



Andreas Holzinger

185.A83 Machine Learning for Health Informatics

2019S, VU, 2.0 h, 3.0 ECTS

Lecture 02 – Dienstag, 19.03.2019



From Clinical Decision Support to Causal Reasoning and explainable AI

andreas.holzinger AT tuwien.ac.at

<https://hci-kdd.org/machine-learning-for-health-informatics-class-2019>

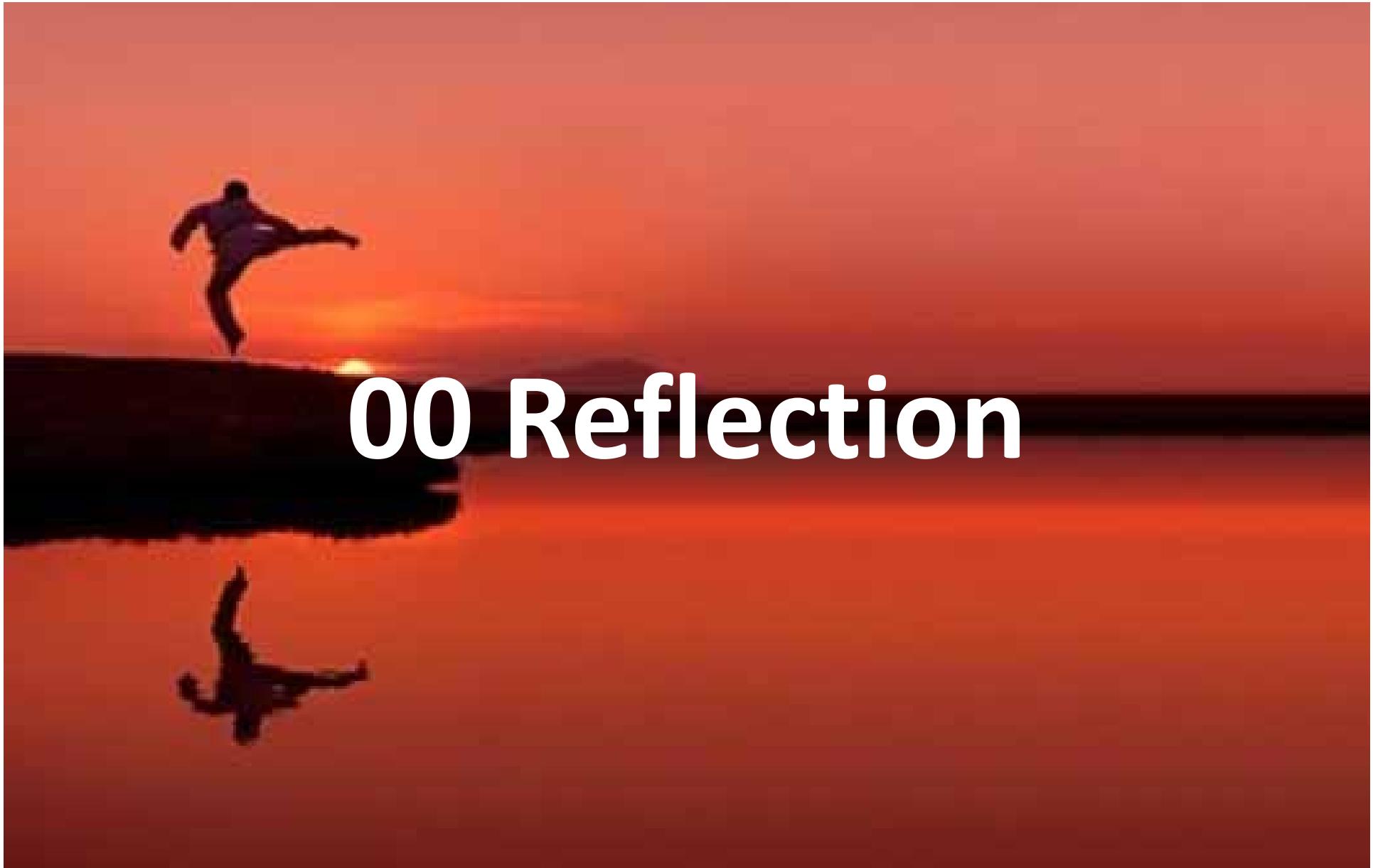


- Decision support system (DSS)
- MYCIN – Rule Based Expert System
- GAMUTS in Radiology
- Reasoning under uncertainty
- Example: Radiotherapy planning
- Example: Case-Based Reasoning
- Explainable Artificial intelligence
- Re-trace > Understand > Explain
- Transparency > Trust > Acceptance
- Fairness > Transparency > Accountability
- Causality > Usability
- (Some) Methods of Explainable AI

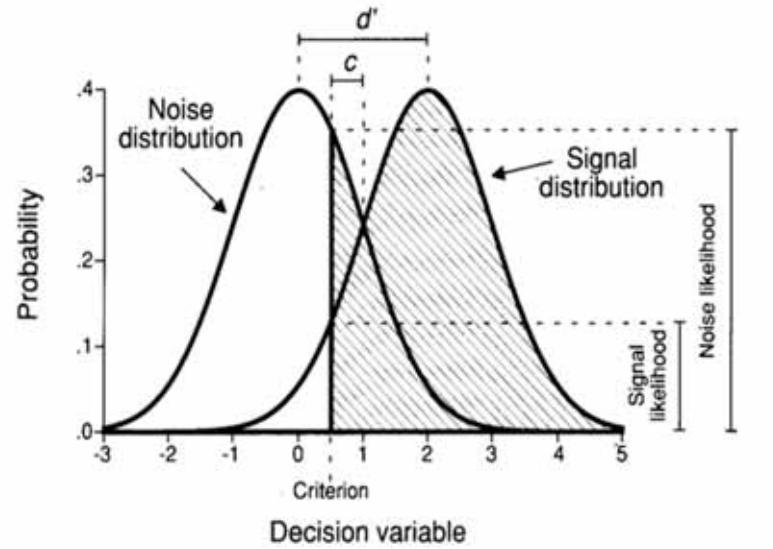
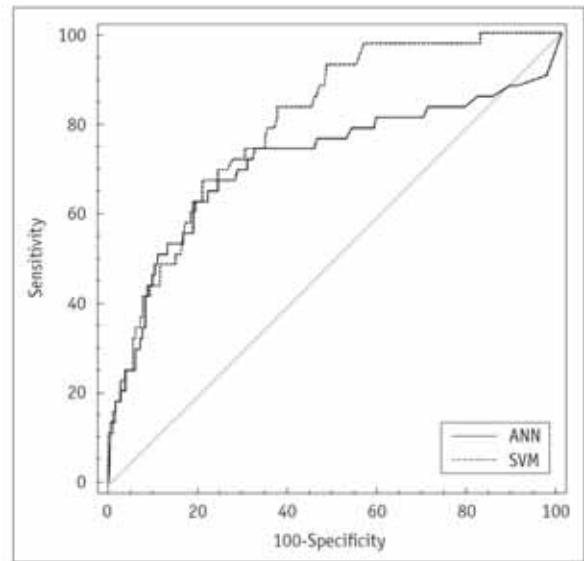
- **Causality** = fundamental relationship between cause and effect
- Causability = similar to the concept of usability the property of a human explanation
- **Case-based reasoning (CBR)** = process of solving new problems based on the solutions of similar past problems;
- **Certainty factor model (CF)** = a method for managing uncertainty in rule-based systems;
- **CLARION** = Connectionist Learning with Adaptive Rule Induction ON-line (CLARION) is a cognitive architecture that incorporates the distinction between implicit and explicit processes and focuses on capturing the interaction between these two types of processes. By focusing on this distinction, CLARION has been used to simulate several tasks in cognitive psychology and social psychology. CLARION has also been used to implement intelligent systems in artificial intelligence applications.
- **Clinical decision support (CDS)** = process for enhancing health-related decisions and actions with pertinent, organized clinical knowledge and patient information to improve health delivery;
- **Clinical Decision Support System (CDSS)** = expert system that provides support to certain reasoning tasks, in the context of a clinical decision;
- **Collective Intelligence** = shared group (symbolic) intelligence, emerging from cooperation/competition of many individuals, e.g. for consensus decision making;
- Counterfactual = relating to or expressing what has not happened or is not the case
- **Crowdsourcing** = a combination of "crowd" and "outsourcing" coined by Jeff Howe (2006), and describes a distributed problem-solving model; example for crowdsourcing is a public software beta-test;
- **Decision Making** = central cognitive process in every medical activity, resulting in the selection of a final choice of action out of several alternatives;
- **Decision Support System (DSS)** = is an IS including knowledge based systems to interactively support decision-making activities, i.e. making data useful;

