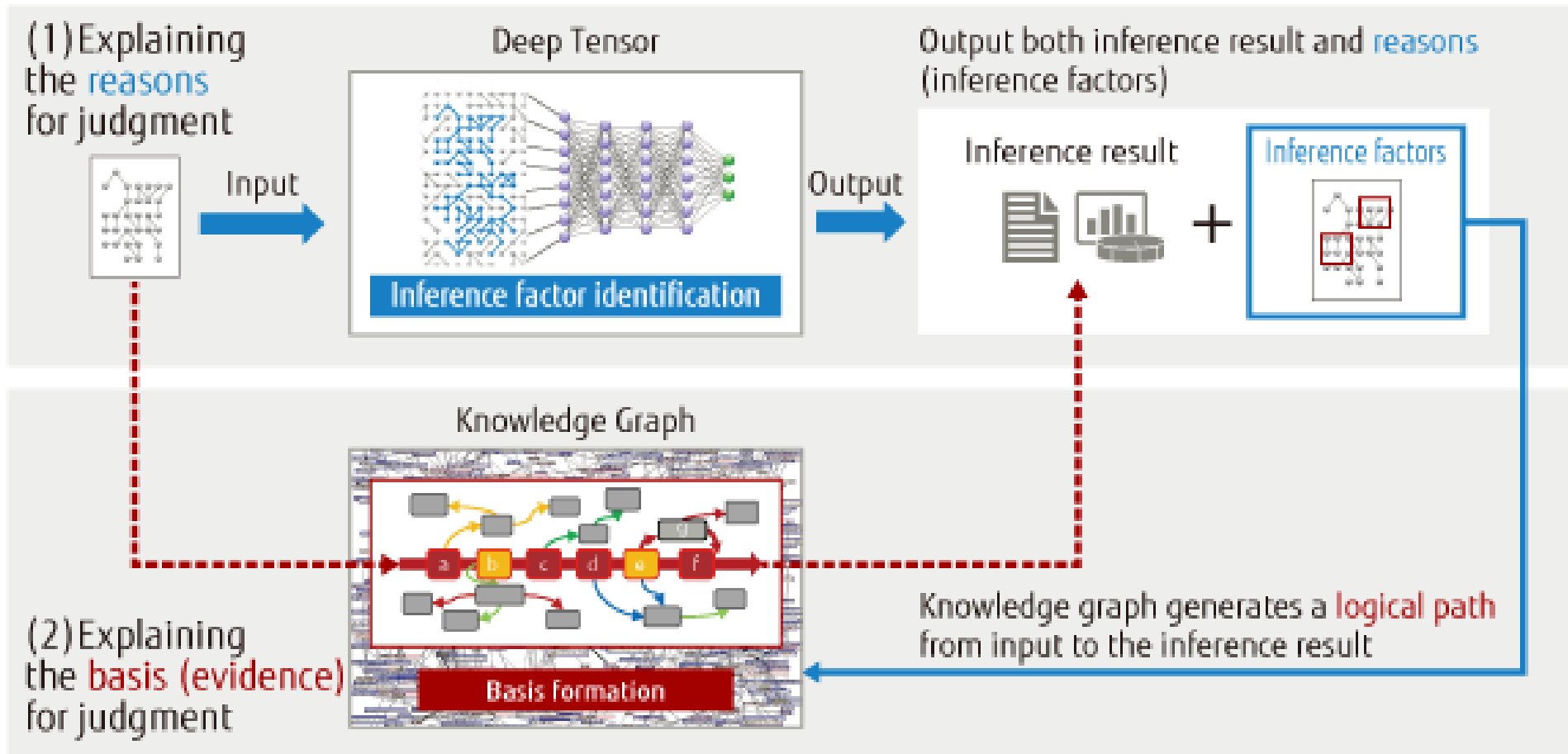




Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller & Wojciech Samek 2015. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. PloS one, 10, (7), e0130140, doi:10.1371/journal.pone.0130140.

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!**

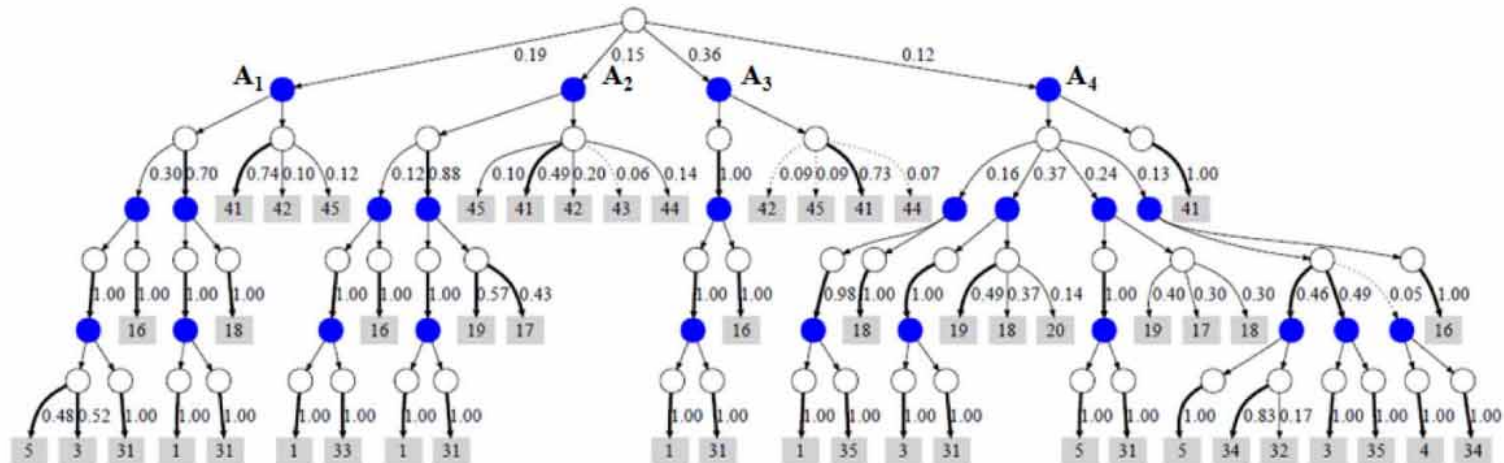


Randy Goebel, Ajay Chander, Katharina Holzinger, Freddy Lecue, Zeynep Akata, Simone Stumpf, Peter Kieseberg & Andreas Holzinger 2018. Explainable AI: the new 42? Springer Lecture Notes in Computer Science LNCS 11015

