

From Machine Learning to Explainable AI

Assoc.Prof. Dr. Andreas Holzinger

Holzinger-Group, HCI-KDD, Institute for Medical Informatics/Statistics

Medical University Graz, Austria

&

Institute of interactive Systems and Data Science

Graz University of Technology, Austria



a.holzinger@hci-kdd.org

<http://hci-kdd.org>

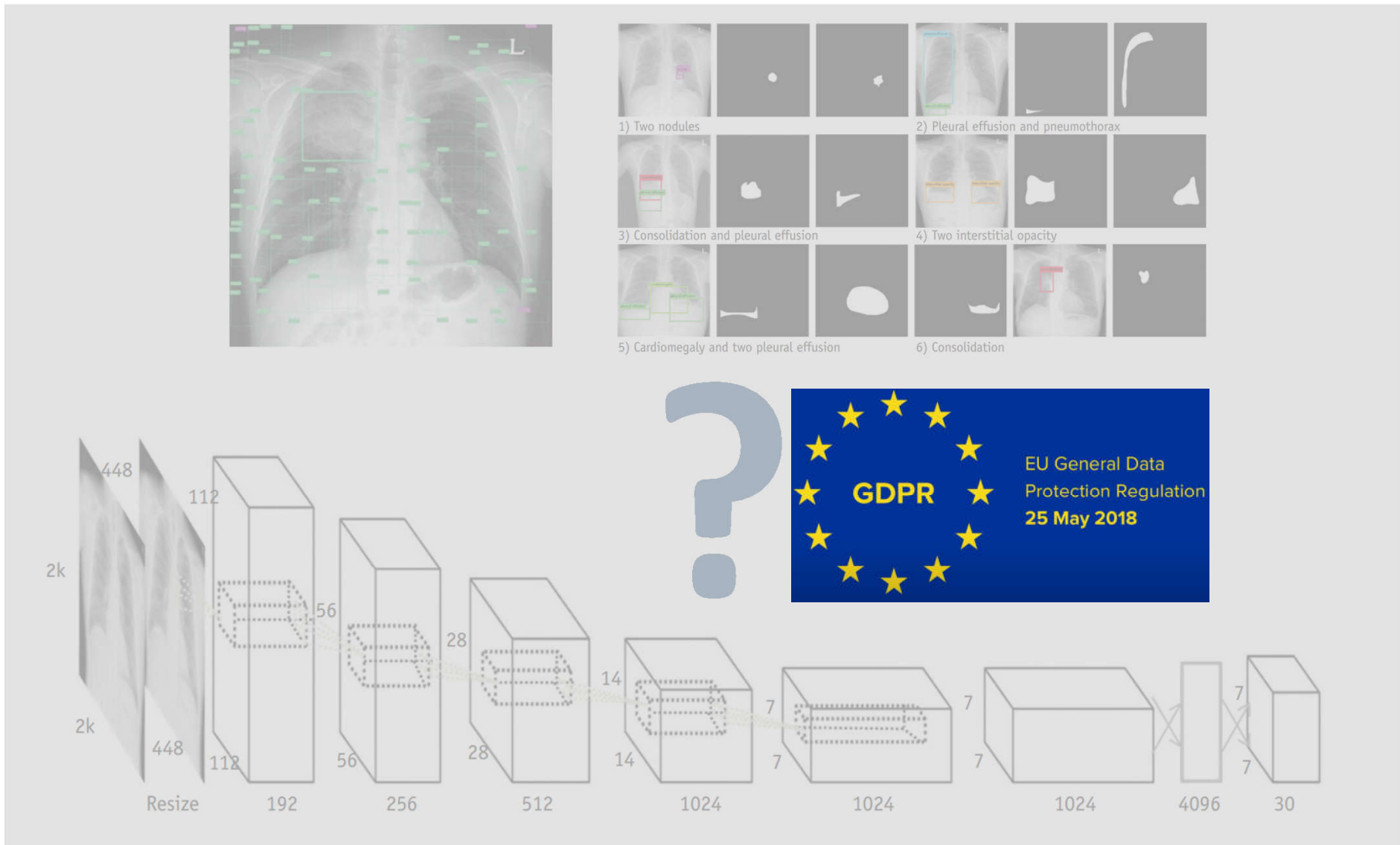


TECHNICKÁ
UNIVERZITA
V KOŠICIACH





Deep Learning is considered as “black-box” approach



June-Goo Lee, Sanghoon Jun, Young-Won Cho, Hyunna Lee, Guk Bae Kim, Joon Beom Seo & Namkug Kim 2017. Deep learning in medical imaging: general overview. Korean journal of radiology, 18, (4), 570-584, doi:10.3348/kjr.2017.18.4.570.

- **01 HCI-KDD – integrative ML**
- **02 Understanding Intelligence**
- **03 Application Area: Health**
- **04 automatic ML (aML)**
- **05 interactive ML (iML)**
- **06 towards explainable AI**

01 What is the



approach?

- **ML is a very practical field – algorithm development is at the core – however, successful ML needs a concerted effort of various topics ...**



MAchine Learning & Knowledge Extraction MAKE

(Safety) 4 - Privacy, Data Protection, Safety & Security



(Space and Time) 5 - Network, 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.



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



<http://hci-kdd.org/international-expert-network>

CD-MAKE

Cross Domain Conference for Machine Learning and Knowledge Extraction
co-located with [ARES 2018](#), Hamburg, Germany, August 27-30, 2018



[About CD-MAKE](#) [Call for Papers](#) [Committee](#) [Special Sessions](#) [Authors Area](#) [Venue & Registration](#) [Call for Sponsors](#) [Contact](#) [Archive](#) [ARES 2018](#) [Q](#)



Welcome

“Augmenting Human Intelligence with Artificial Intelligence”

International IFIP Cross Domain (CD) Conference for
Machine Learning & Knowledge Extraction (MAKE)

CD-MAKE 2018



*machine learning and
knowledge extraction*



[General Information](#)

[About CD-MAKE](#)

[ARES 2018](#)

CD-MAKE is held in conjunction with the
International Conference on Availability,

02 Understanding Intelligence

“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>



Humanoid AI

≠

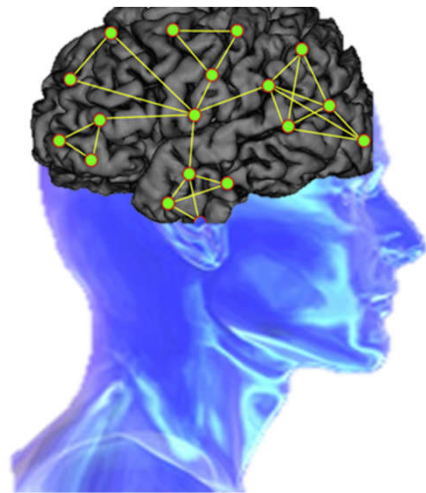
Human-level AI

Lotfi A. Zadeh 2008. Toward Human Level Machine Intelligence - Is It Achievable? The Need for a Paradigm Shift. IEEE Computational Intelligence Magazine, 3, (3), 11-22, doi:10.1109/MCI.2008.926583.

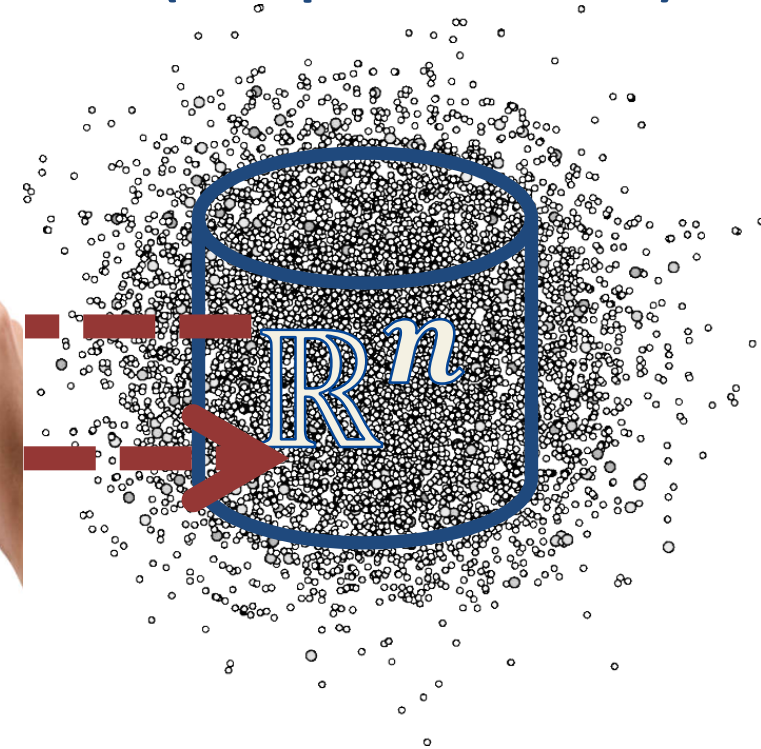
Understanding Context

- 1) learn from prior data
- 2) extract knowledge
- 2) 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


Human intelligence (Cognitive Science)



Artificial intelligence (Computer Science)



Holzinger, A. (2013). Human–Computer Interaction & Knowledge Discovery (HCI-KDD): What is the benefit of bringing those two fields to work together? In: Lecture Notes in Computer Science 8127 (pp. 319-328)



03 Application Area Health Informatics

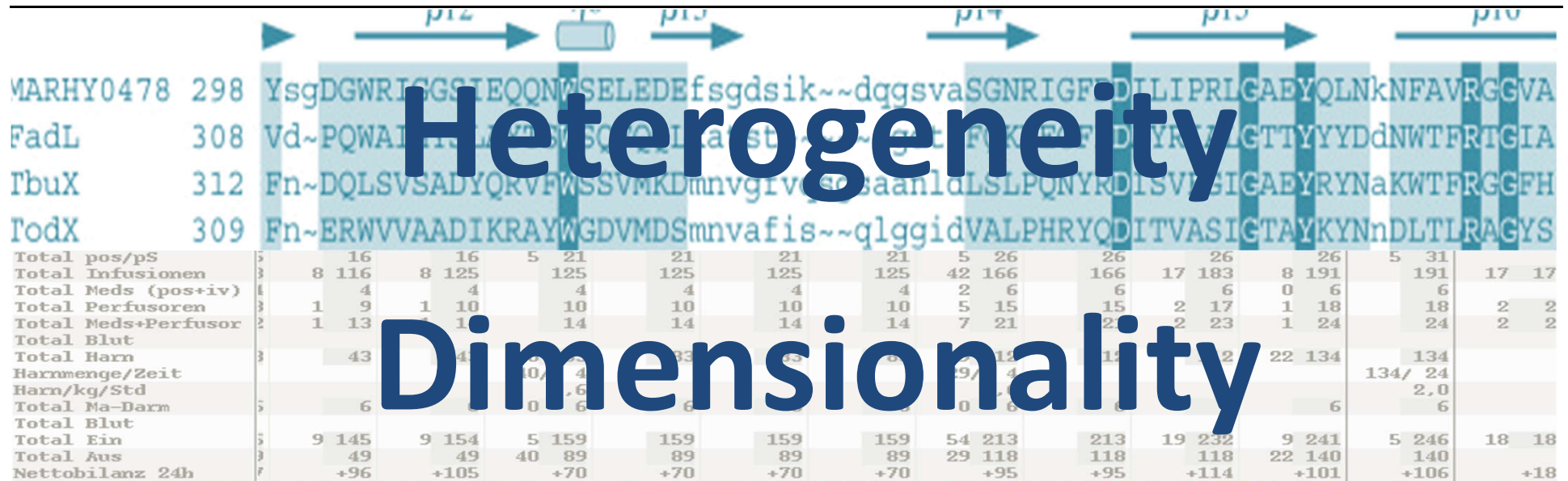
Why is this application area complex ?