- **DXplain** = a DSS from the Harvard Medical School, to assist making a diagnosis (clinical consultation), and also as an instructional instrument (education); provides a description of diseases, etiology, pathology, prognosis and up to 10 references for each disease;
- **Etiology** = in medicine (many) factors coming together to cause an illness (see causality)
- **Explainable AI** = Explainability = upcoming fundamental topic within recent AI; answering e.g. **why** a decision has been made
- **Expert-System** = emulates the decision making processes of a human expert to solve complex problems;
- **GAMUTS** in Radiology = Computer-Supported list of common/uncommon differential diagnoses;
- **ILIAD** = medical expert system, developed by the University of Utah, used as a teaching and testing tool for medical students in problem solving. Fields include Pediatrics, Internal Medicine, Oncology, Infectious Diseases, Gynecology, Pulmonology etc.
- **Interpretability** = there is no formal technical definition yet, but it is considered as a prerequisite for trust
- **MYCIN** = one of the early medical expert systems (Shortliffe (1970), Stanford) to identify bacteria causing severe infections, such as bacteremia and meningitis, and to recommend antibiotics, with the dosage adjusted for patient's body weight;
- **Reasoning** = cognitive (thought) processes involved in making medical decisions (clinical reasoning, medical problem solving, diagnostic reasoning);
- **Transparency** = opposite of opacity of black-box approaches, and connotes the ability to understand how a model works (that does not mean that it should always be understood, but that – in the case of necessity – it can be re-enacted

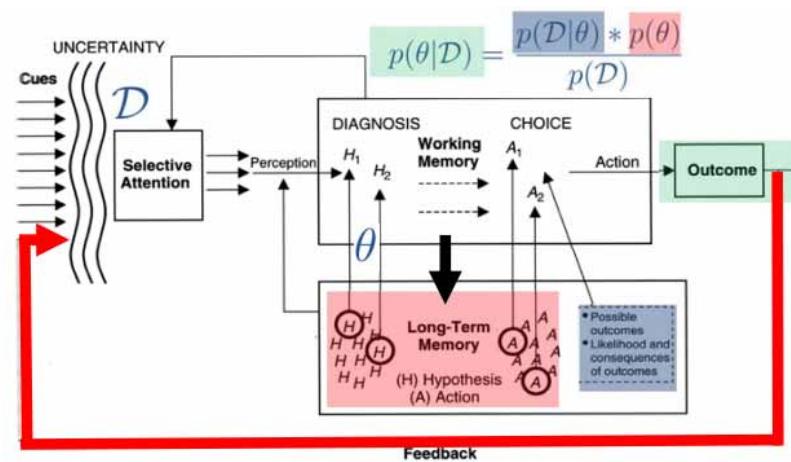
- 00 Reflection – follow-up from last lecture
- 01 Decision Support Systems (DSS)
- 02 History of DSS = History of AI
- 03 Example: Towards Personalized Medicine
- 04 Example: Case Based Reasoning (CBR)
- 05 Causal Reasoning
- 06 Explainability – Causability
- 07 (Some) Methods of Explainable AI