Input images



Stochastic AOT



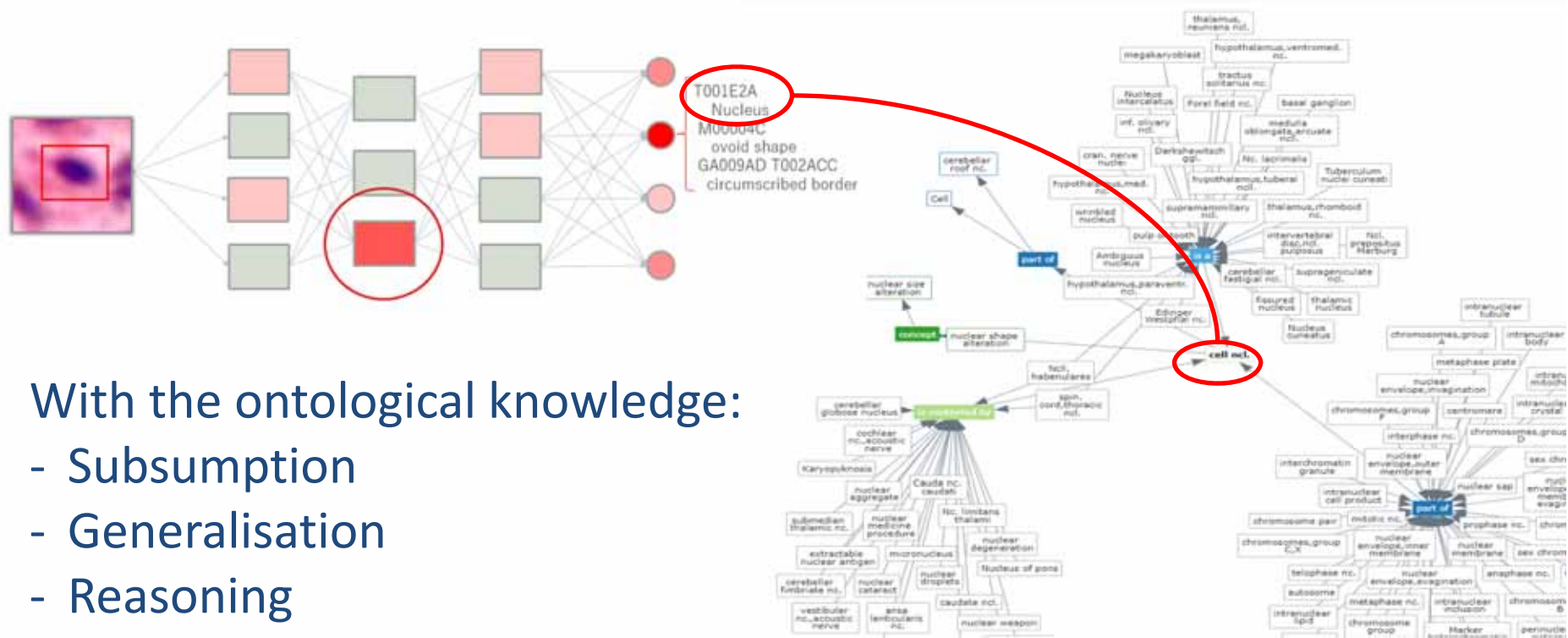
Part dictionary (terminal nodes)

	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.

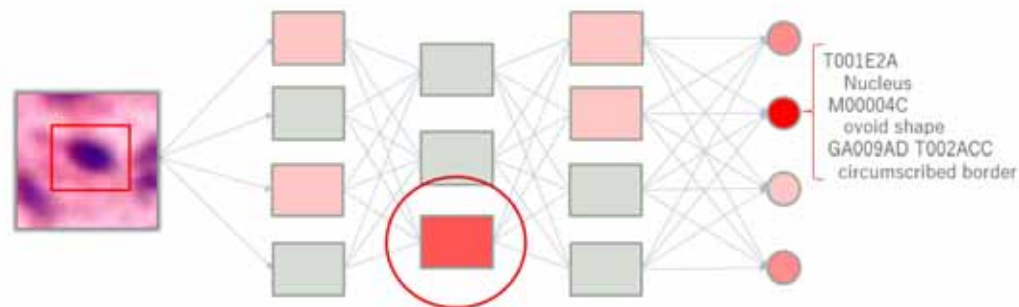


With the ontological knowledge:

- Subsumption
- Generalisation
- Reasoning
- Validation
- Context !

Image Source: unpublished, from our current Project with ID-Berlin

Which nodes have which relevant task?
Which semantic categories have been used at all?