Our central hypothesis: Information may bridge this gap

Holzinger, A. & Simonic, K.-M. (eds.) 2011. *Information Quality in e-Health. Lecture Notes in Computer Science LNCS 7058, Heidelberg, Berlin, New York: Springer.*



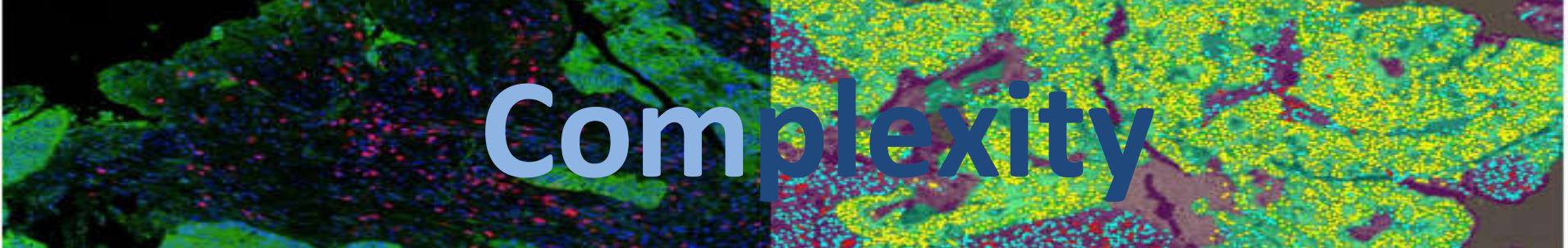


The image shows a protein sequence alignment with amino acid sequences for MARHY0478, FadL, TbuX, and TodX. Below the sequences is a summary table with 14 columns and 14 rows of data.

Total pos/pS	5	16	16	5	21	21	21	21	5	26	26	26	5	31				
Total Infusionen	3	8	116	8	125	125	125	125	42	166	166	17	183	8	191	17	17	
Total Meds (pos+iv)	4	4	4	4	4	4	4	4	2	6	6	0	6	6	6	6	6	
Total Perfusoren	3	1	9	1	10	10	10	10	5	15	15	2	17	1	18	18	2	2
Total Meds+Perfusor	2	1	13	1	14	14	14	14	7	21	21	2	23	1	24	24	2	2
Total Blut																		
Total Harn	3	43	43	43	43	43	43	43	13	134	134	2	134	134	134	134	134	
Harnmenge/Zeit																		
Harn/kg/Std																		
Total Ma-Darm	5	6	6	6	6	6	6	6	0	6	6	6	6	6	6	6	6	
Total Blut																		
Total Ein	5	9	145	9	154	159	159	159	54	213	213	19	232	9	241	241	18	18
Total Aus	3	49	49	40	89	89	89	89	29	118	118	118	118	22	140	140	140	140
Nettobilanz 24h	7	+96	+105	+70	+70	+70	+70	+70	+95	+95	+95	+114	+101	+106	+106	+106	+18	+18

Heterogeneity

Dimensionality



Complexity

Uncertainty

Probabilistic Information $p(x)$

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



Pierre Simon de Laplace (1749-1827)

$$\mathcal{D} = x_{1:n} = \{x_1, x_2, \dots, x_n\} \quad p(\mathcal{D}|\theta)$$



$$p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta) * p(\theta)}{p(\mathcal{D})}$$

$$posterior = \frac{likelihood * prior}{evidence}$$

The “inverse probability” allows to learn, to infer unknowns and to make predictions

d ... data

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

$\forall h, d \dots$

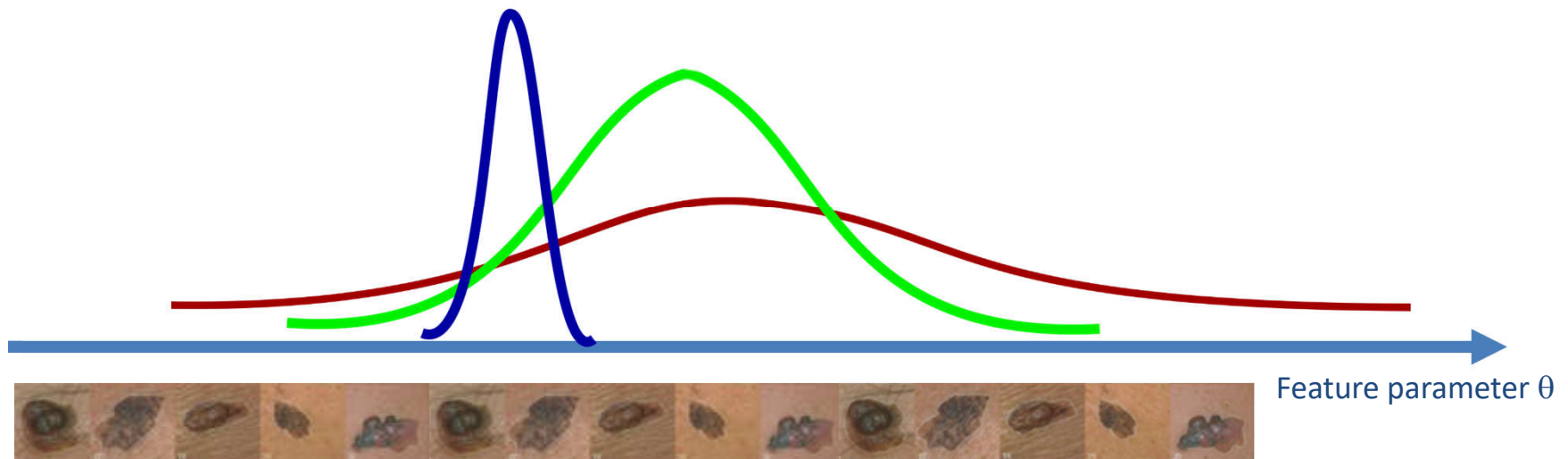
h ... hypotheses

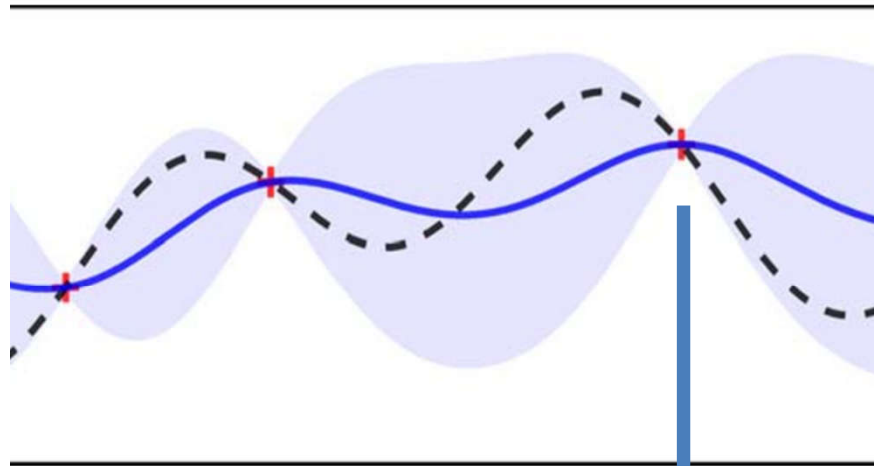
$$p(h|d) = \frac{p(d|h) * p(h)}{\sum_{h \in \mathcal{H}} p(d|h') p(h')}$$

Likelihood Prior Probability

Posterior Probability

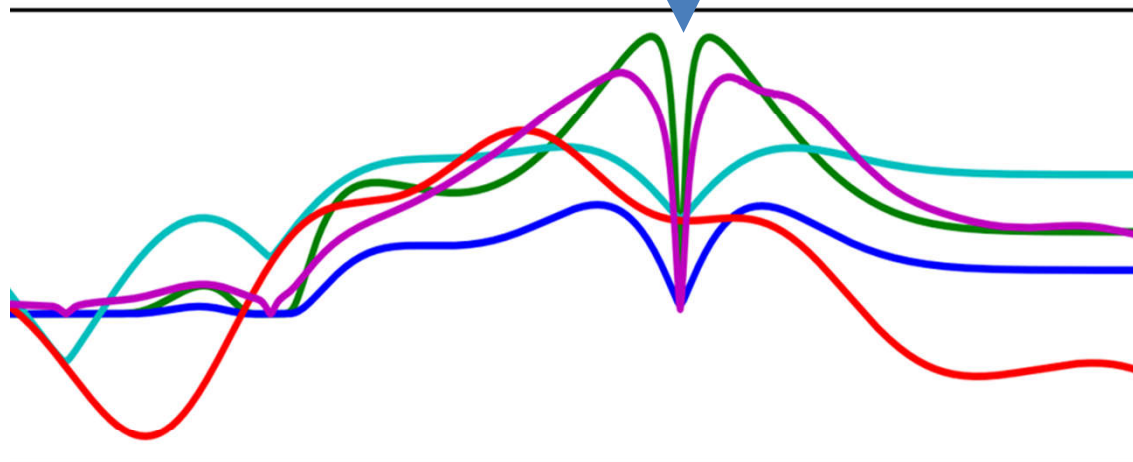
Problem in $\mathbb{R}^n \rightarrow$ complex





Algorithm 1 Bayesian optimization

- 1: **for** $n = 1, 2, \dots$ **do**
- 2: select new \mathbf{x}_{n+1} by optimizing acquisition function α
$$\mathbf{x}_{n+1} = \arg \max_{\mathbf{x}} \alpha(\mathbf{x}; \mathcal{D}_n)$$
- 3: query objective function to obtain y_{n+1}
- 4: augment data $\mathcal{D}_{n+1} = \{\mathcal{D}_n, (\mathbf{x}_{n+1}, y_{n+1})\}$
- 5: update statistical model
- 6: **end for**



- PI Probability of Improvement
- 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

amazon.co.uk Try Prime

Shop by Department Your Amazon.co.uk Today's Deals Gift Cards & Top Up Sell Help

Amazon.co.uk Today's Deals Warehouse Deals Outlet Subscribe & Save Vouchers Amazon Family Amazon Prime Amazon Video Amazon Student Mobile Apps An

Showing results for "glass cutter circular"

Show results for

DIY & Tools >

- Glass Cutters
- Cold Chisels
- Power Tools

Sports & Outdoors >

- Compasses

+ See All 13 I

Refine by

Delivery Opt

- Prime
- Free UK D

Brand

- sourcingm
- SODIAL(R



Silverline 101228 Circular Glass Cutter with 65-300 mm Diameter

by Silverline

£7.81 ~~£10.02~~ Prime
Get it by **Tomorrow, Sep 5**
Eligible for FREE UK Delivery

★★★★☆ 42

DIY & Tools: See all 162 items

More buying choices
£6.40 new (22 offers)



Highlander 3 Hole Thinsulate Balaclava

by Highlander

£1.99 - £7.00 Prime

More buying choices
£1.99 new (5 offers)

★★★★☆ 163

Sports & Outdoors: See all 5,918 items



Sanwood® Outdoor Motorcycle Cycling Ski Neck Protecting Lycra Balaclava Full Face Mask

by Phoenix B2C UK

£1.74 - £3.57

More buying choices
£0.01 new (4 offers)

★★★★☆ 73

Sports & Outdoors: See all 5,918 items



Guizzo, E. 2011. How google’s self-driving car works. IEEE Spectrum Online, 10, 18.