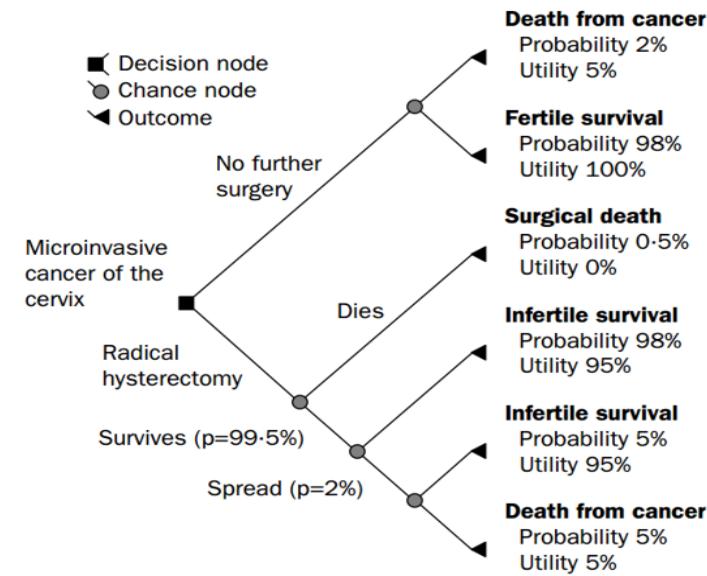
00 Reflection



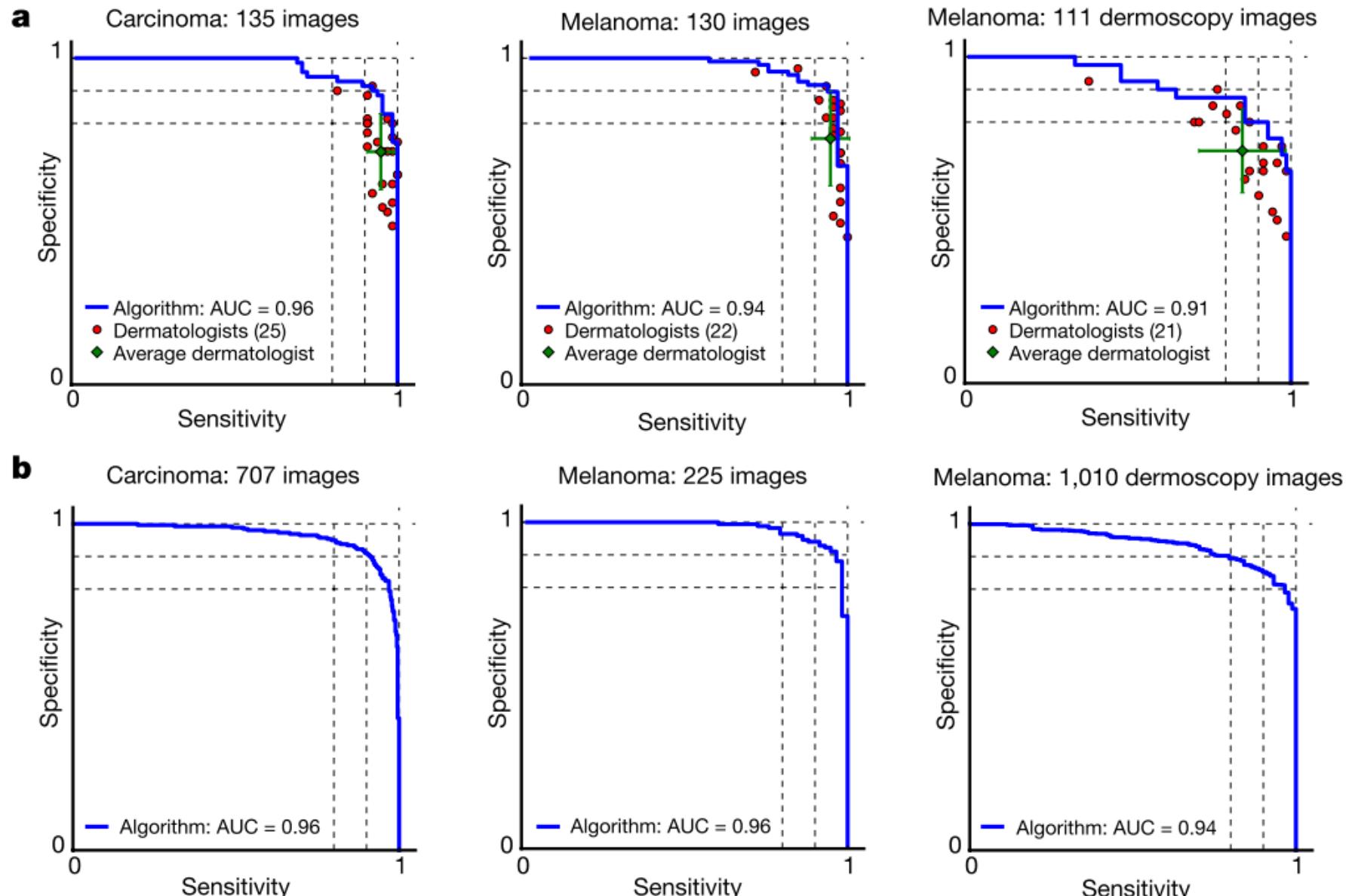
3



2



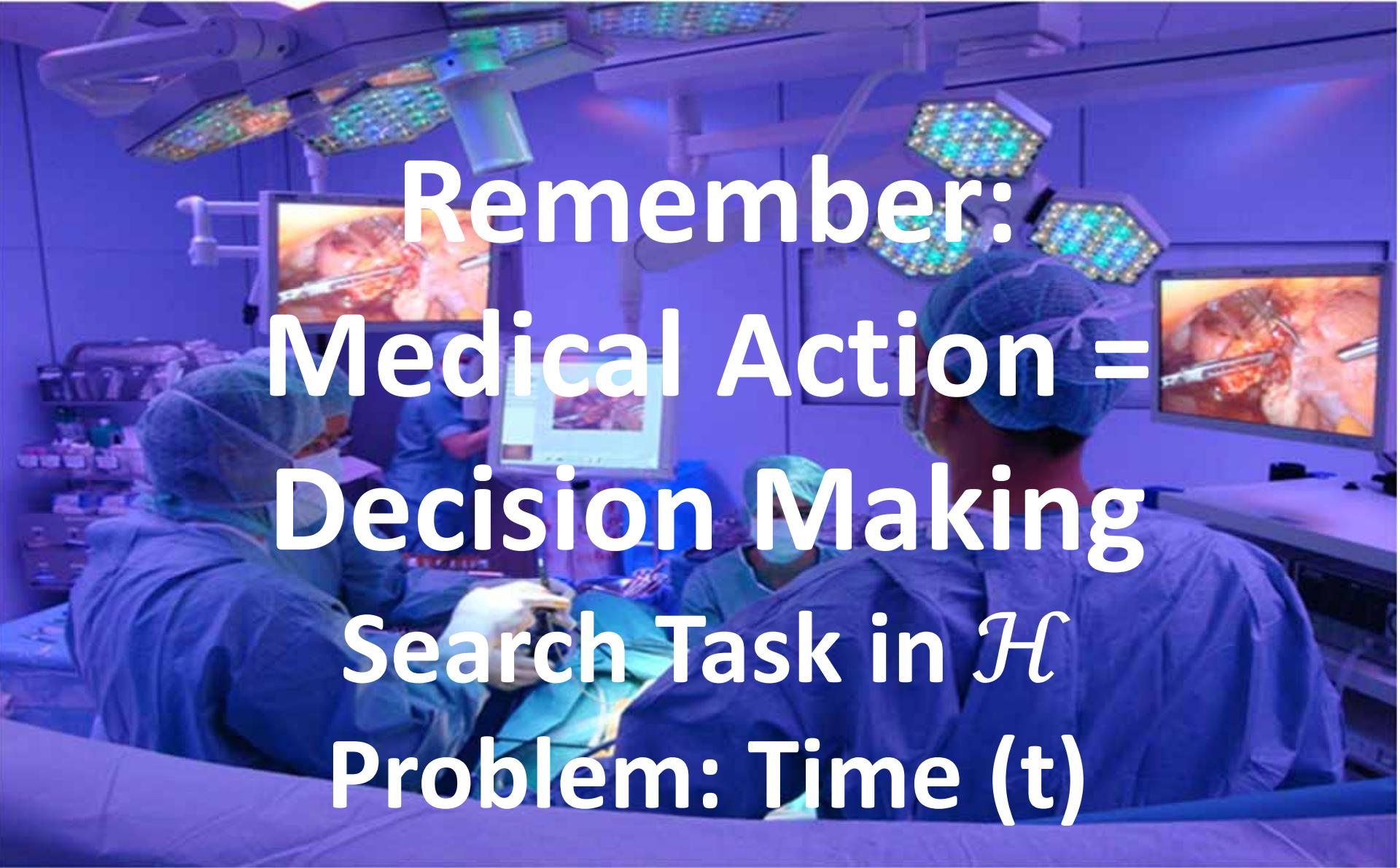
4



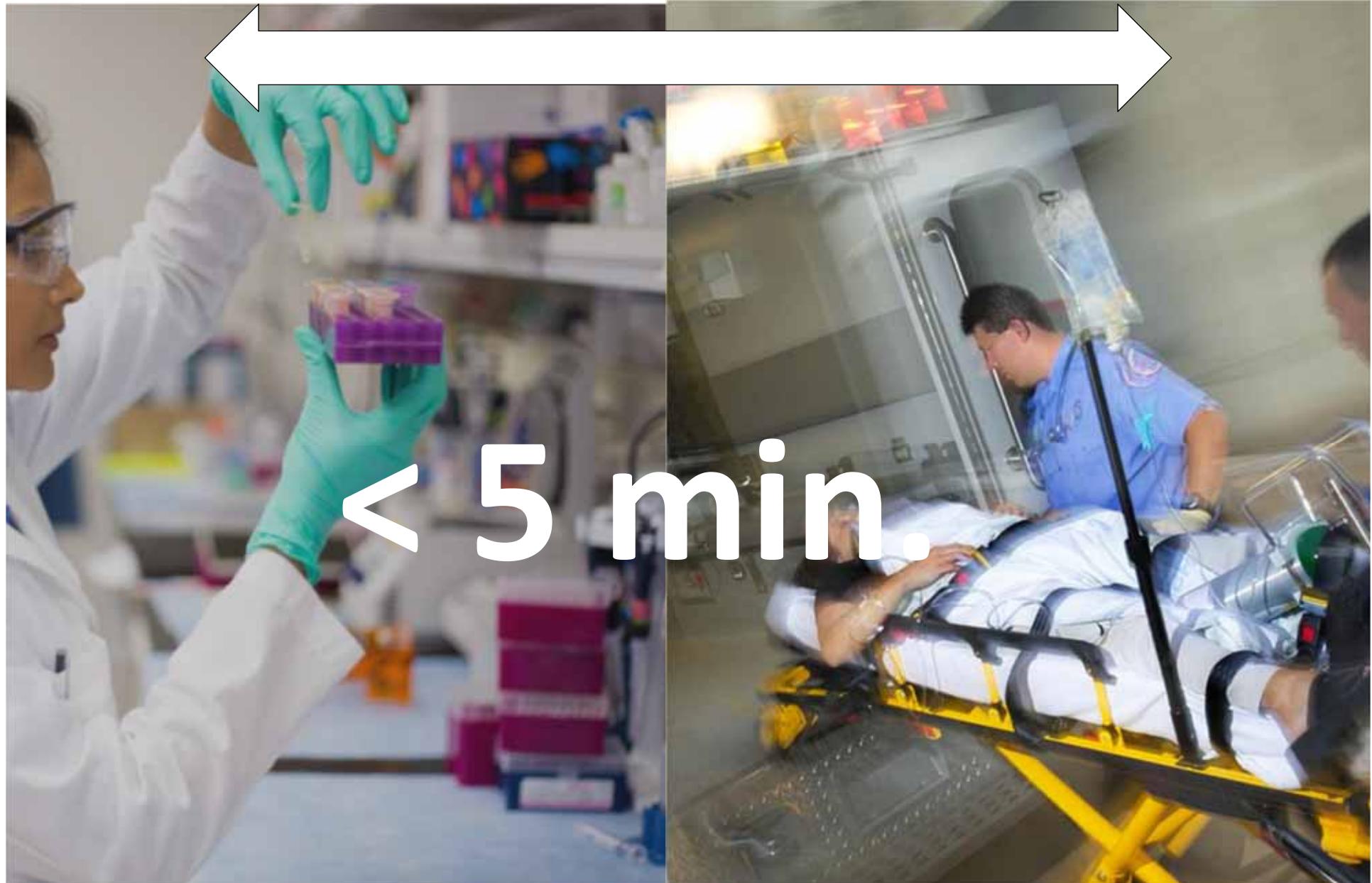
Andre Esteva, Brett Kuprel, Roberto A. Novoa, Justin Ko, Susan M. Swetter, Helen M. Blau & Sebastian Thrun 2017.
Dermatologist-level classification of skin cancer with deep neural networks. Nature, 542, (7639), 115-118