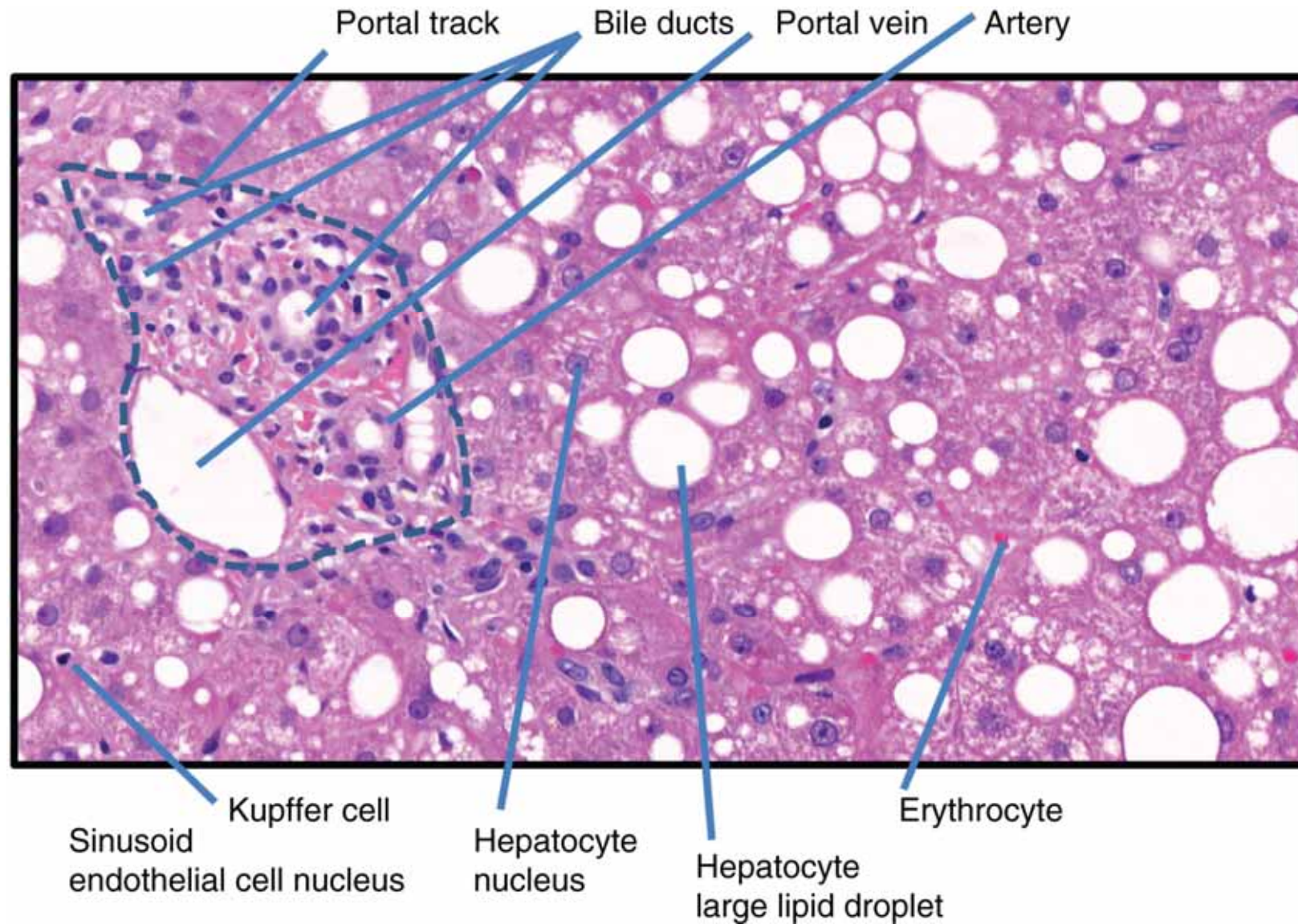


Moreover, ontologies can be used to

- analyze the internal structure
- describe the network
- filter nodes

Image Source: unpublished, from our current Project with ID-Berlin

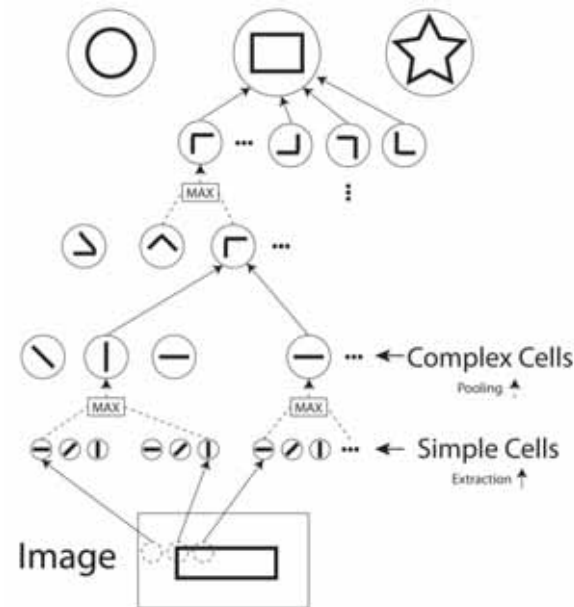
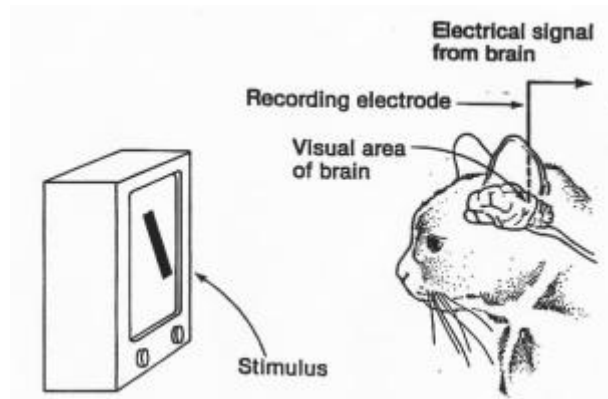
08 Measuring Machine Intelligence KANDINSKY-Project



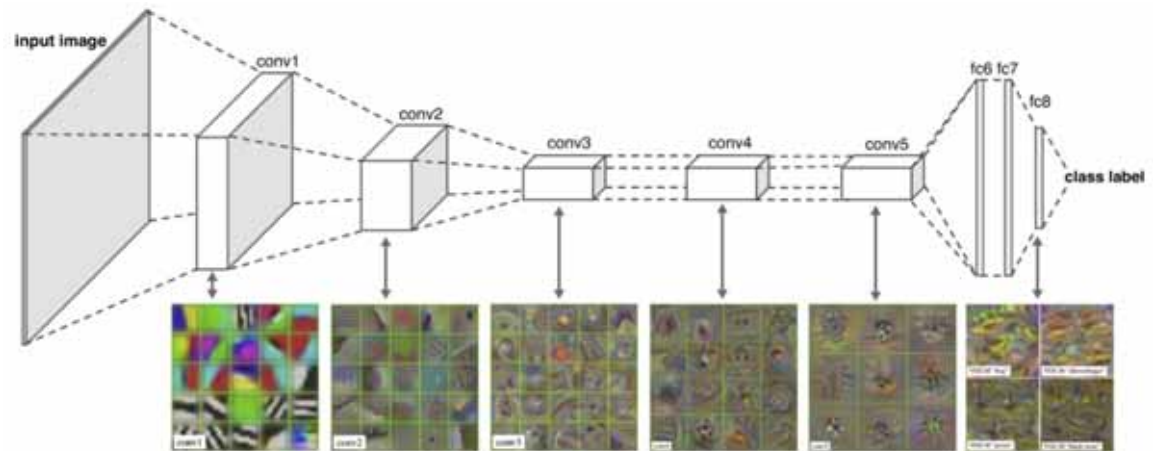
Andreas Holzinger, Georg Langs, Helmut Denk, Kurt Zatloukal & Heimo Mueller 2019. Causability and Explainability of AI in Medicine. Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery, doi:10.1002/widm.1312.