Image Source: <http://www.businessinsider.de/who-is-responsible-when-a-driverless-car-crashes-2016-2?r=US&IR=T>

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

Automotive



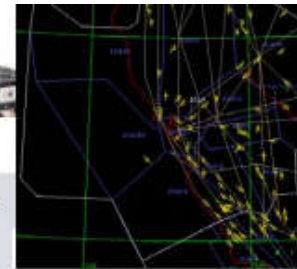
E-Corner, Siemens

Building Systems

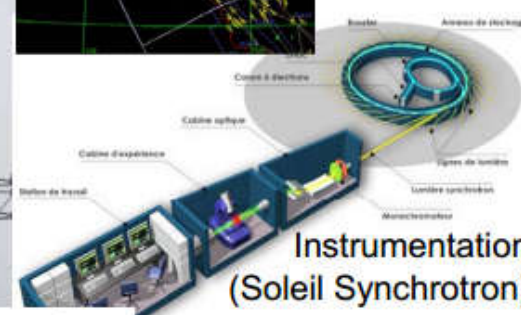


Avionics

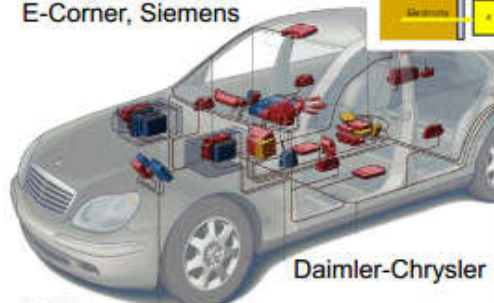
Telecommunications



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.

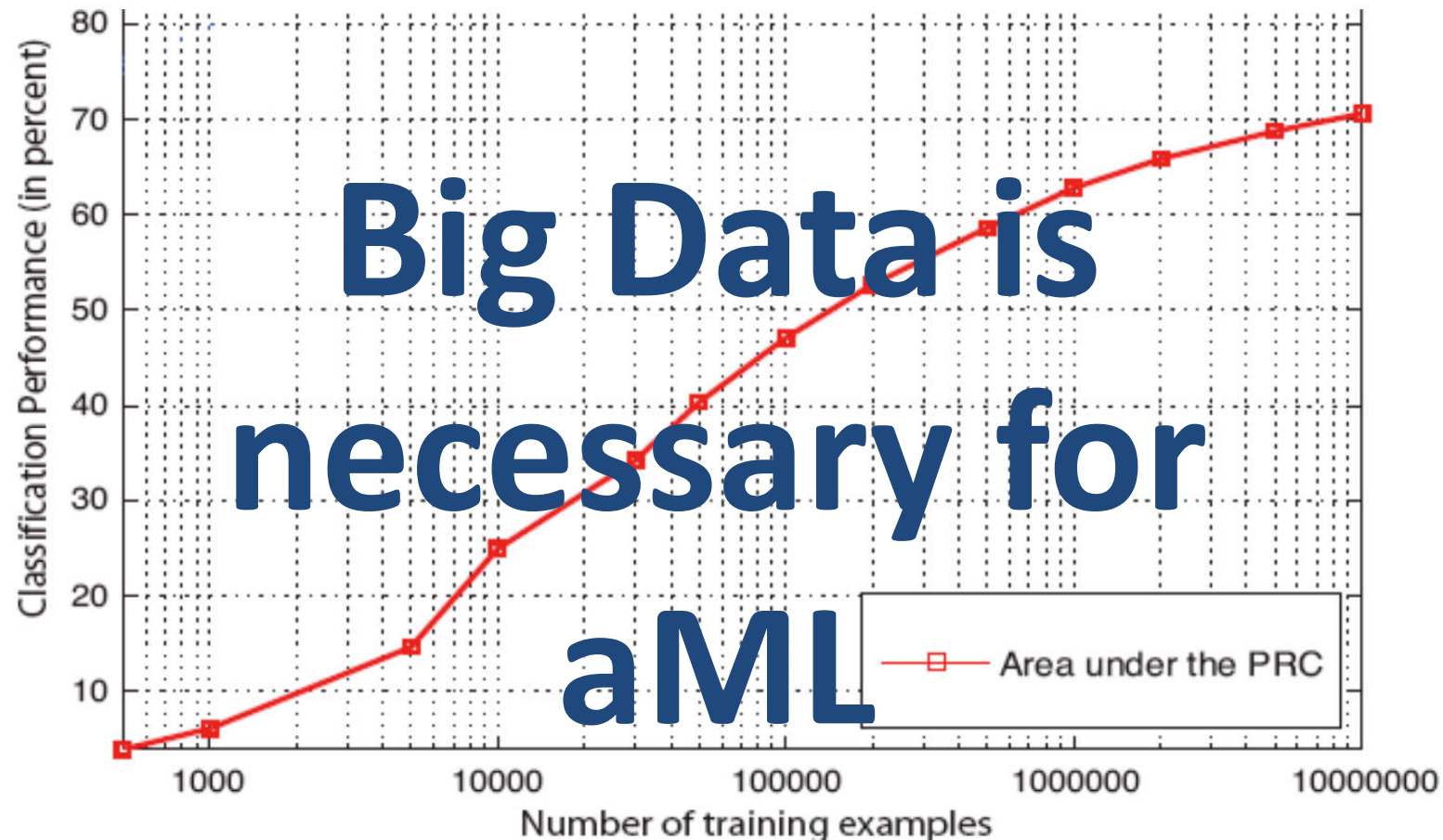
Military systems:



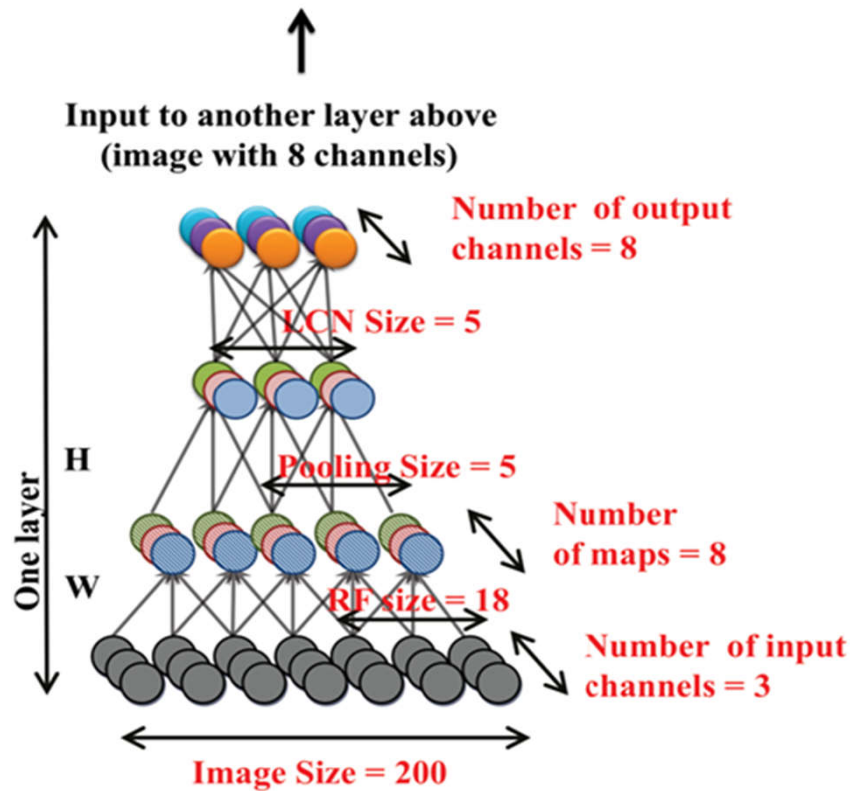
Courtesy of Doug Schmidt



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.



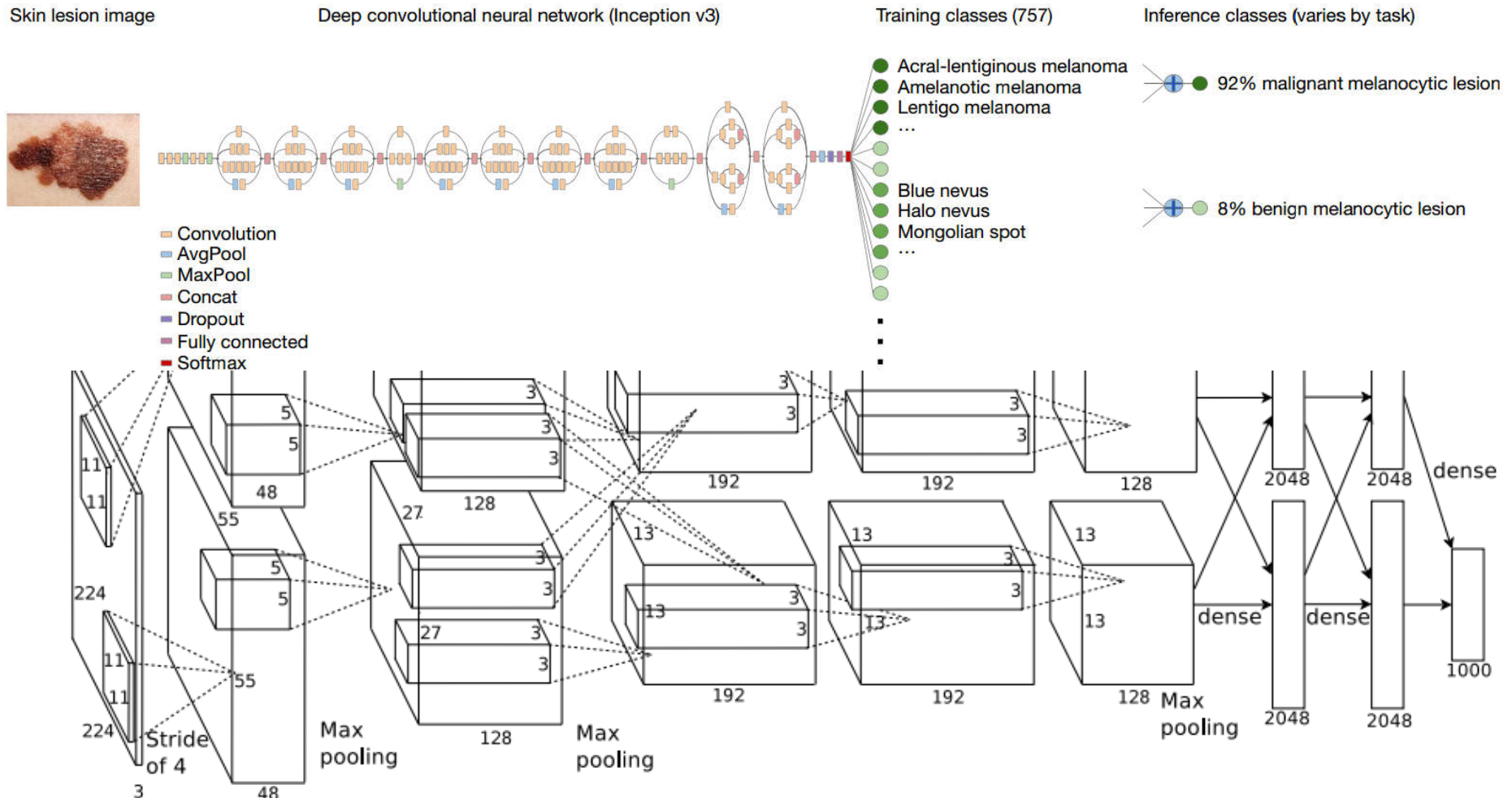
$$x^* = \arg \min_x f(x; W, H), \text{ subject to } \|x\|_2 = 1.$$

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.

1,28 million images ...