- Remember: Medicine is an complex application domain – dealing most of the time with **probable information!**
- Some challenges include:
- (a) defining hospital system architectures in terms of generic tasks such as diagnosis, therapy planning and monitoring to be executed for (b) medical reasoning in (a);
- (c) patient information management with (d) minimum uncertainty.
- Other challenges include: (e) knowledge acquisition and encoding, (f) human-ai interface and ai-interaction; and (g) system integration into existing clinical legacy and proprietary environments, e.g. the enterprise hospital information system; to mention only a few.

01 Decision Support Systems



Remember:
Medical Action =
Decision Making
Search Task in \mathcal{H}
Problem: Time (t)



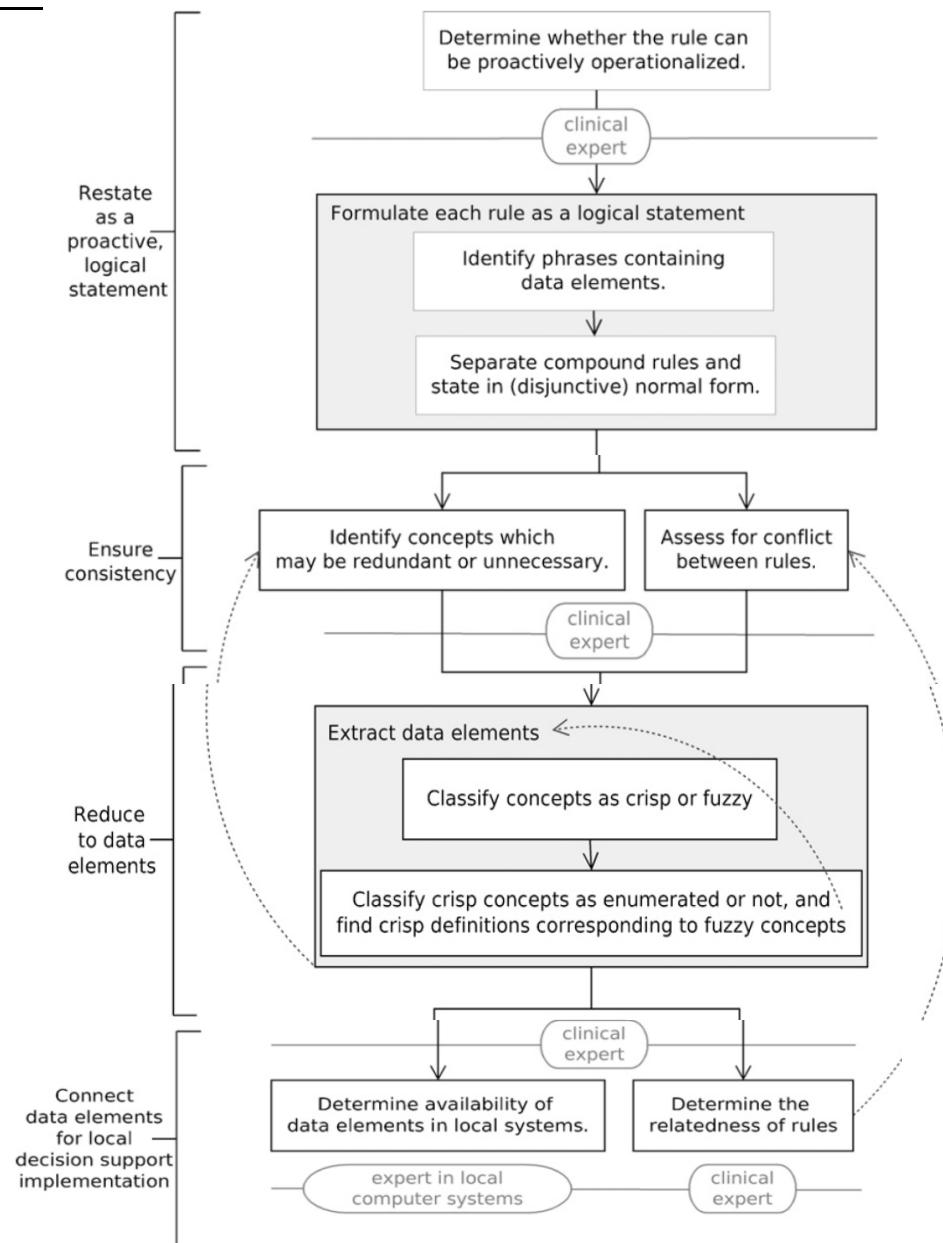


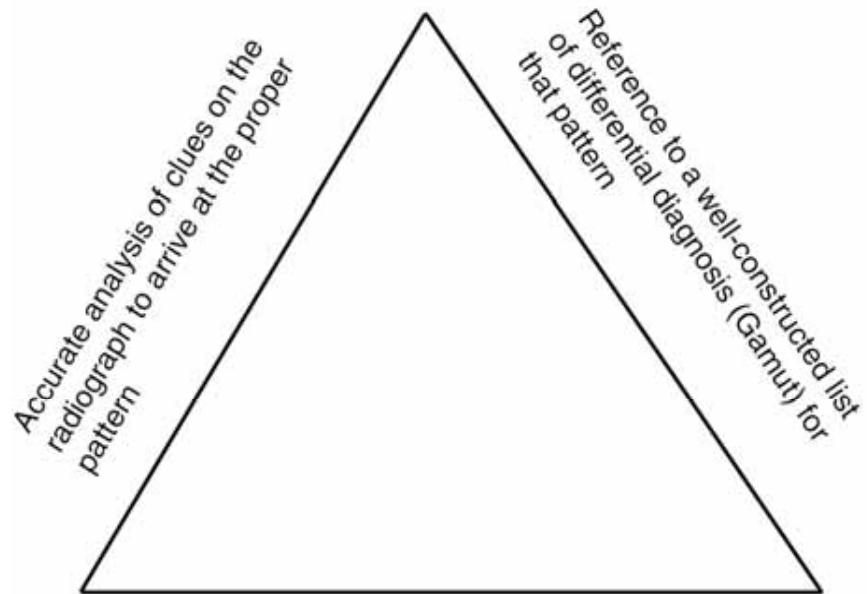
Source: Cisco (2008).
Cisco Health Presence
Trial at Aberdeen Royal
Infirmary in Scotland

- 400 BC Hippocrates (460-370 BC), father of western medicine:
 - A medical record should accurately reflect the course of a disease
 - A medical record should indicate the probable cause of a disease
- 1890 William Osler (1849-1919), father of modern western medicine
 - **Medicine is a science of uncertainty and an art of probabilistic decision making**
- Today
 - Prediction models are based on data features, patient health status is modelled as high-dimensional feature vectors ...