Radiologischer Befund		angelegt am 06.05.2006/20:26 geschr. von [REDACTED] gedruckt am 17.11.2006/08:24 Anfo: NCHIN
Kurzanamnese:	St.p. SHT	
Fragestellung:	-	
Untersuchung:	Thorax eine Ebene liegend [REDACTED]	
SB		
Bewegungsartefakte. Zustand nach Schädelhirntrauma.		
Das Cor in der Größennorm, keine akuten Stauungszeichen. 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. Der re. Rezessus frei.		
Mit kollegialen Grüßen		
[REDACTED]		
*** Elektronische Freigabe durch [REDACTED] am 09.05.2006 ***		

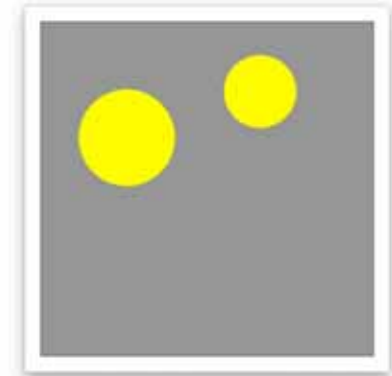
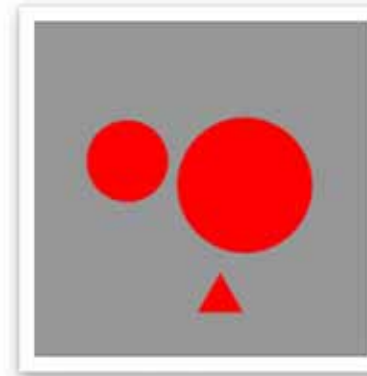
Holzinger, A., Geierhofer, R. & Errath, M. 2007. Semantische Informationsextraktion in medizinischen Informationssystemen. *Informatik Spektrum*, 30, (2), 69-78.



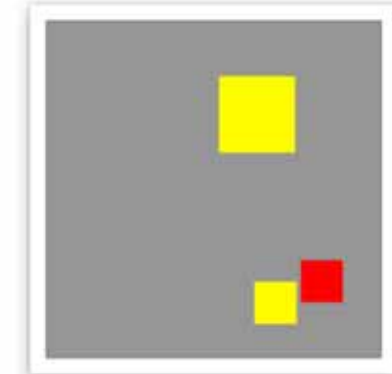
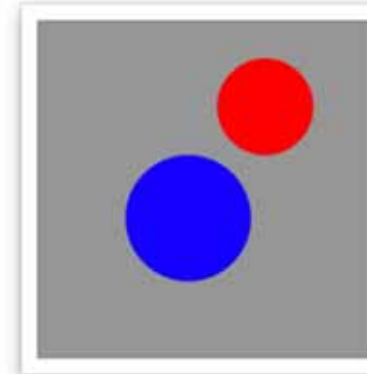
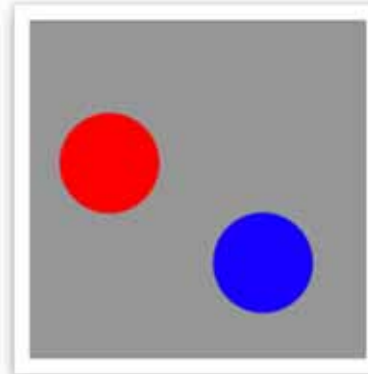
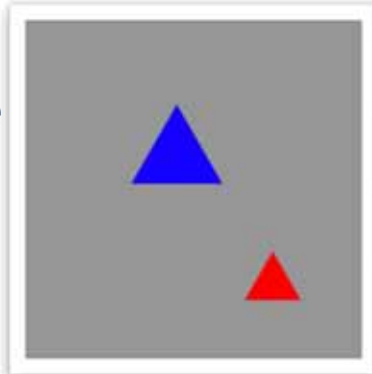
David H. Hubel & Torsten N. Wiesel
 1962. Receptive fields, binocular
 interaction and functional
 architecture in the cat's visual cortex.
 The Journal of Physiology, 160, (1),
 106-154,
 doi:10.1113/jphysiol.1962.sp006837.



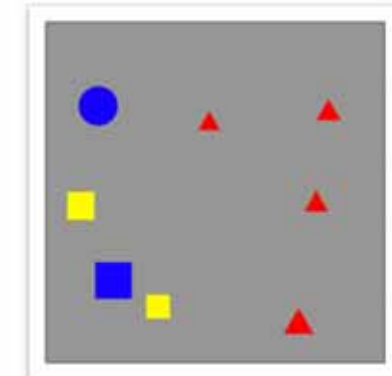
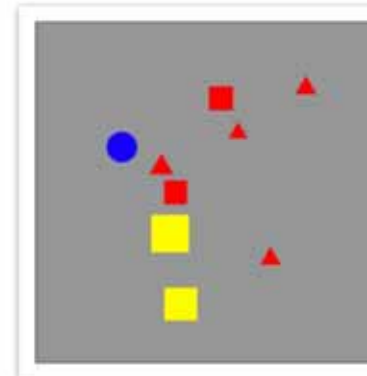
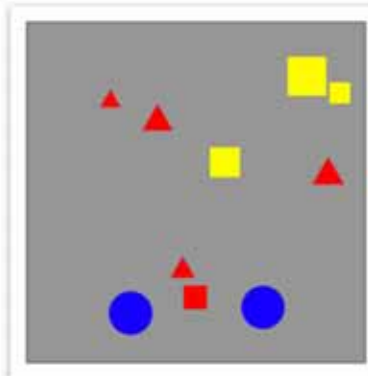
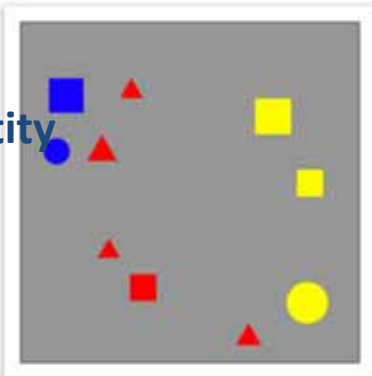
A
Colour