Esteva, A., Kuprel, B., Novoa, R. A., Ko, J., Swetter, S. M., Blau, H. M. & Thrun, S. 2017. Dermatologist-level classification of skin cancer with deep neural networks. Nature, 542, (7639), 115-118, doi:10.1038/nature21056.



Krizhevsky, A., Sutskever, I. & Hinton, G. E. Imagenet classification with deep convolutional neural networks. In: Pereira, F., Burges, C. J. C., Bottou, L. & Weinberger, K. Q., eds. Advances in neural information processing systems (NIPS 2012), 2012 Lake Tahoe. 1097-1105.

- 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
(still) need a
human-in-the-loop**

05 iML

- iML := algorithms which interact with agents*) and can optimize their learning behaviour through this interaction

***) where the agents can be human**

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



Image Source: 10 Ways Technology is Changing Healthcare <http://newhealthpost.com> Posted online on April 22, 2018



Image Source: Cisco (2008). Cisco Health Presence Trial at Aberdeen Royal Infirmary in Scotland



Image is in the public domain

- **Humans can generalize even from few examples ...**
 - They learn relevant representations
 - Can disentangle the explanatory factors
 - Find 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.

Even Children can make inferences from little data ...



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.



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

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

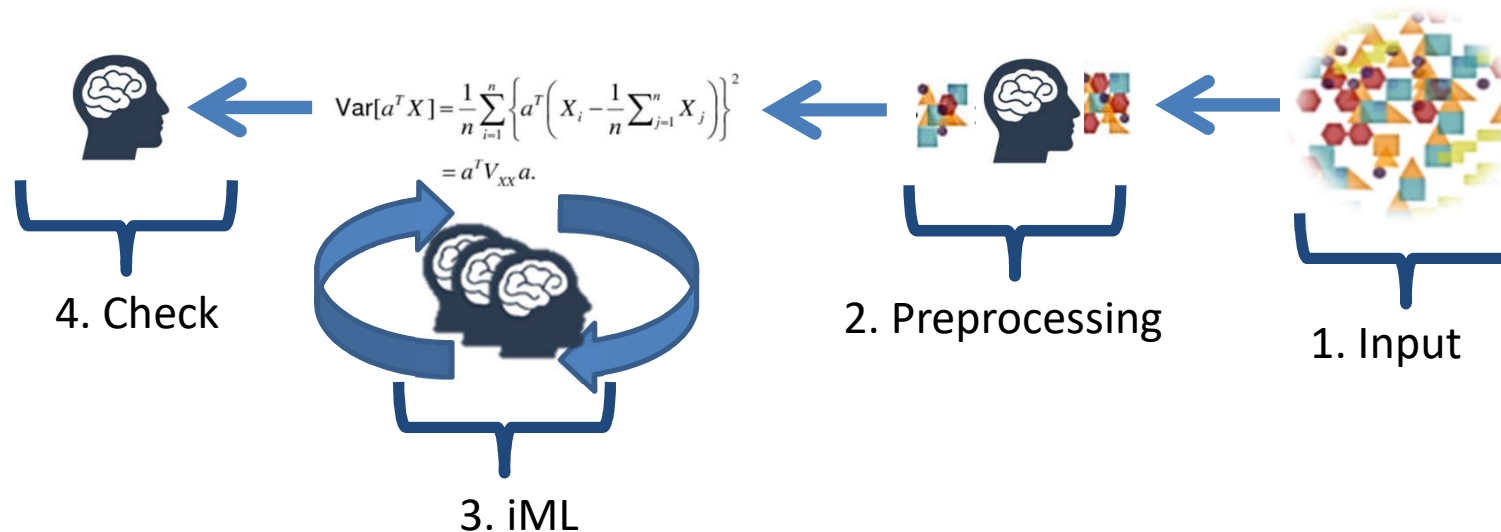
Jascha Sohl-Dickstein
Google Brain
jaschasd@google.com

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.

Interactive Machine Learning: Human is seen as an agent involved in the actual learning phase, step-by-step influencing measures such as distance, cost functions ...



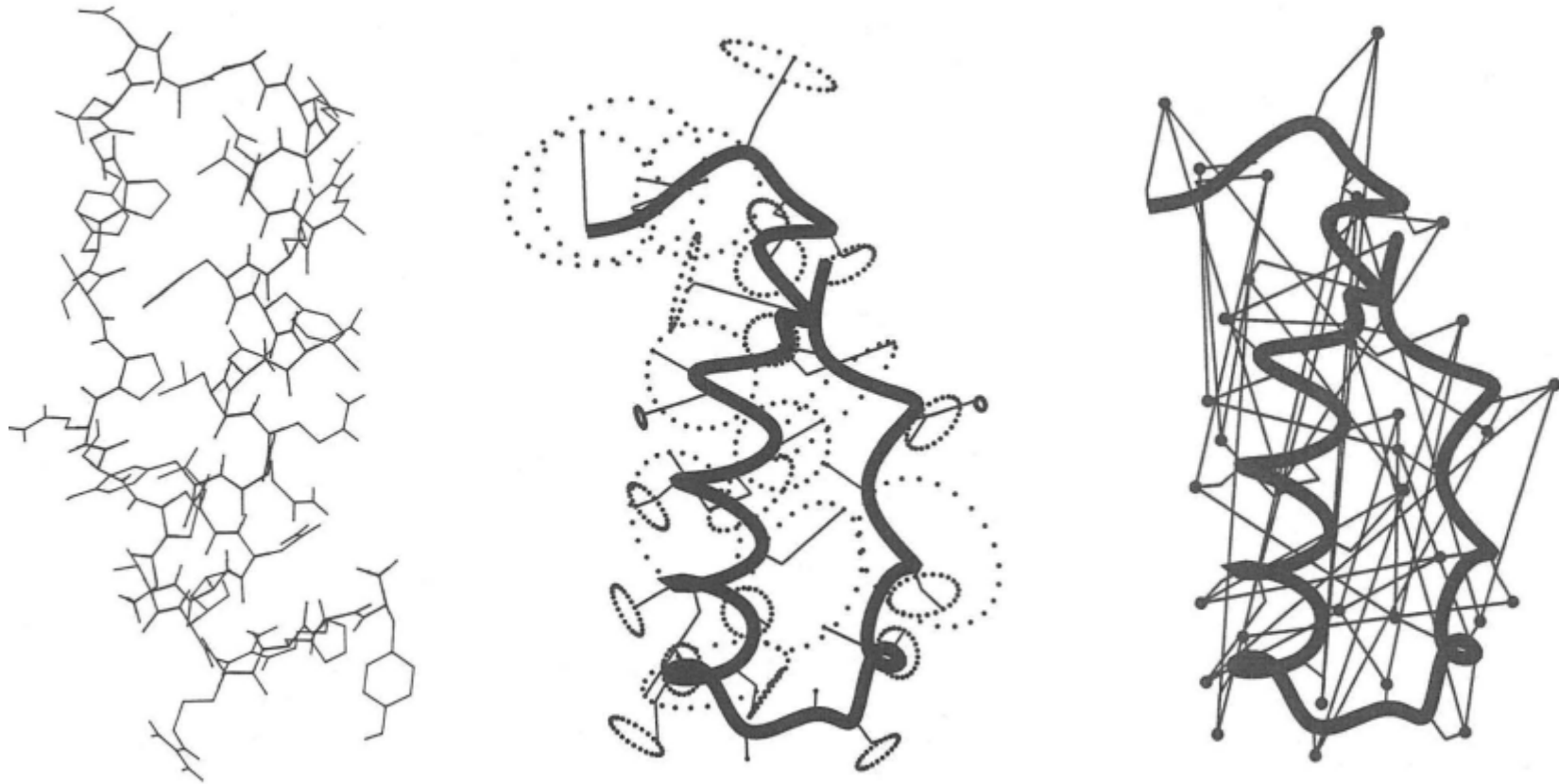
Holzinger, A. 2016. Interactive Machine Learning for Health Informatics: When do we need the human-in-the-loop? Brain Informatics (BRIN), 3, (2), 119-131, doi:10.1007/s40708-016-0042-6.

- **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.



Bohr, H. & Brunak, S. 1989. A travelling salesman approach to protein conformation. *Complex Systems*, 3, 9-28

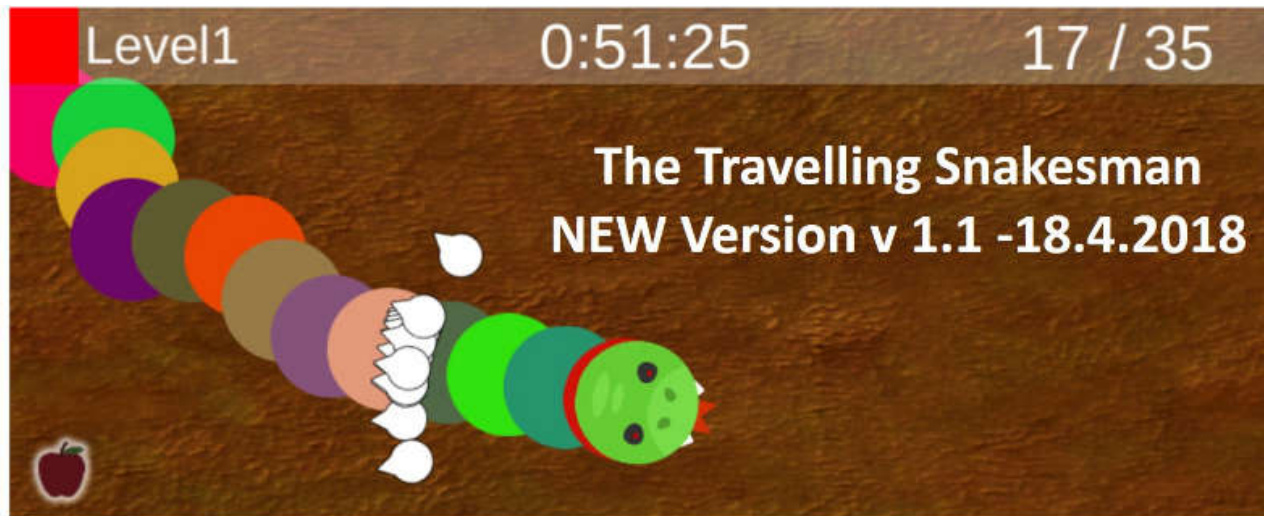
<http://hci-kdd.org/projects/iml-experiment>

```
Input : ProblemSize,  $m$ ,  $\beta$ ,  $\rho$ ,  $\sigma$ ,  $q_0$ 
Output:  $P_{best}$ 
 $P_{best} \leftarrow$  CreateHeuristicSolution(ProblemSize);
 $P_{best\_cost} \leftarrow$  Cost( $P_{best}$ );
 $Pheromone_{init} \leftarrow \frac{1.0}{ProblemSize \times P_{best\_cost}}$ ;
 $Pheromone \leftarrow$  InitializePheromone( $Pheromone_{init}$ );
while  $\neg StopCondition()$  do
  for  $i = 1$  to  $m$  do
     $S_i \leftarrow$  ConstructSolution(Pheromone, ProblemSize,  $\beta$ ,  $q_0$ );
     $S_{i\_cost} \leftarrow$  Cost( $S_i$ );
    if  $S_{i\_cost} \leq P_{best\_cost}$  then
       $P_{best\_cost} \leftarrow S_{i\_cost}$ ;
       $P_{best} \leftarrow S_i$ ;
    end
    LocalUpdateAndDecayPheromone(Pheromone,  $S_i$ ,  $S_{i\_cost}$ ,  $\rho$ );
  end
  GlobalUpdateAndDecayPheromone(Pheromone,  $P_{best}$ ,  $P_{best\_cost}$ ,  $\rho$ );
  while  $isUserInteraction()$  do
    GlobalAddAndRemovePheromone(Pheromone,  $P_{best}$ ,  $P_{best\_cost}$ ,  $\rho$ );
  end
end
return  $P_{best}$ ;
```