- Clinical guidelines are **systematically** developed documents to assist doctors and patient decisions about appropriate care;
- In order to build DS, based on a guideline, it is **formalized** (transformed from natural language to a logical algorithm), and
- **implemented** (using the algorithm to program a DSS);
- To increase the quality of care, they must be linked to a process of care, for example:
 - “80% of diabetic patients should have an HbA1c below 7.0” could be linked to processes such as:
 - “All diabetic patients should have an annual HbA1c test” and
 - “Patients with values over 7.0 should be rechecked within 2 months.”
- **Condition-action rules** specify one or a few conditions which are linked to a specific action, in contrast to narrative guidelines which describe a series of branching or iterative decisions unfolding over time.
- Narrative guidelines and clinical rules are two ends of a continuum of clinical care standards.

Medlock, S., Oondo, D.,
Eslami, S., Askari, M.,
Wierenga, P., de Rooij, S. E. &
Abu-Hanna, A. (2011) LERM
(Logical Elements Rule
Method): A method for
assessing and formalizing
clinical rules for decision
support. *International Journal
of Medical Informatics*, 80, 4,
286-295.





Reeder, M. M. & Felson, B. 2003.
Reeder and Felson's gamuts in radiology: comprehensive lists of roentgen differential diagnosis, New York, Springer Verlag.

Gamut F-137

PHRENIC NERVE PARALYSIS OR DYSFUNCTION

COMMON

1. Iatrogenic (eg, surgical injury; chest tube; therapeutic avulsion or injection; subclavian vein puncture)
2. Infection (eg, tuberculosis; fungus disease; abscess)
3. Neoplastic invasion or compression (esp. carcinoma of lung)

UNCOMMON

1. Aneurysm_g, aortic or other
2. Birth trauma (Erb's palsy)
3. Herpes zoster
4. Neuritis, peripheral (eg, diabetic neuropathy)
5. Neurologic disease_g (eg, hemiplegia; encephalitis; polio; Guillain-Barré S.)
6. Pneumonia
7. Trauma

Reference

1. Prasad S, Athreya BH: Transient paralysis of the phrenic nerve associated with head injury. JAMA 1976;236:2532–2533

REEDER AND FELSON'S

GAMUTS IN RADIOLOGY

GAMUT G-25

EROSIVE GASTRITIS*

COMMON

1. Acute gastritis (eg, alcohol abuse)
2. Crohn's disease  
3. Drugs (eg, aspirin  
5. Idiopathic
6. [Normal areae gastricae 
7. Peptic ulcer; hyperacidity

UNCOMMON

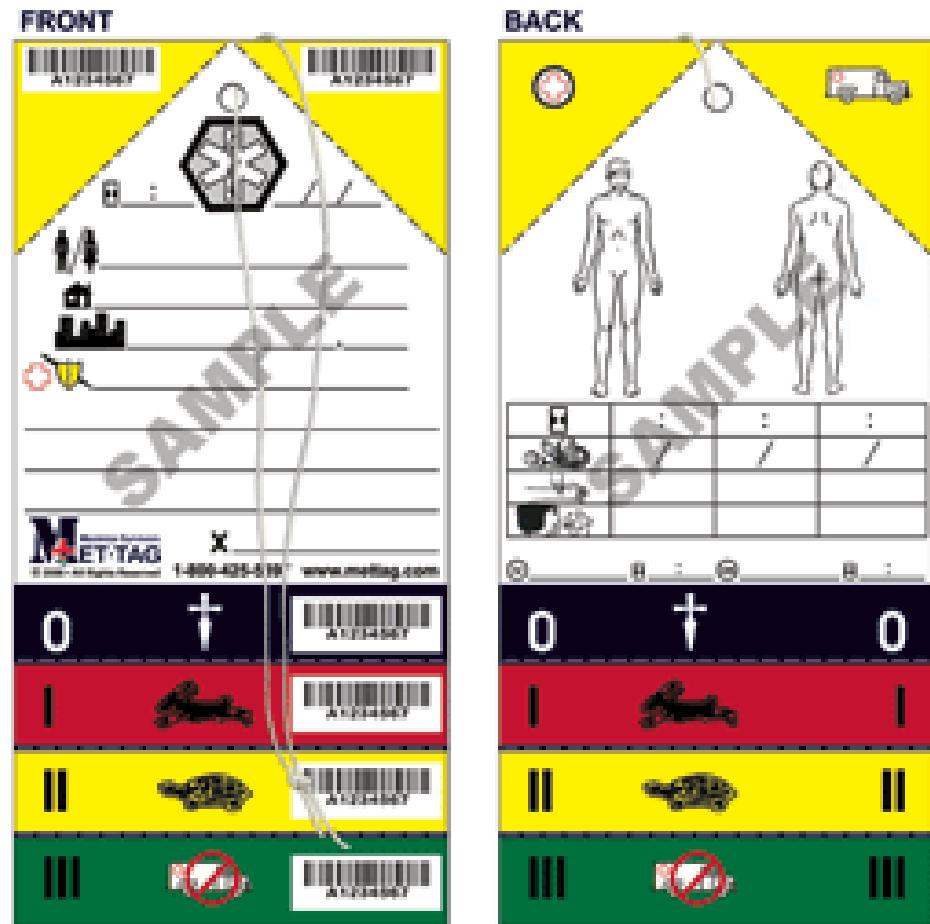
1. Corrosive gastritis 
2. *Cryptosporidium* antritis
3. [Lymphoma]
4. Opportunistic infection (eg, candidiasis {moniliasis} ; multiple endocrine neoplasia (MEN) S.

* Superficial erosions or aphthoid ulcerations seen especially with double contrast technique.