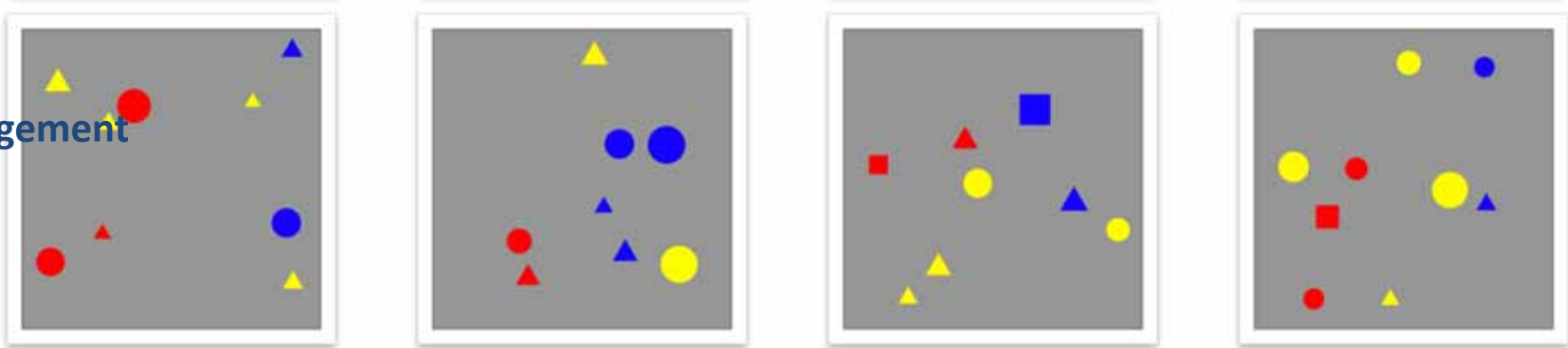
B
Shape



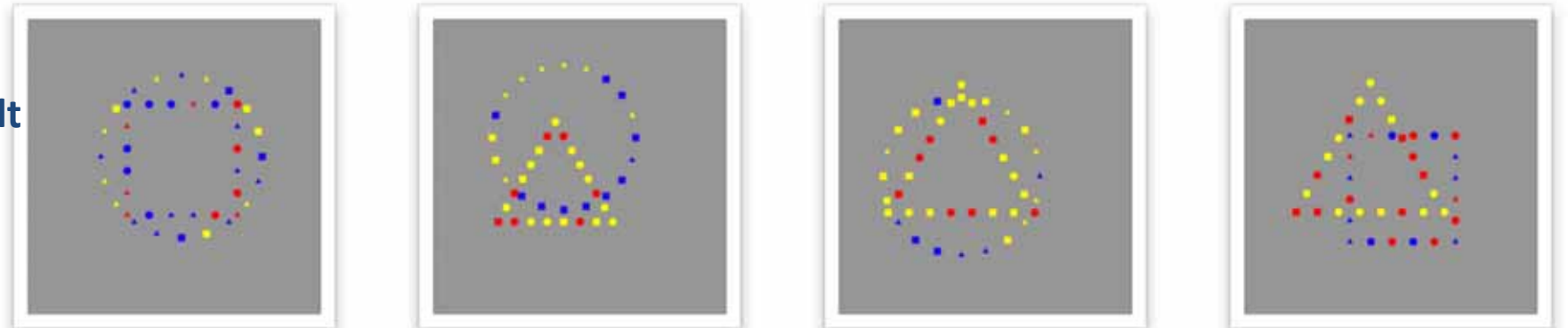
C
Quantity



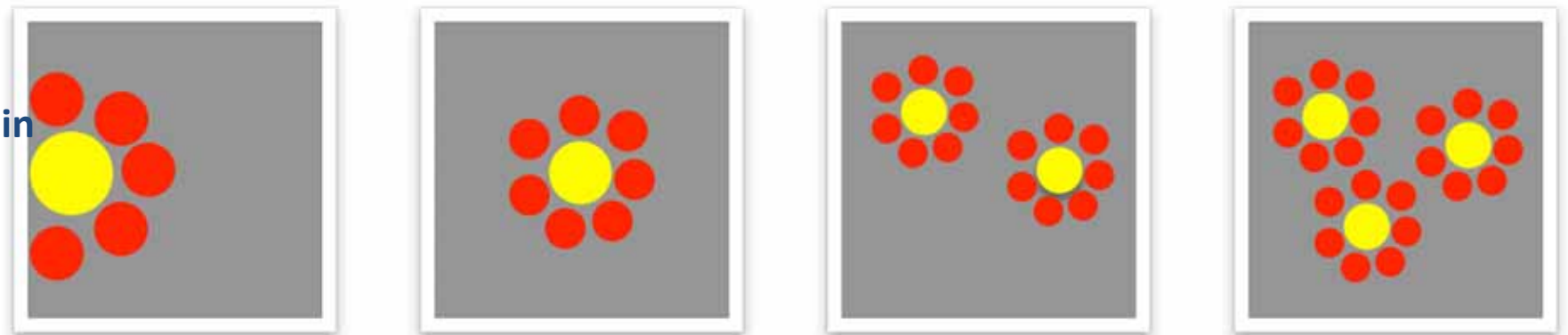
D
Arrangement



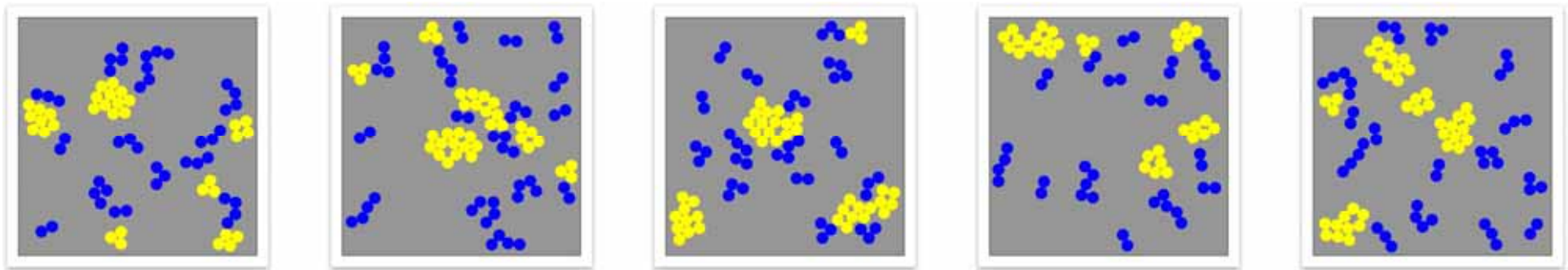
E
Gestalt



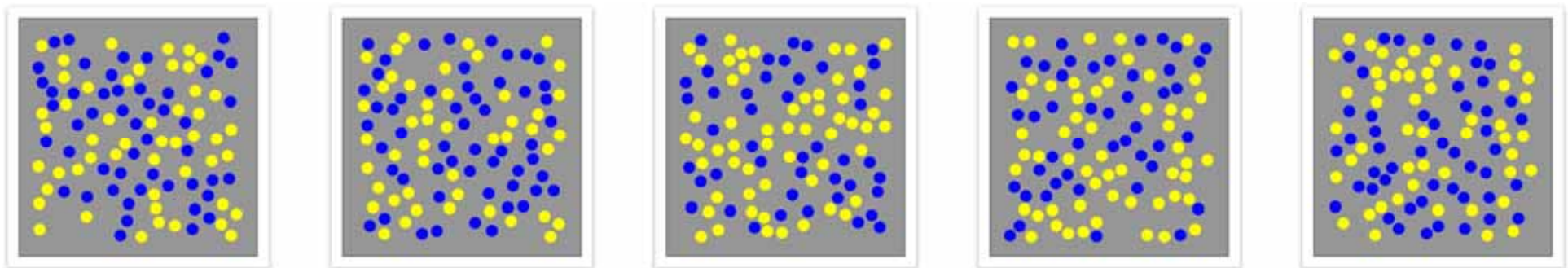
F
Domain