Holzinger, A., Plass, M., Holzinger, K., Crisan, G., Pintea, C. & Palade, V. 2016. Towards interactive Machine Learning (iML): Applying Ant Colony Algorithms to solve the Traveling Salesman Problem with the Human-in-the-Loop approach. Springer Lecture Notes in Computer Science LNCS 9817. 81-95, doi:10.1007/978-3-319-45507-56.

$$p_{ij} = \frac{[\tau_{ij}]^{\alpha} \cdot [\eta_{ij}]^{\beta}}{\sum_{l \in J_i^k} [\tau(t)]^{\alpha} \cdot [\eta]^{\beta}}$$

- p_{ij} ... **probability** of ants that they, at a particular node i , select the route from node $i \rightarrow j$ (“**heuristic desirability**”)
- $\alpha > 0$ and $\beta > 0$... the **influence parameters** (α ... history coefficient, β ...heuristic coefficient) usually $\alpha \approx \beta \approx 2 < 5$
- τ_{ij} ... the **pheromone value** for the components, i.e. the amount of pheromone on edge (i, j)
- k ... the set of usable components
- J_i ... the set of nodes that ant k can reach from v_i (tabu list)
- $\eta_{ij} = \frac{1}{d_{ij}}$... attractiveness computed by a heuristic, indicating the “a-priori **desirability**” of the move



Instruction to the Travelling Snakesman NEW versions v1.1 and v2 (as of 18.April 2018)

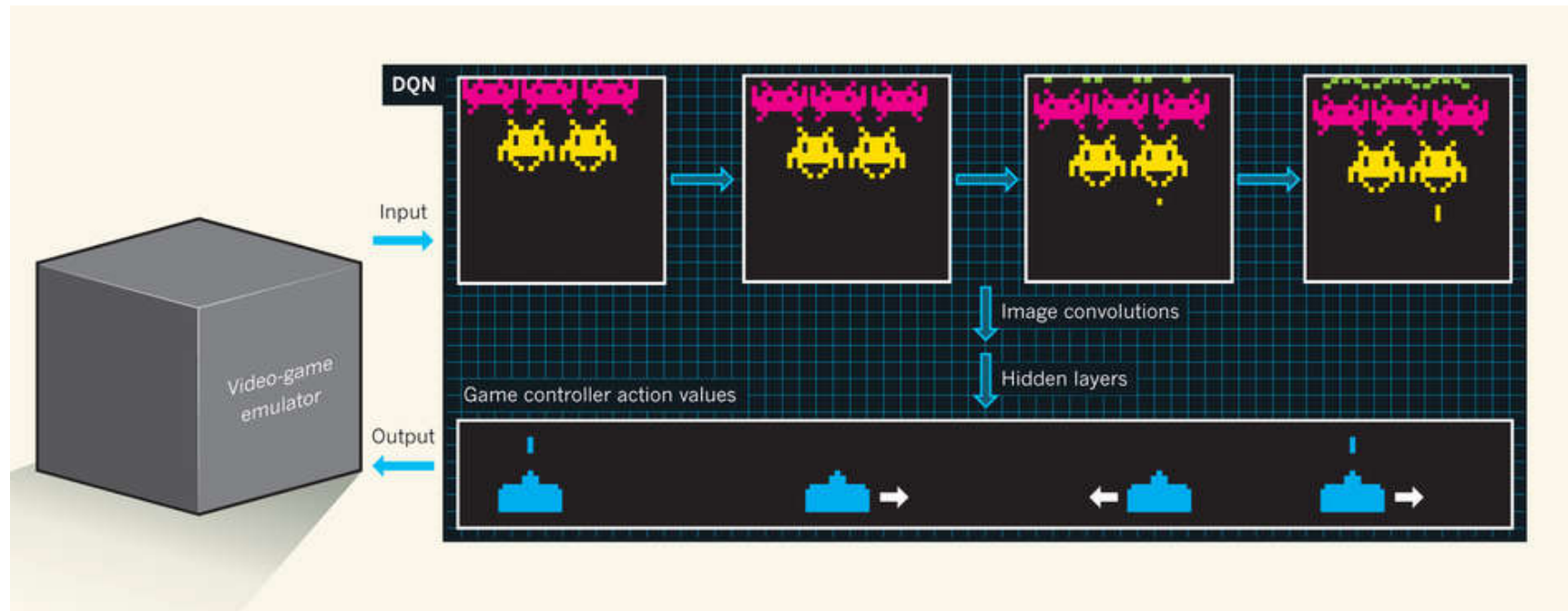
This page is current as of May, 11, 2018 13:15 CEST

This game uses an iML algorithm for computations in the background. We try to measure, if human interaction with this algorithm leads to better solutions than the algorithm running automatically without any interaction.

<https://hci-kdd.org/gamification-interactive-machine-learning>

**YOU ARE A SNAKE AND YOUR GOAL IS TO
EAT ALL APPLES AS FAST AS POSSIBLE!
ENJOY PLAYING BOTH VARIANTS!**

- You find the links for the Browser and for Android below (just click)
- 1) Enter a name
- 2) Select the level (1 = easy, 3=difficult)
- 3) Press "Play!" With your mouse/touch you direct the snake and your goal is to eat all apples as fast as possible!



Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M., Fidjeland, A. K., Ostrovski, G., Petersen, S., Beattie, C., Sadik, A., Antonoglou, I., King, H., Kumaran, D., Wierstra, D., Legg, S. & Hassabis, D. 2015. Human-level control through deep reinforcement learning. **Nature**, 518, (7540), 529-533, doi:10.1038/nature14236



06 Towards Explainable AI

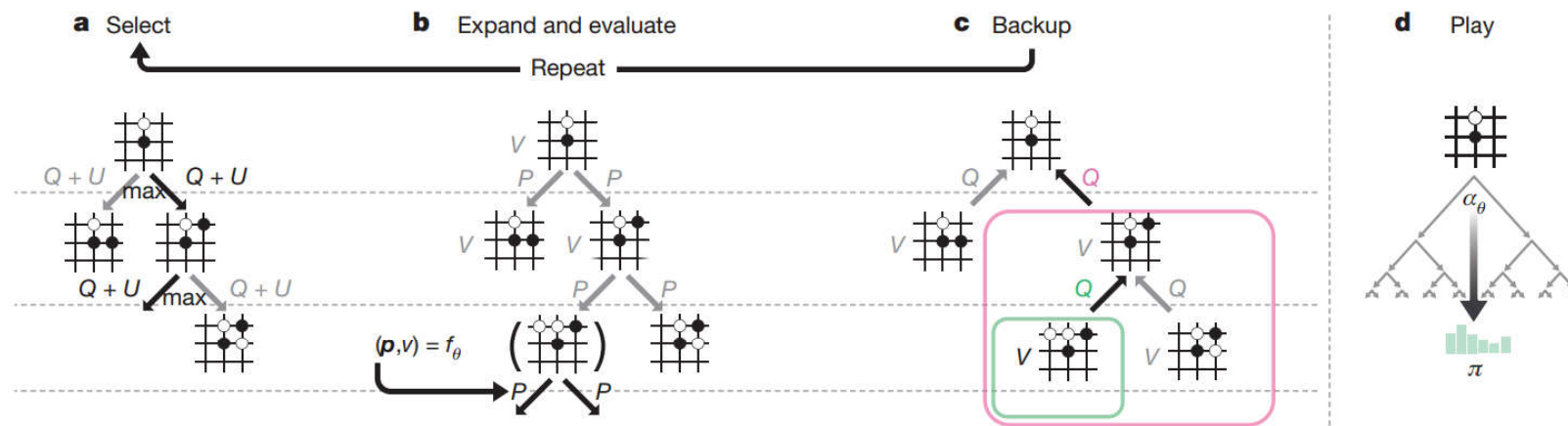


Figure 2 | MCTS in AlphaGo Zero. **a**, Each simulation traverses the tree by selecting the edge with maximum action value Q , plus an upper confidence bound U that depends on a stored prior probability P and visit count N for that edge (which is incremented once traversed). **b**, The leaf node is expanded and the associated position s is evaluated by the neural network $(\mathbf{p}, v) = f_{\theta}(s)$; the vector of P values are stored in

the outgoing edges from s . **c**, Action value Q is updated to track the mean of all evaluations V in the subtree below that action. **d**, Once the search is complete, search probabilities π are returned, proportional to $N^{1/\tau}$, where N is the visit count of each move from the root state and τ is a parameter controlling temperature.

19 OCTOBER 2017 | VOL 550 | NATURE | 355

$$(\mathbf{p}, v) = f_{\theta}(s) \quad \text{and} \quad l = (z - v)^2 - \pi^T \log \mathbf{p} + c \|\theta\|^2$$

David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, Yutian Chen, Timothy Lillicrap, Fan Hui, Laurent Sifre, George Van Den Driessche, Thore Graepel & Demis Hassabis 2017. Mastering the game of go without human knowledge. Nature, 550, (7676), 354-359, doi:doi:10.1038/nature24270.



David Silver, Aja Huang, Chris J. Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, Sander Dieleman, Dominik Grewe, John Nham, Nal Kalchbrenner, Ilya Sutskever, Timothy Lillicrap, Madeleine Leach, Koray Kavukcuoglu, Thore Graepel & Demis Hassabis 2016. Mastering the game of Go with deep neural networks and tree search. Nature, 529, (7587), 484-489, doi:10.1038/nature16961.



a woman riding a horse on a dirt road



an airplane is parked on the tarmac at an airport



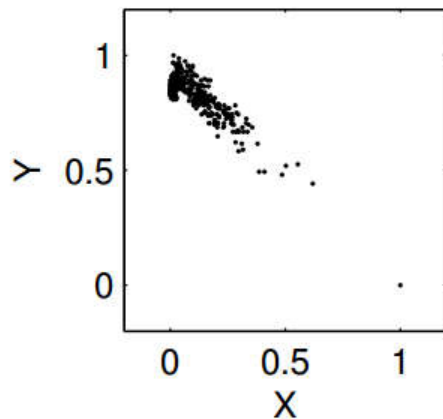
a group of people standing on top of a beach

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.

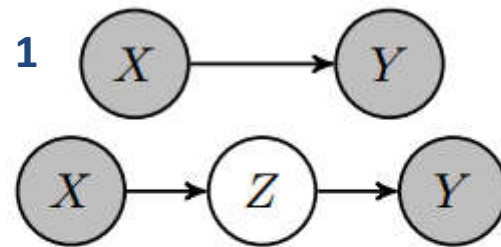
Image Captions by deep learning : github.com/karpathy/neuraltalk2

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

Decide if $X \rightarrow Y$, or $Y \rightarrow X$ using only observed data

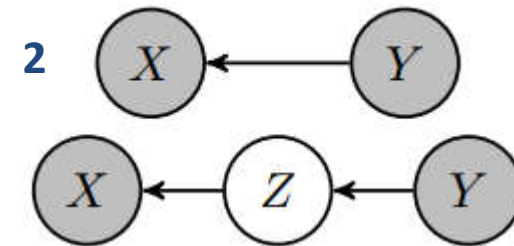


Joris M. Mooij, Jonas Peters, Dominik Janzing, Jakob Zscheischler & Bernhard Schölkopf 2016. Distinguishing cause from effect using observational data: methods and benchmarks. The Journal of Machine Learning Research, 17, (1), 1103-1204.