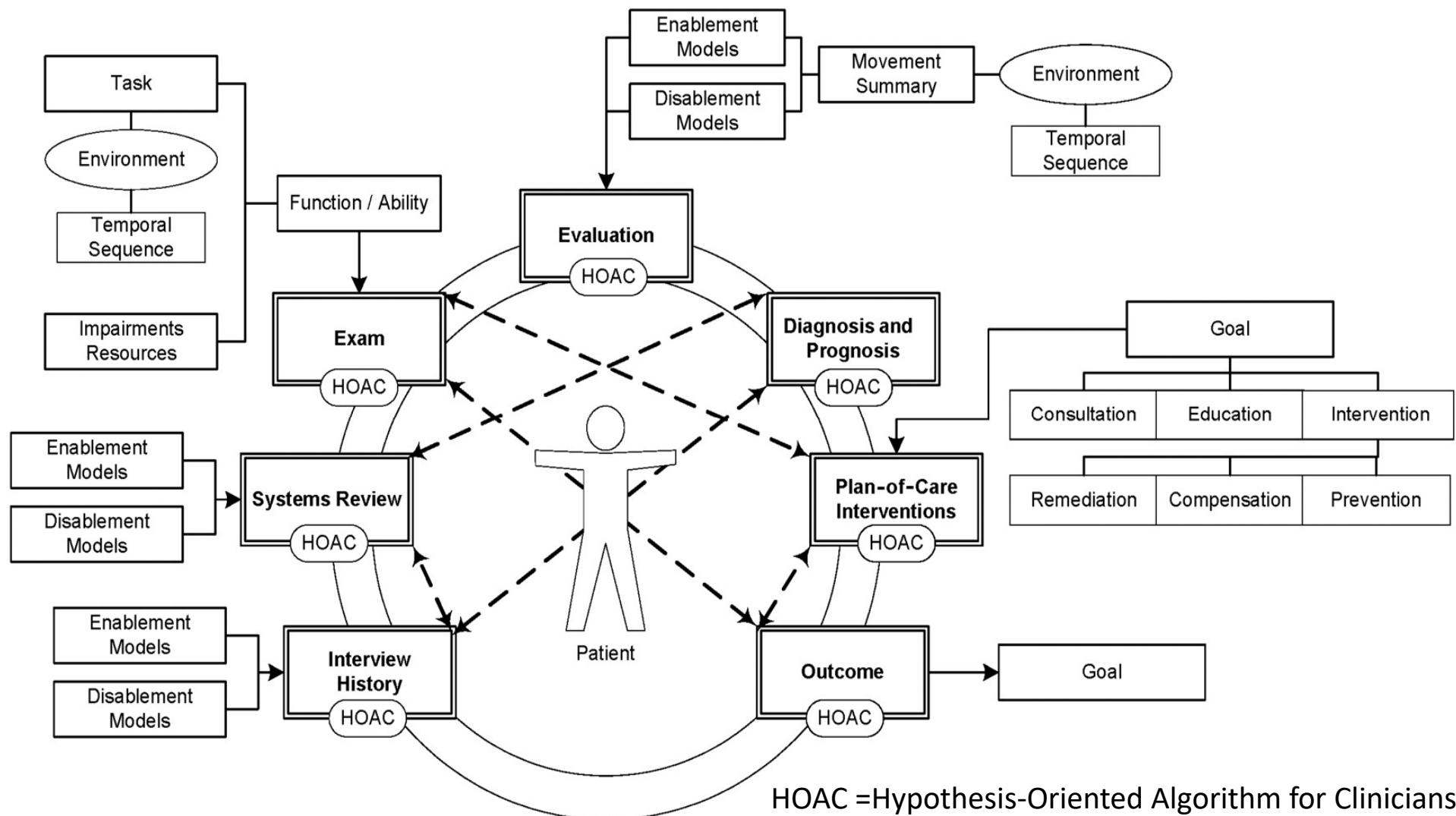
[] This condition does not actually cause the gamuted imaging finding, but can produce imaging changes that simulate it.

<http://rfs.acr.org/gamuts/data/G-25.htm>

No 565959		No 565959	
EVAC-AIR® TRIAGE TAG		CONTAMINATION: NO YES	
Respiration	Normal	Normal	Normal
Personality	<2 SEC	2-5 SEC	5-10 SEC
Mental Status	Calm	Alert	Confused
Wear GLOVES	DISINFECTED	UNDISINFECTED	
Treatment	Pulse	B/P	Respiratory
Treatment			
Time		Drug Taken	
Allergies		Prescription Medication	
Personal Information			
Name:		Address:	
City:		St:	Zip:
Phone:		Weight:	
Male		Female	
Age:		Height:	
Occupation:		Destination:	
DECEASED		DECEASED	
IMMEDIATE		IMMEDIATE	
DELAYED		DELAYED	
MINOR		MINOR	



Iserson, K. V. & Moskop, J. C. 2007. Triage in Medicine, Part I: Concept, History, and Types. Annals of Emergency Medicine, 49, (3), 275-281.



Schenkman, M., Deutsch, J. E. & Gill-Body, K. M. (2006) An Integrated Framework for Decision Making in Neurologic Physical Therapist Practice. *Physical Therapy*, 86, 12, 1681-1702.

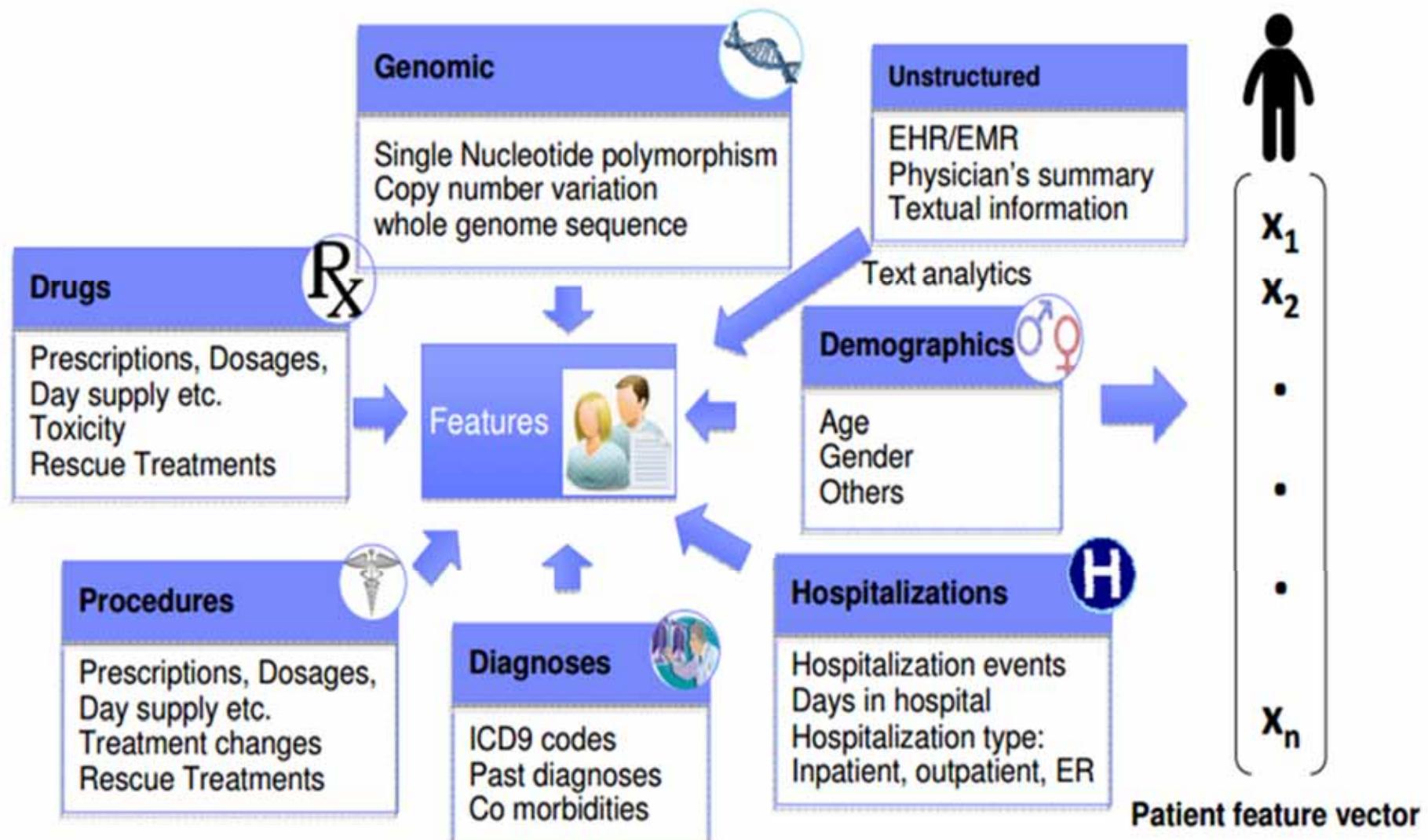
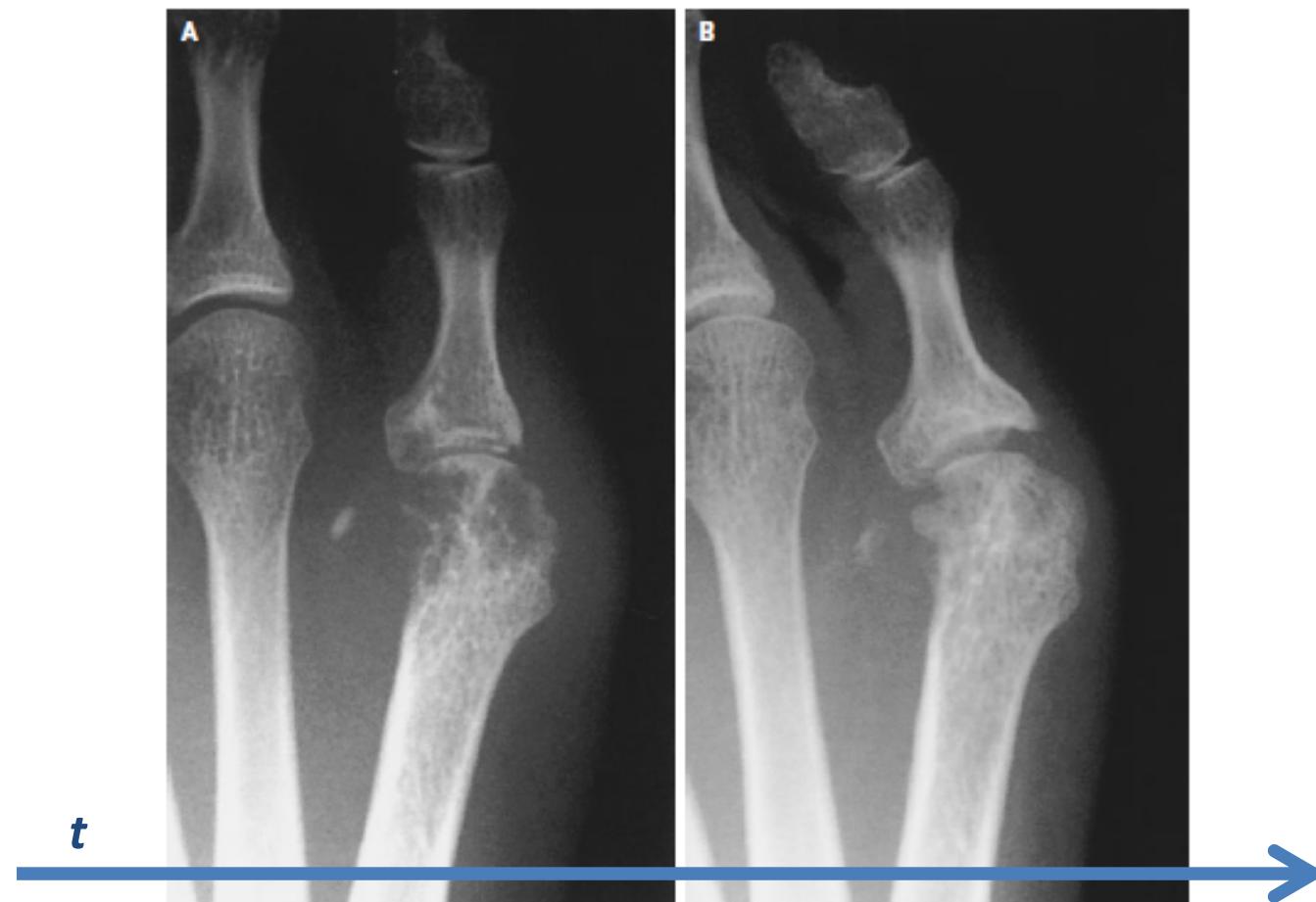