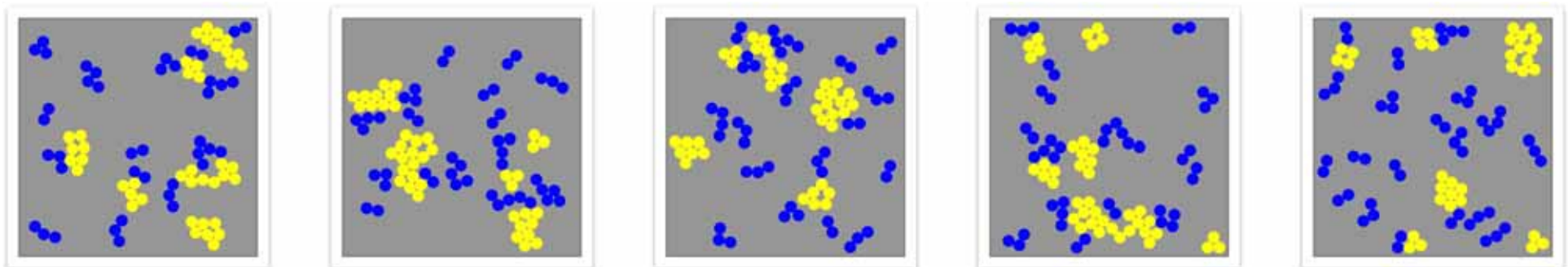
A) True



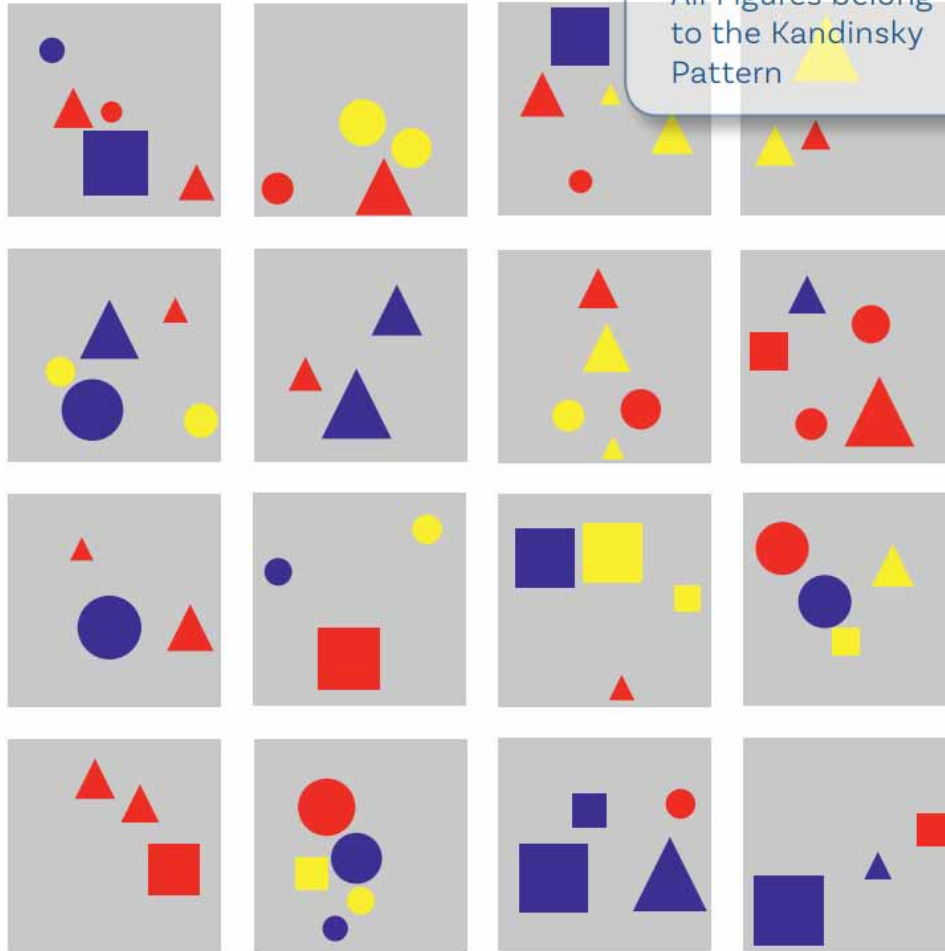
B) False



C) Counterfactual



☰ Part of the pattern



All Figures belong to the Kandinsky Pattern

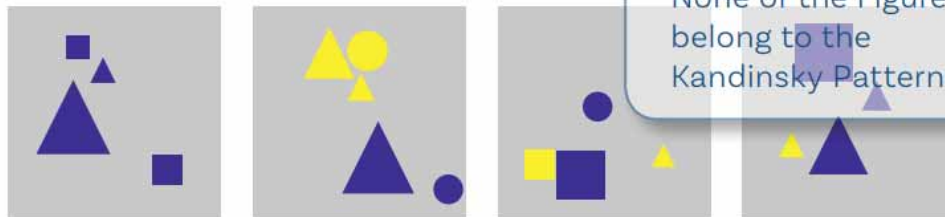
Hypothesis 1 _____

It only contains circles and triangles.