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

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



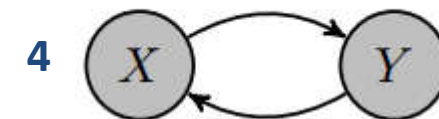
$$\mathbb{P}_Y = \mathbb{P}_{Y|\text{do}(x)} \neq \mathbb{P}_{Y|x}$$

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



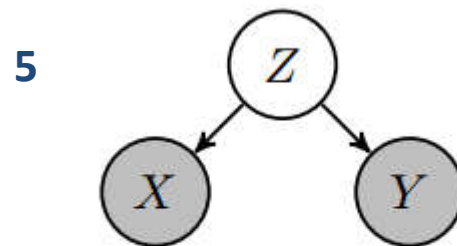
$$\mathbb{P}_Y = \mathbb{P}_{Y|\text{do}(x)} = \mathbb{P}_{Y|x}$$

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



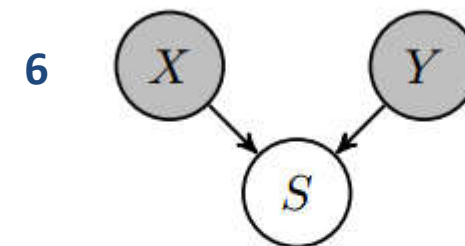
$$\mathbb{P}_Y \neq \mathbb{P}_{Y|\text{do}(x)} \neq \mathbb{P}_{Y|x}$$

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



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

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



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

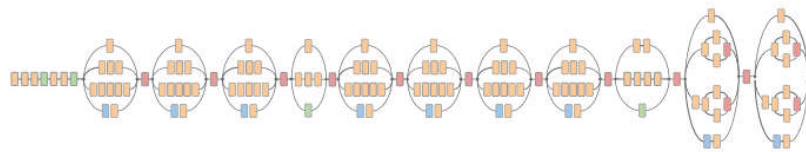
$$\mathbb{P}_{X|s} \neq \mathbb{P}_{X|\text{do}(y),s} = \mathbb{P}_{X|y,s}$$

Why?

Skin lesion image



Deep convolutional neural network (Inception v3)



- Convolution
- AvgPool
- MaxPool
- Concat
- Dropout
- Fully connected
- Softmax

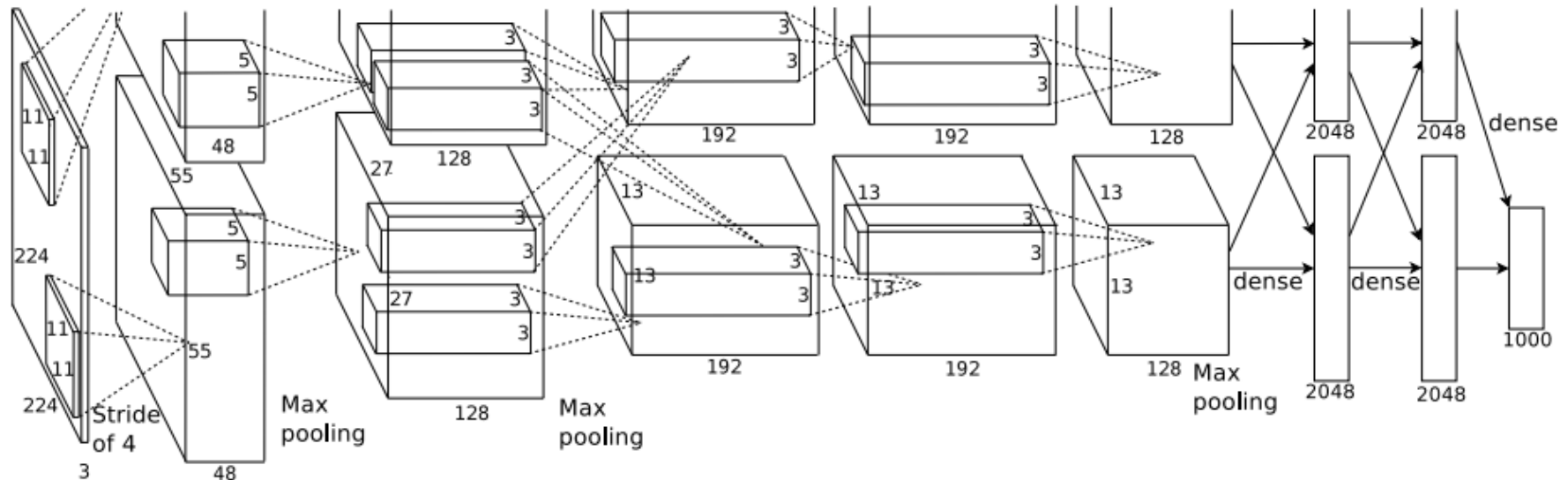
Training classes (757)

- Acral-lentiginous melanoma
- Amelanotic melanoma
- Lentigo melanoma
- ...
- Blue nevus
- Halo nevus
- Mongolian spot
- ...

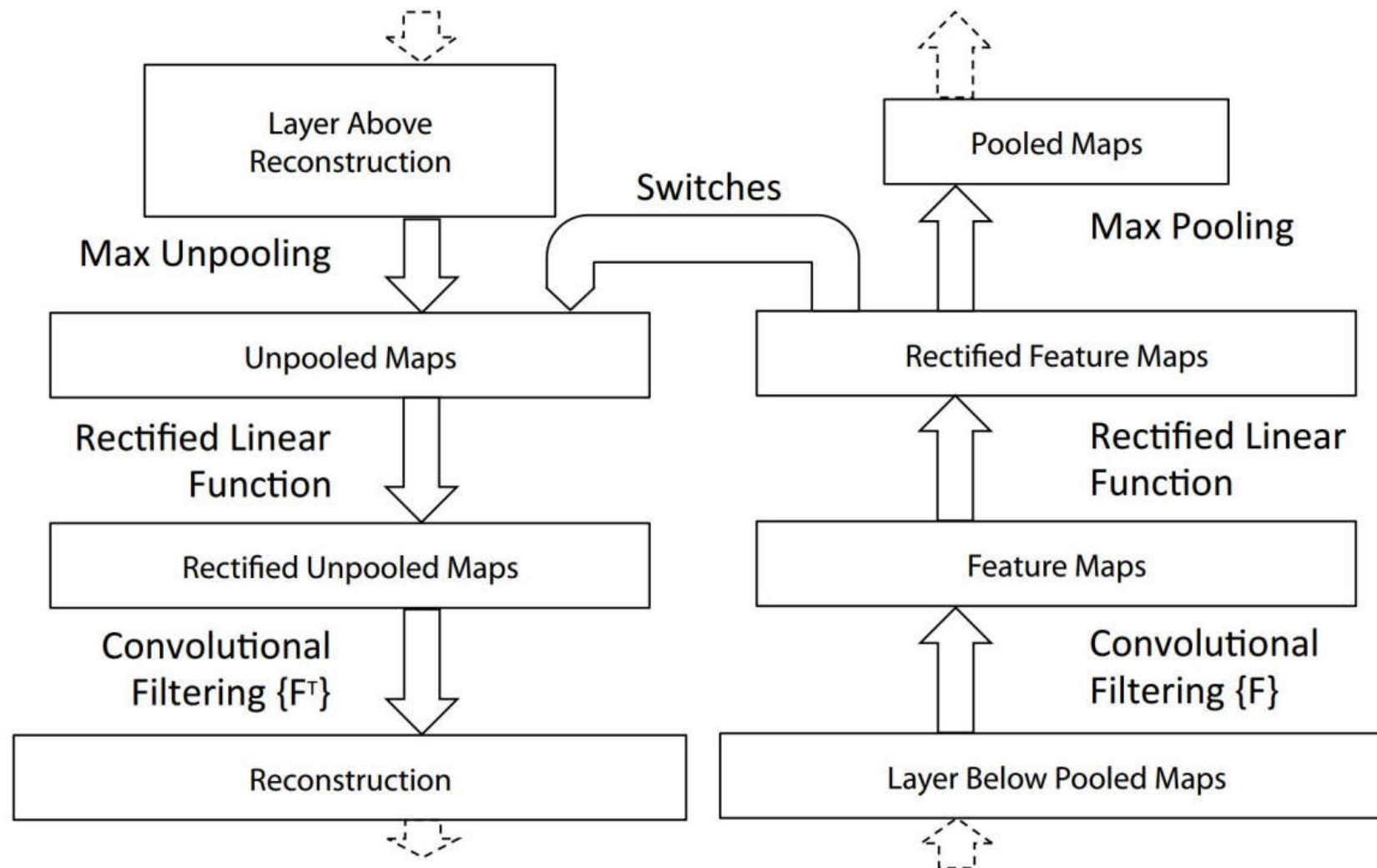
Inference classes (varies by task)

- ⊕ ● 92% malignant melanocytic lesion
- ⊕ ● 8% benign melanocytic lesion

Esteva, A., Kuprel, B., Novoa, R. A., Ko, J., Swetter, S. M., Blau, H. M. & Thrun, S. 2017. Dermatologist-level classification of skin cancer with deep neural networks. *Nature*, 542, (7639), 115-118, doi:10.1038/nature21056.



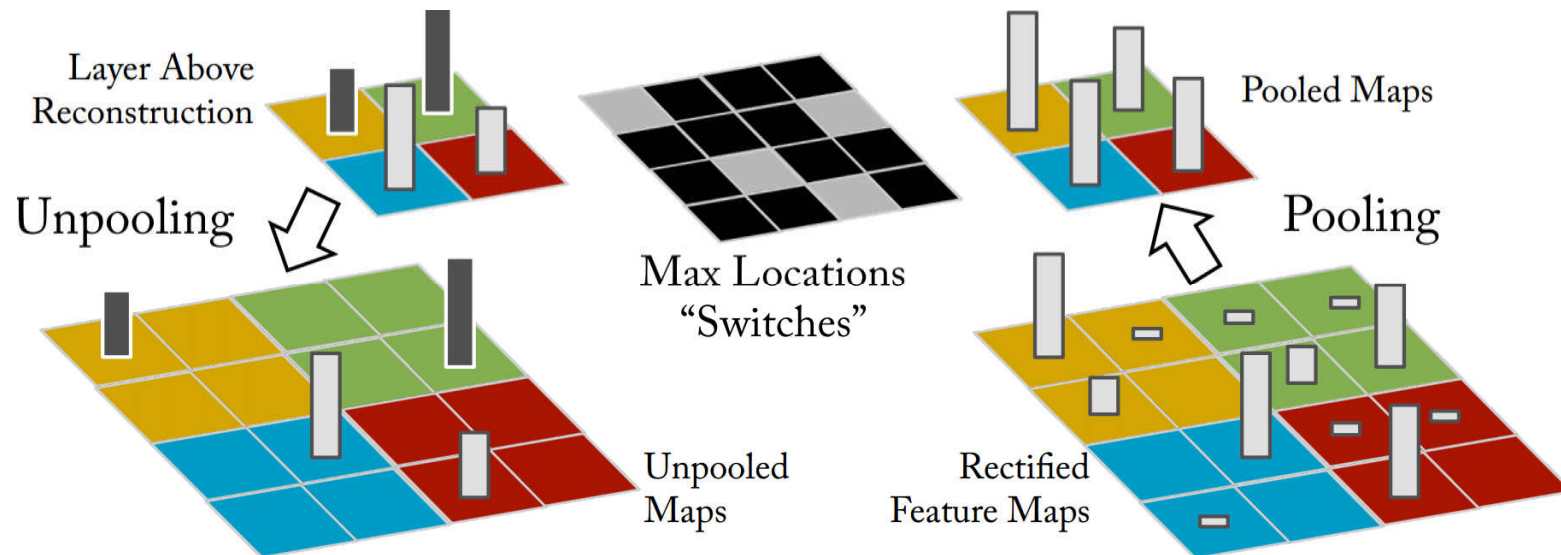
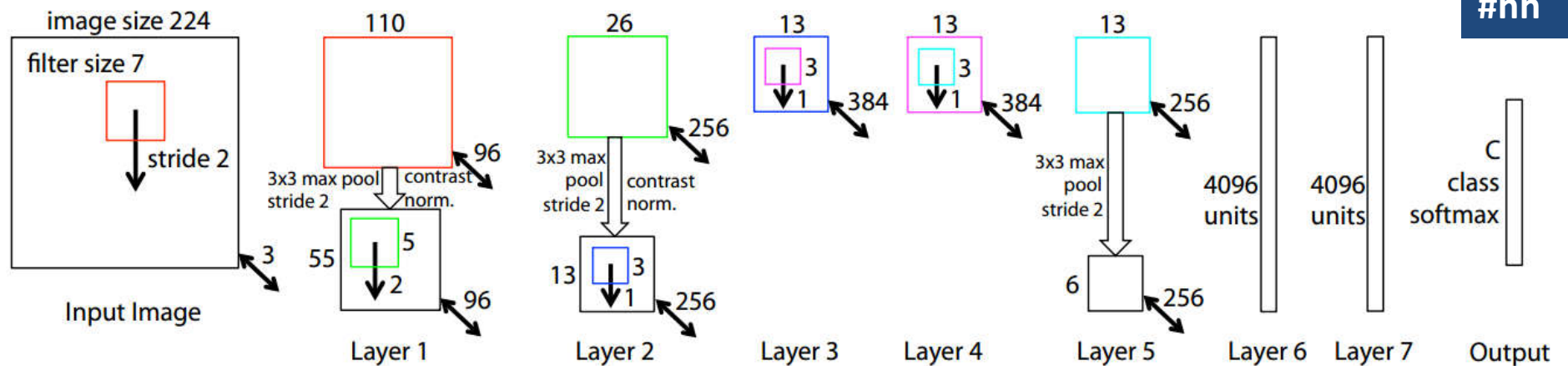
Krizhevsky, A., Sutskever, I. & Hinton, G. E. Imagenet classification with deep convolutional neural networks. In: Pereira, F., Burges, C. J. C., Bottou, L. & Weinberger, K. Q., eds. *Advances in neural information processing systems (NIPS 2012)*, 2012 Lake Tahoe. 1097-1105.