Image credit to Michal Rosen-Zvi

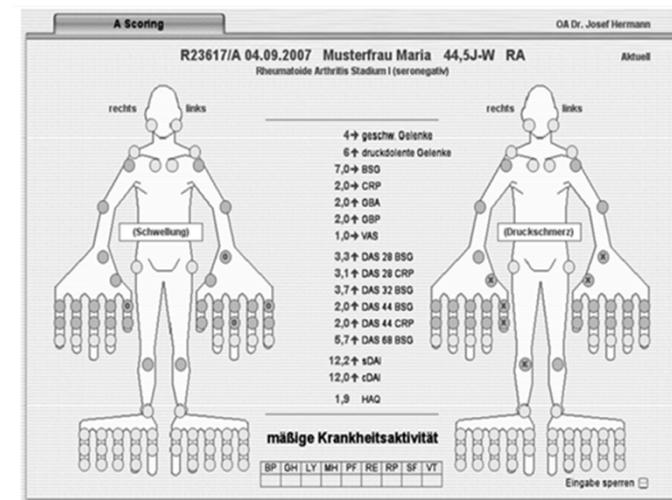
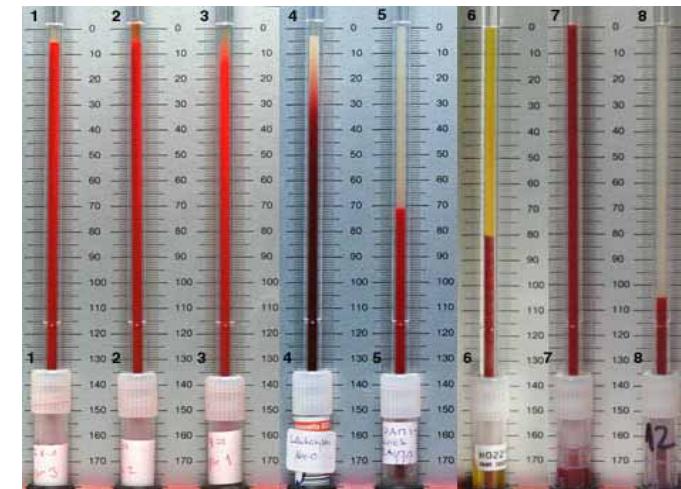


Chao, J., Parker, B. A. & Zvaifler, N. J. (2009) Accelerated Cutaneous Nodulosis Associated with Aromatase Inhibitor Therapy in a Patient with Rheumatoid Arthritis. *The Journal of Rheumatology*, 36, 5, **1087-1088**.