Hypothesis 2 _____

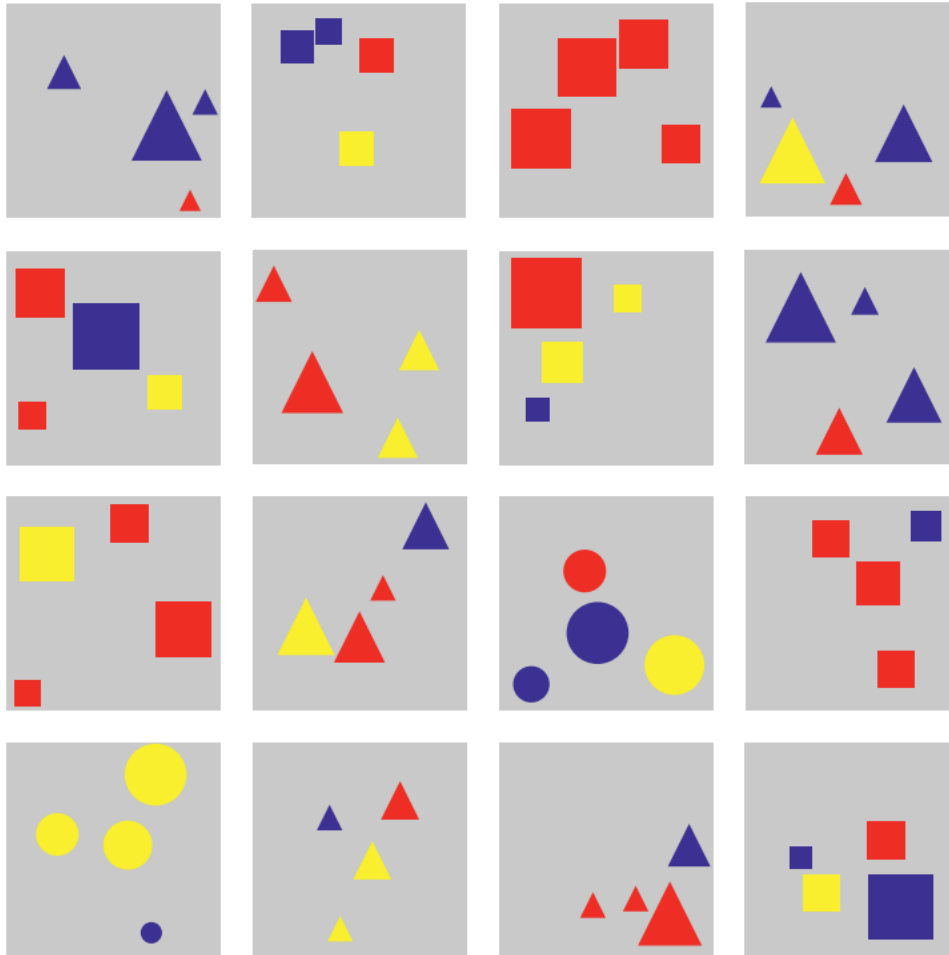
It contains at least a red object. ✓

≠ Not part of the pattern



None of the Figures belong to the Kandinsky Pattern

☰ Part of the pattern

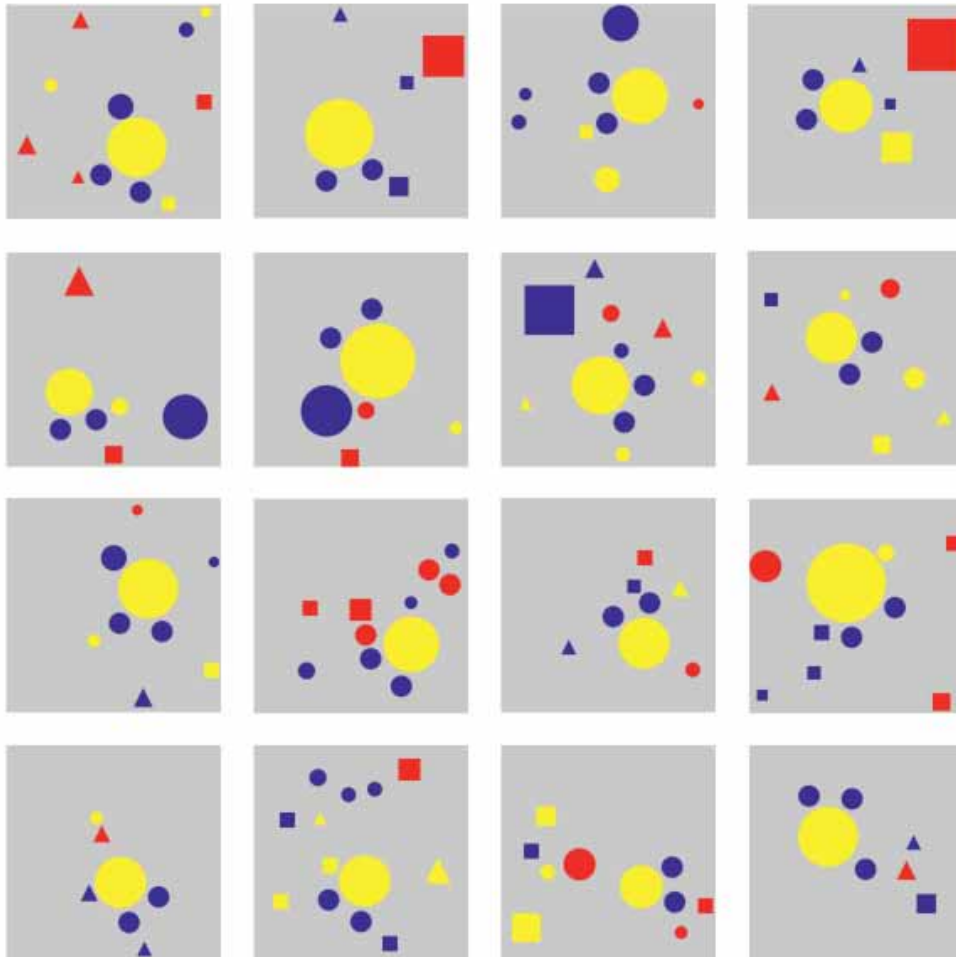


⚡ Not part of the pattern

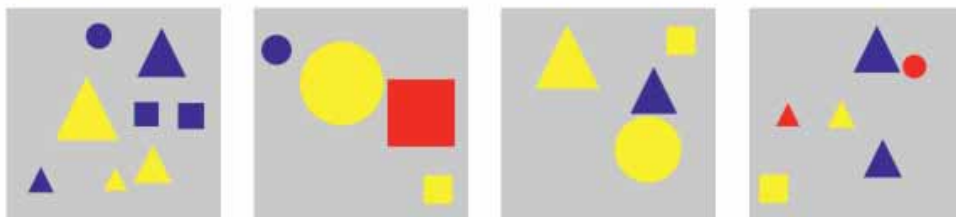


S2 Basic Pattern 2
 Title: **All of Same Shape** ->
 All objects have the same shape.