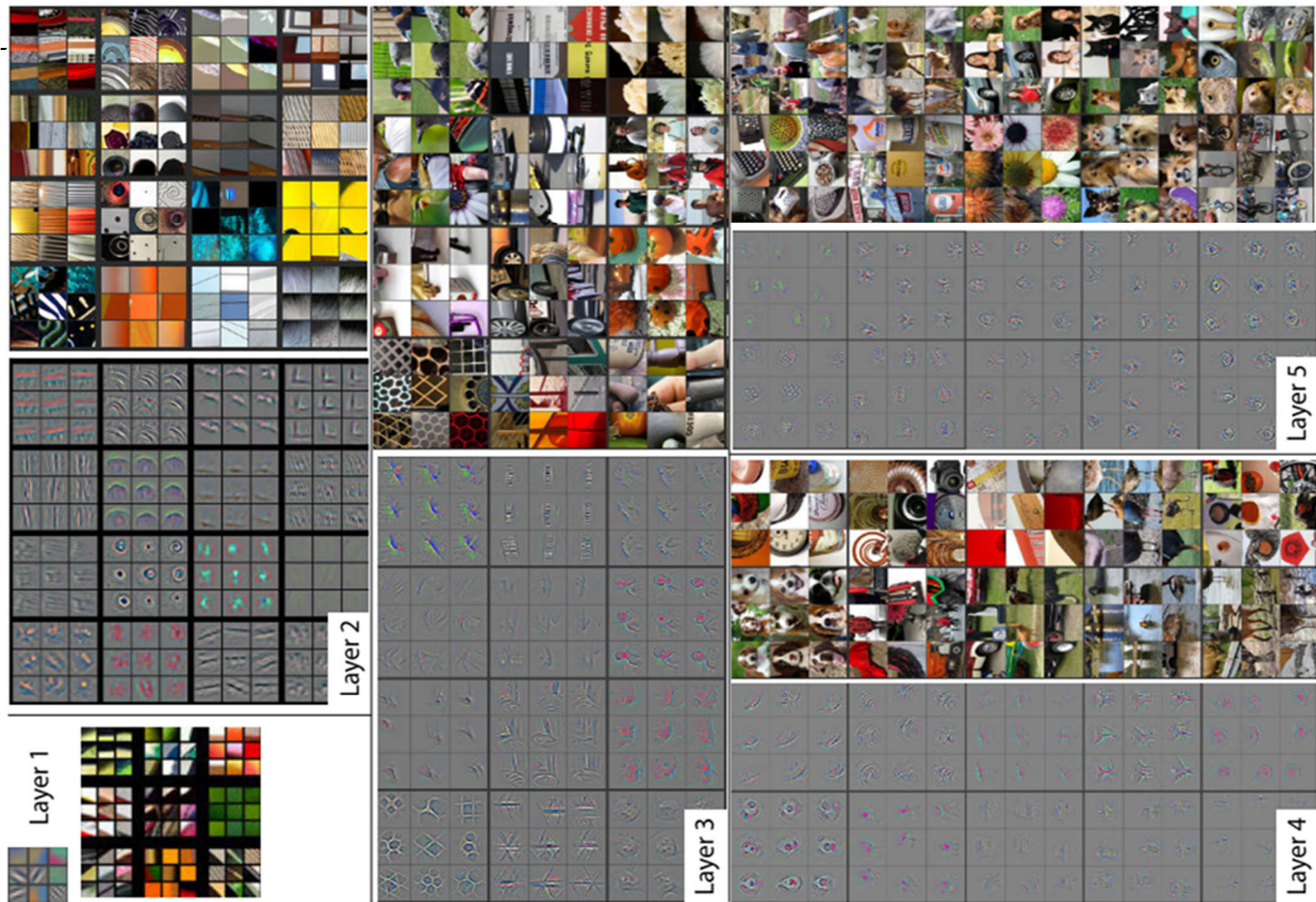
Matthew D. Zeiler & Rob Fergus 2013. Visualizing and Understanding Convolutional Networks. arXiv:1311.2901.

Visualizing a Conv Net with a De-Conv Net

#hh

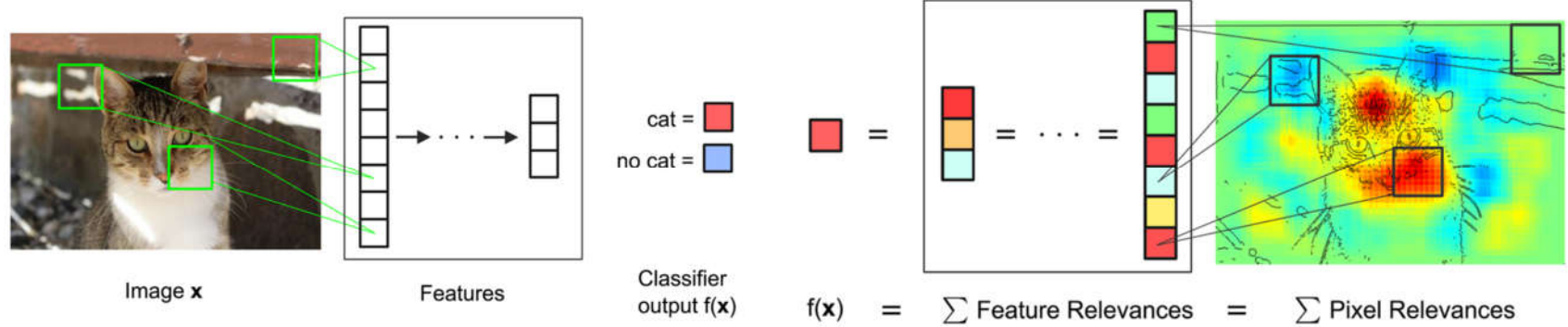


Matthew D. Zeiler & Rob Fergus 2014. Visualizing and understanding convolutional networks. In: D., Fleet, T., Pajdla, B., Schiele & T., Tuytelaars (eds.) ECCV, Lecture Notes in Computer Science LNCS 8689. Cham: Springer, pp. 818-833, doi:10.1007/978-3-319-10590-1_53.



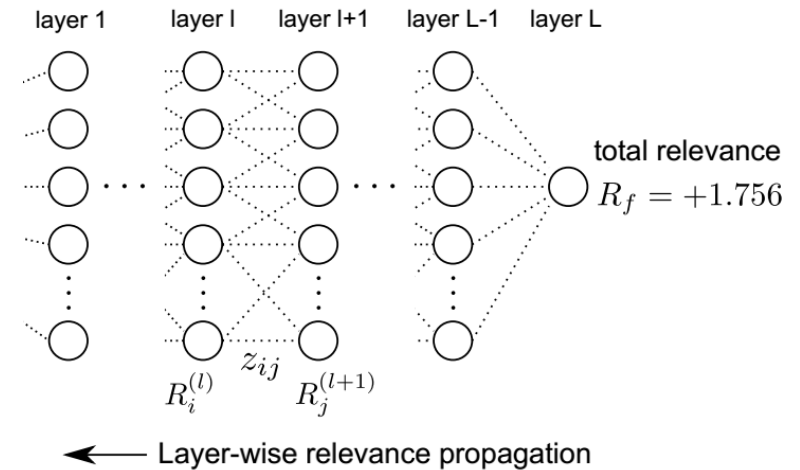
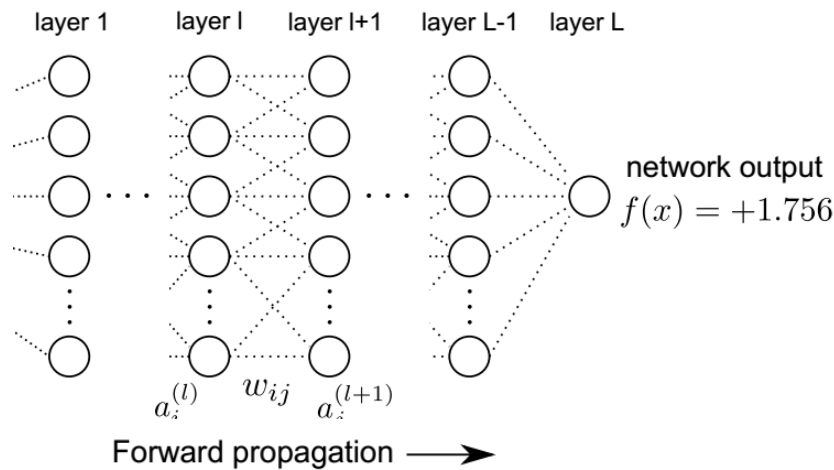
LRP Layer-Wise Relevance Propagation

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.