 Hint: Don't be distracted by the colors

☰ Part of the pattern



≠ Not part of the pattern



S8

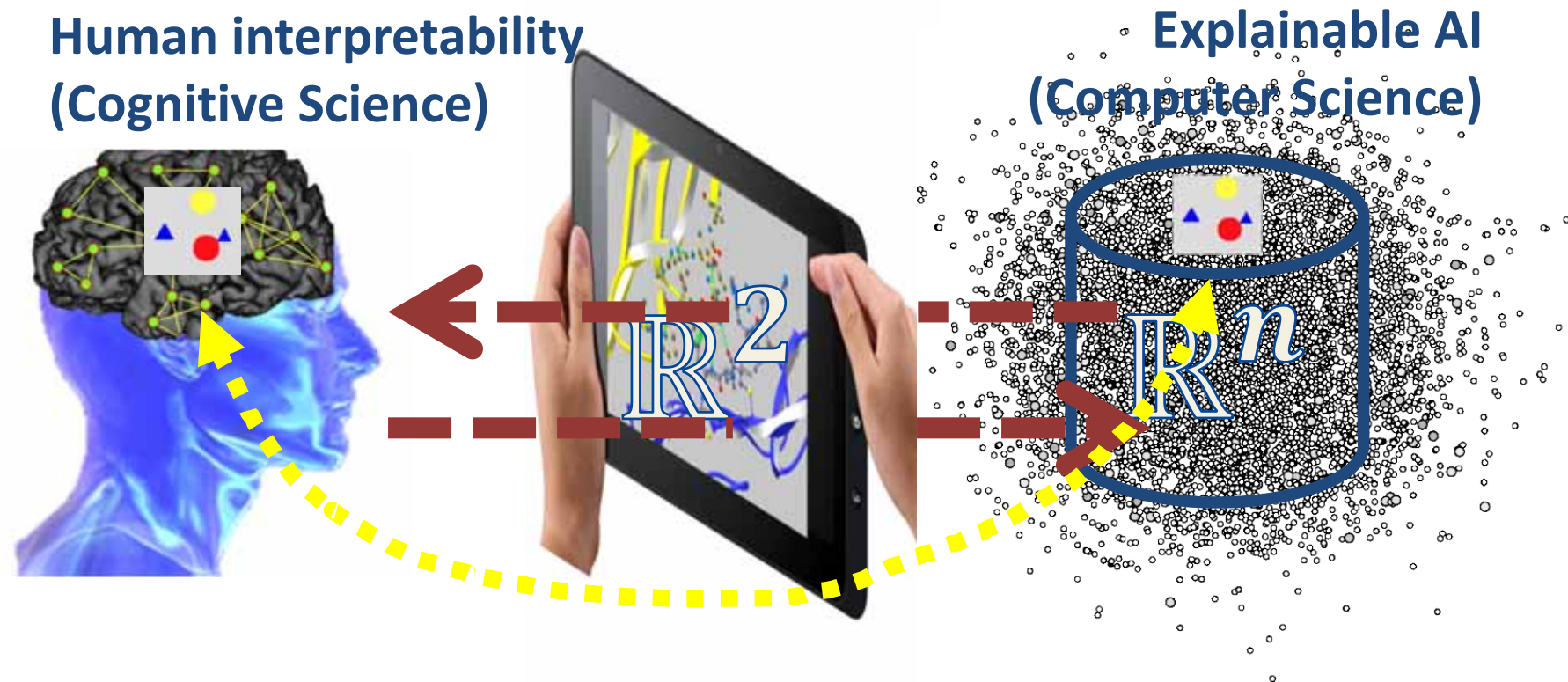
Basic Pattern 8

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.

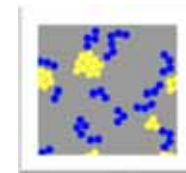
Conclusion

- 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.

Take part in the explainable-AI challenge:
<https://human-centered.ai/kandinsky-challenge>



Thank you !