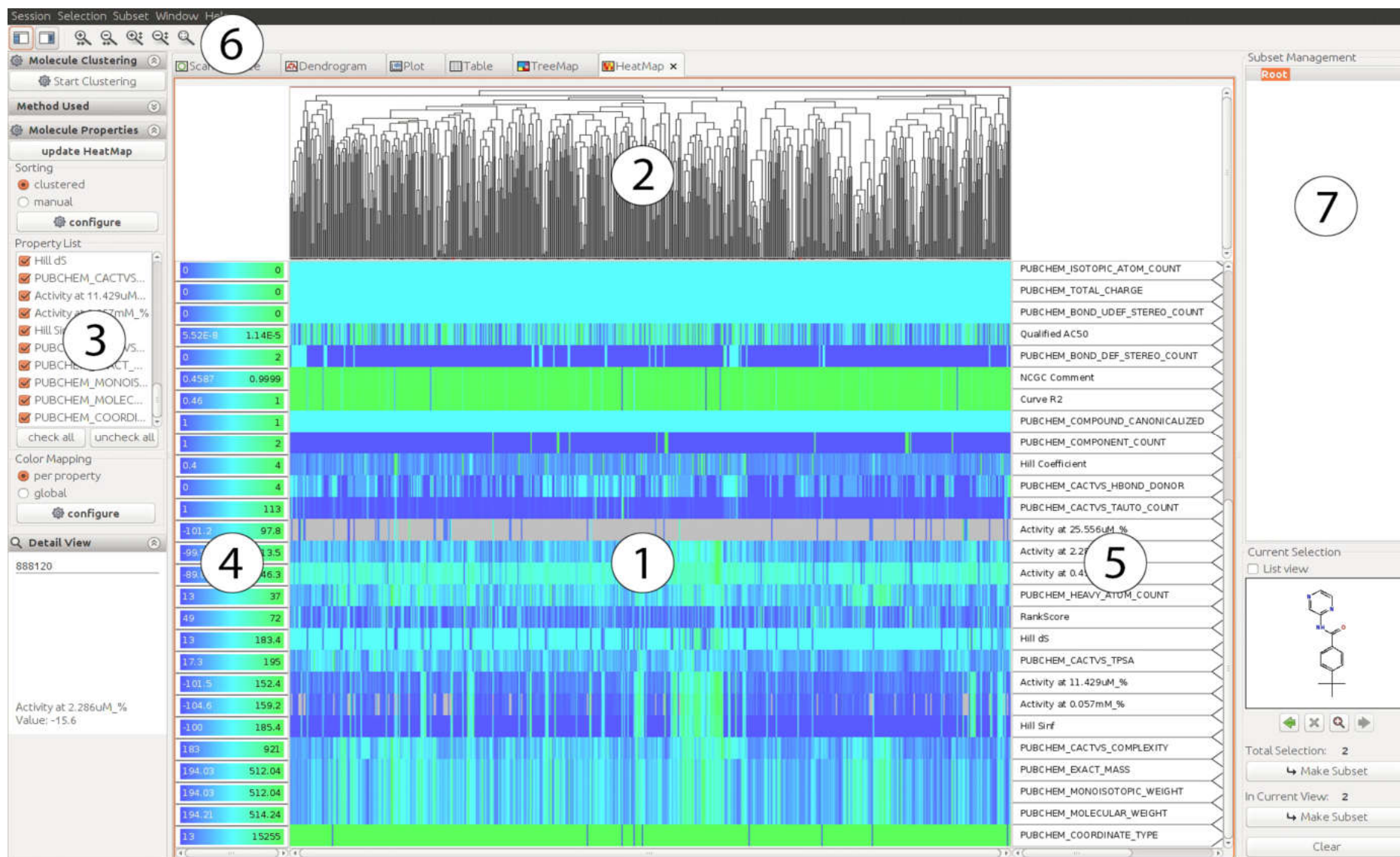


$$a_j^{(l+1)} = \sigma \left(\sum_i a_i^{(l)} w_{ij} + b_j^{(l+1)} \right)$$

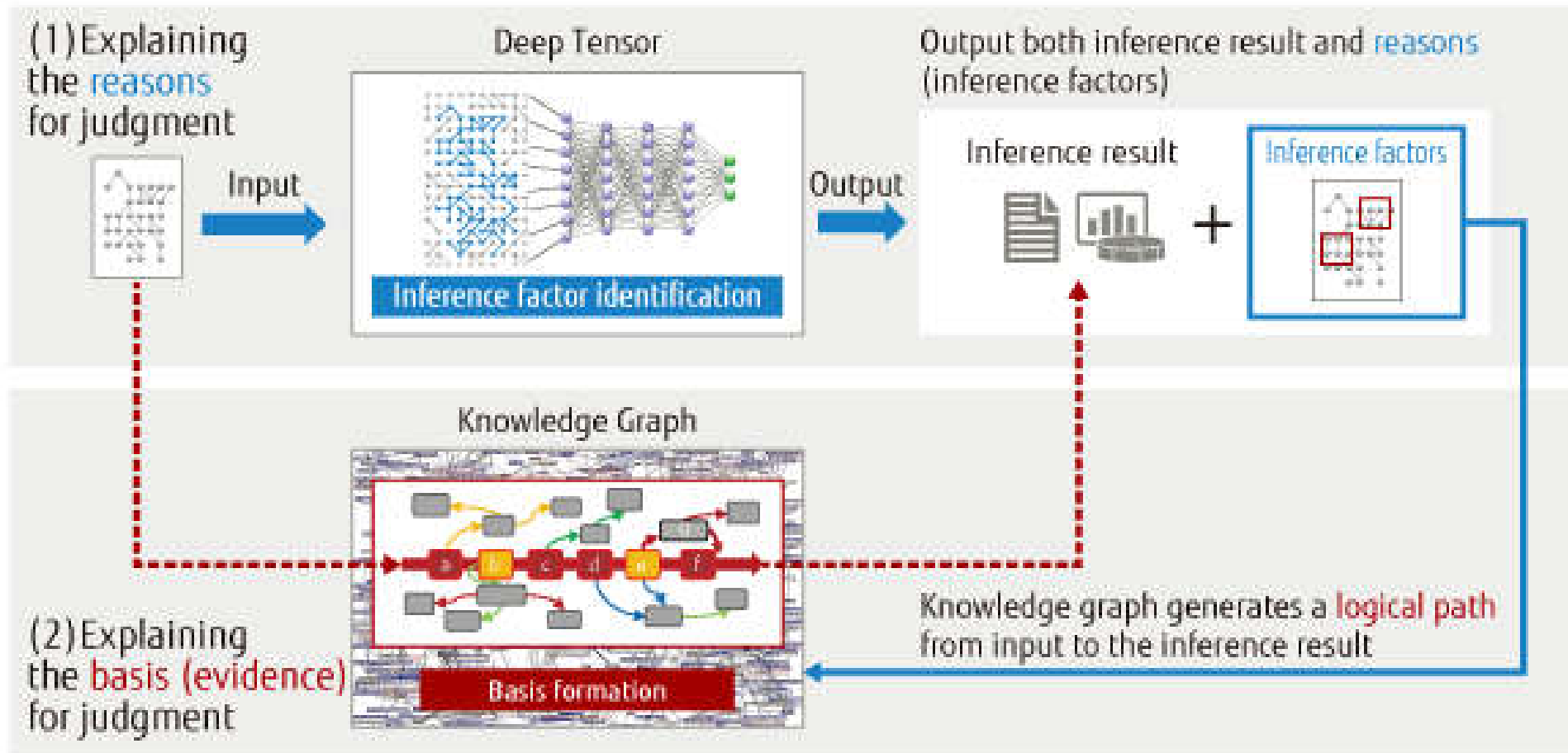
$$R_i^{(l)} = \sum_j \frac{z_{ij}}{\sum_{i'} z_{i'j}} R_j^{(l+1)}$$



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



Werner Sturm, Till Schaefer, Tobias Schreck, Andreas Holzinger & Torsten Ullrich. Extending the Scaffold Hunter Visualization Toolkit with Interactive Heatmaps In: Borgo, Rita & Turkay, Cagatay, eds. EG UK Computer Graphics & Visual Computing CGVC 2015, 2015 University College London (UCL). Euro Graphics (EG), 77-84, doi:10.2312/cgvc.20151247.



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. Cham: Springer, pp 295-303

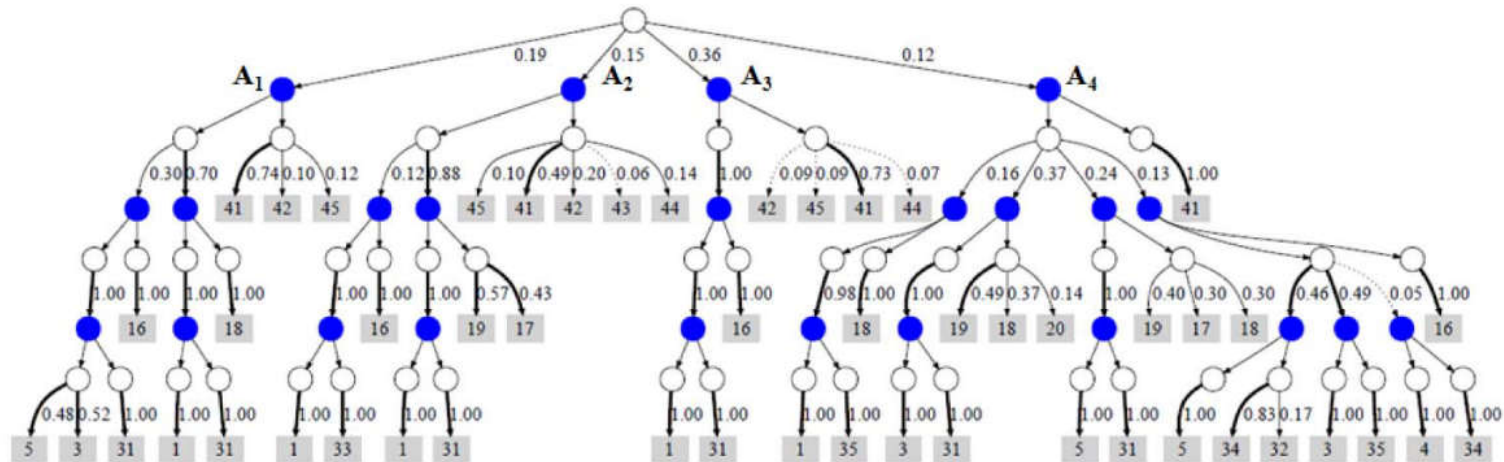
Stochastic AND-OR Templates for visual objects

#cc

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.



- **1** Computational approaches can find in \mathbb{R}^N what no human would be able to see
- **2** Complexity – reduction of search space
augment Human intelligence with AI & v.v.
- **3** Human expert can understand the **context**, need **effective** mapping $\mathbb{R}^N \rightarrow \mathbb{R}^2$
- **4** Black box approaches can not explain **WHY** a decision has been made ...

Multi-Task Learning ...

help to reduce **catastrophic forgetting**

Transfer learning ...

is not easy: learning to perform a task by exploiting knowledge acquired when solving previous tasks:

A solution to this problem would have major impact to AI research generally and ML specifically!

Multi-Agent-Hybrid Systems ...

collective intelligence and crowdsourcing
client-side federated machine learning – ensures
privacy, data protection, safety & security ...

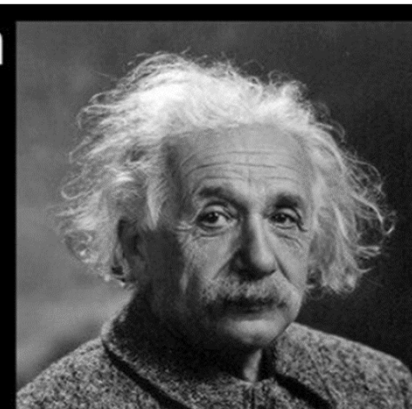
- Computers are fast, accurate and stupid,
- humans are slow, inaccurate and brilliant,
- **together** they are powerful beyond imagination

(Einstein never said that)

<https://www.benshoemate.com/2008/11/30/einstein-never-said-that>

„Das Dumme an Zitaten
aus dem Internet ist,
dass man nie weiß,
ob sie echt sind“

Albert Einstein





Ďakujem!