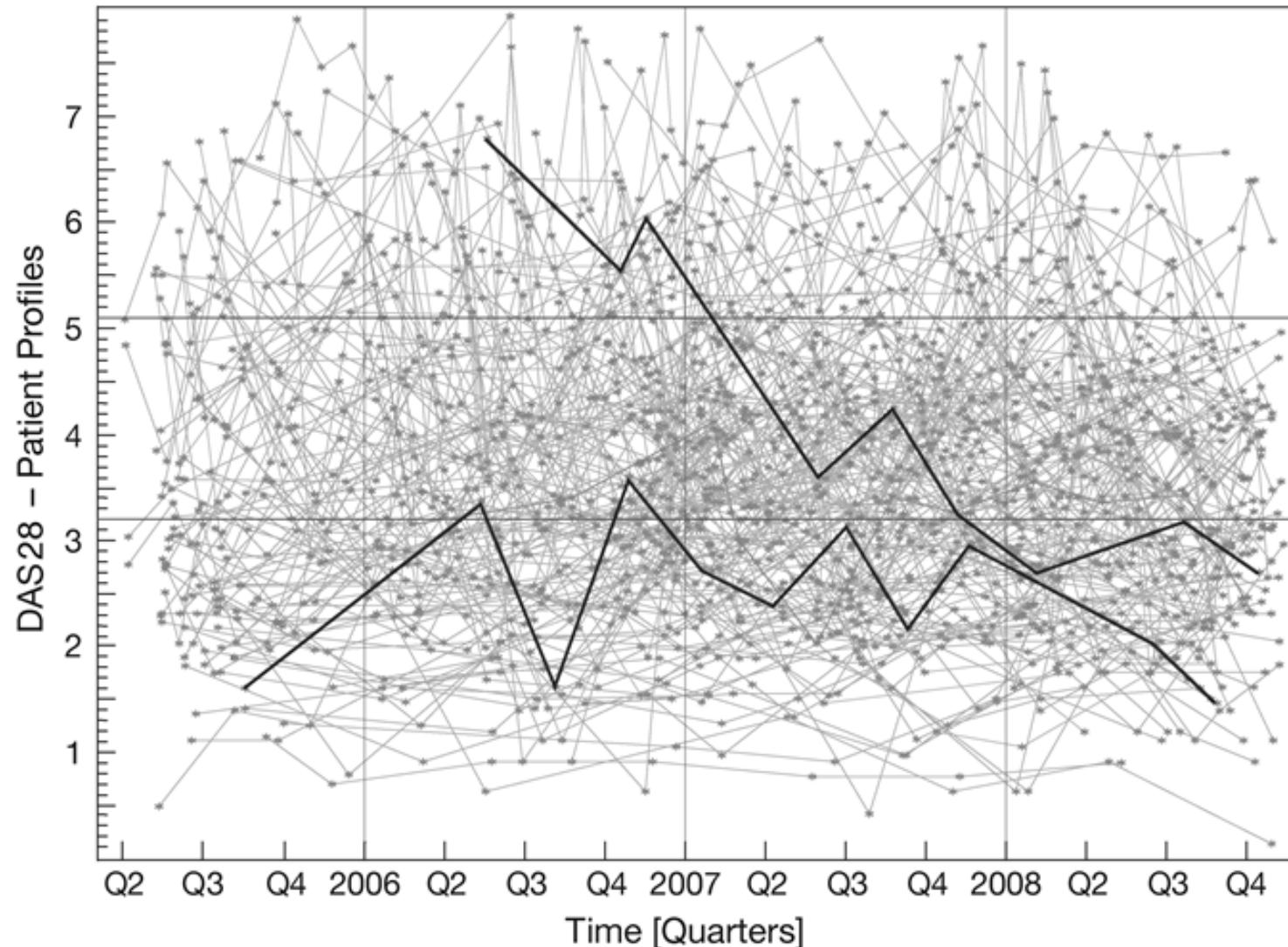


Ikari, K. & Momohara, S. (2005) Bone Changes in Rheumatoid Arthritis.
New England Journal of Medicine, 353, 15, e13.

- 50+ Patients per day ~ 5000 data points per day ...
- Aggregated with specific scores (Disease Activity Score, DAS)
- Current patient status is related to previous data
- = convolution over time
- ⇒ **time-series data**



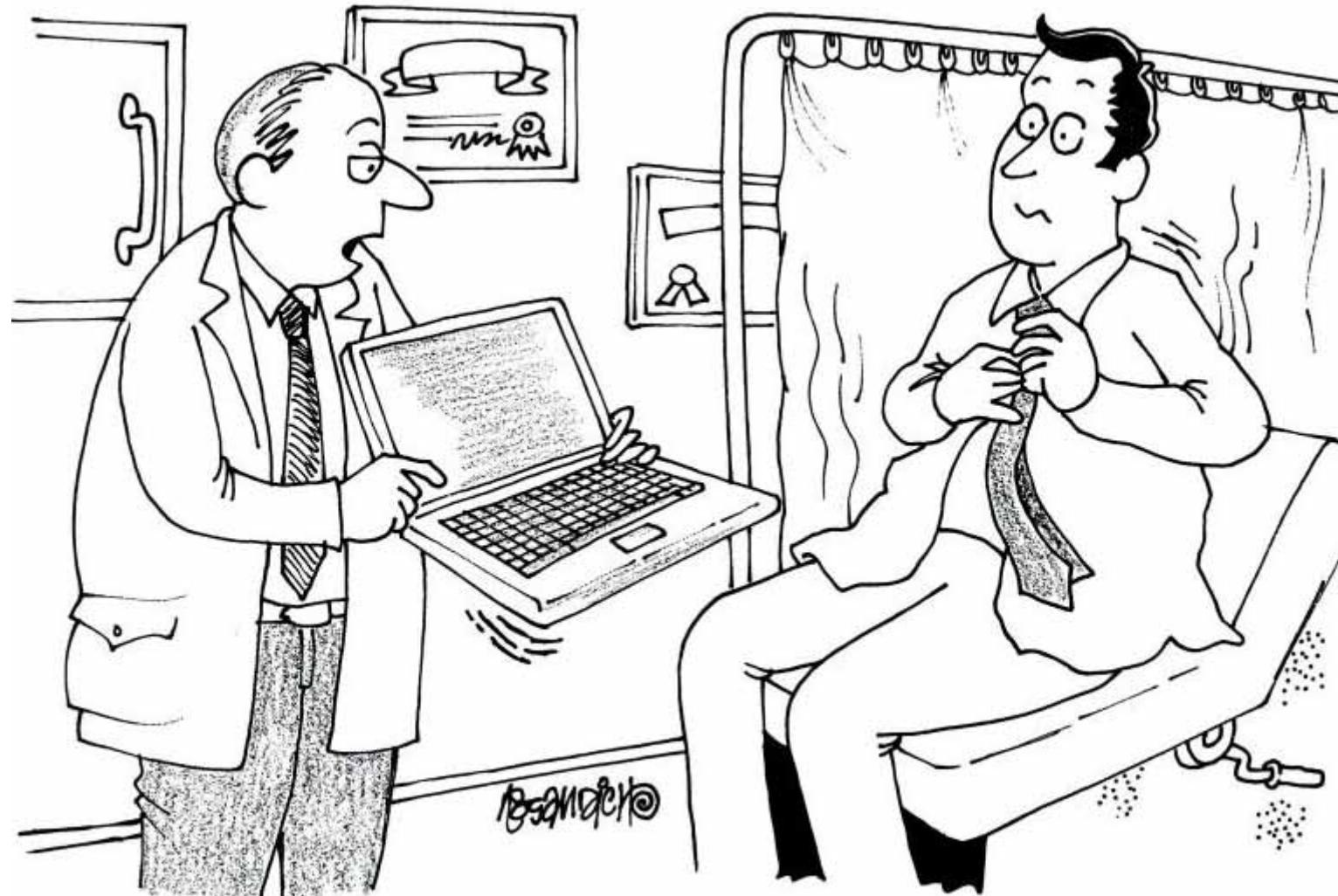
Simonic, K. M., Holzinger, A., Bloice, M. & Hermann, J. (2011). *Optimizing Long-Term Treatment of Rheumatoid Arthritis with Systematic Documentation*. Pervasive Health - 5th International Conference on Pervasive Computing Technologies for Healthcare, Dublin, IEEE, 550-554.



Simonic, K. M., Holzinger, A., Bloice, M. & Hermann, J. (2011). *Optimizing Long-Term Treatment of Rheumatoid Arthritis with Systematic Documentation*. Pervasive Health - 5th International Conference on Pervasive Computing Technologies for Healthcare, Dublin, IEEE, 550-554.



Can Computers help doctors to make better decisions?



"If you want a second opinion, I'll ask my computer."

<http://biomedicalcomputationreview.org/content/clinical-decision-support-providing-quality-healthcare-help-computer>



- **Type 1 Decisions:** related to the diagnosis, i.e. AI/ML is used to assist in diagnosing a disease on the basis of the individual patient data. Questions include:
 - What is the probability that this patient has a myocardial infarction on the basis of given data (patient history, ECG, ...)?
 - What is the probability that this patient has acute appendicitis, given the signs and symptoms concerning abdominal pain?
- **Type 2 Decisions:** related to therapy, i.e. AI/ML is used to select the best therapy on the basis of clinical evidence, e.g.:
 - What is the best therapy for patients of age x and risks y, if an obstruction of more than z % is seen in the left coronary artery?
 - What amount of insulin should be prescribed for a patient during the next 5 days, given the blood sugar levels and the amount of insulin taken during the recent weeks?

Bemmel, J. H. V. & Musen, M. A. 1997. *Handbook of Medical Informatics*, Heidelberg, Springer.



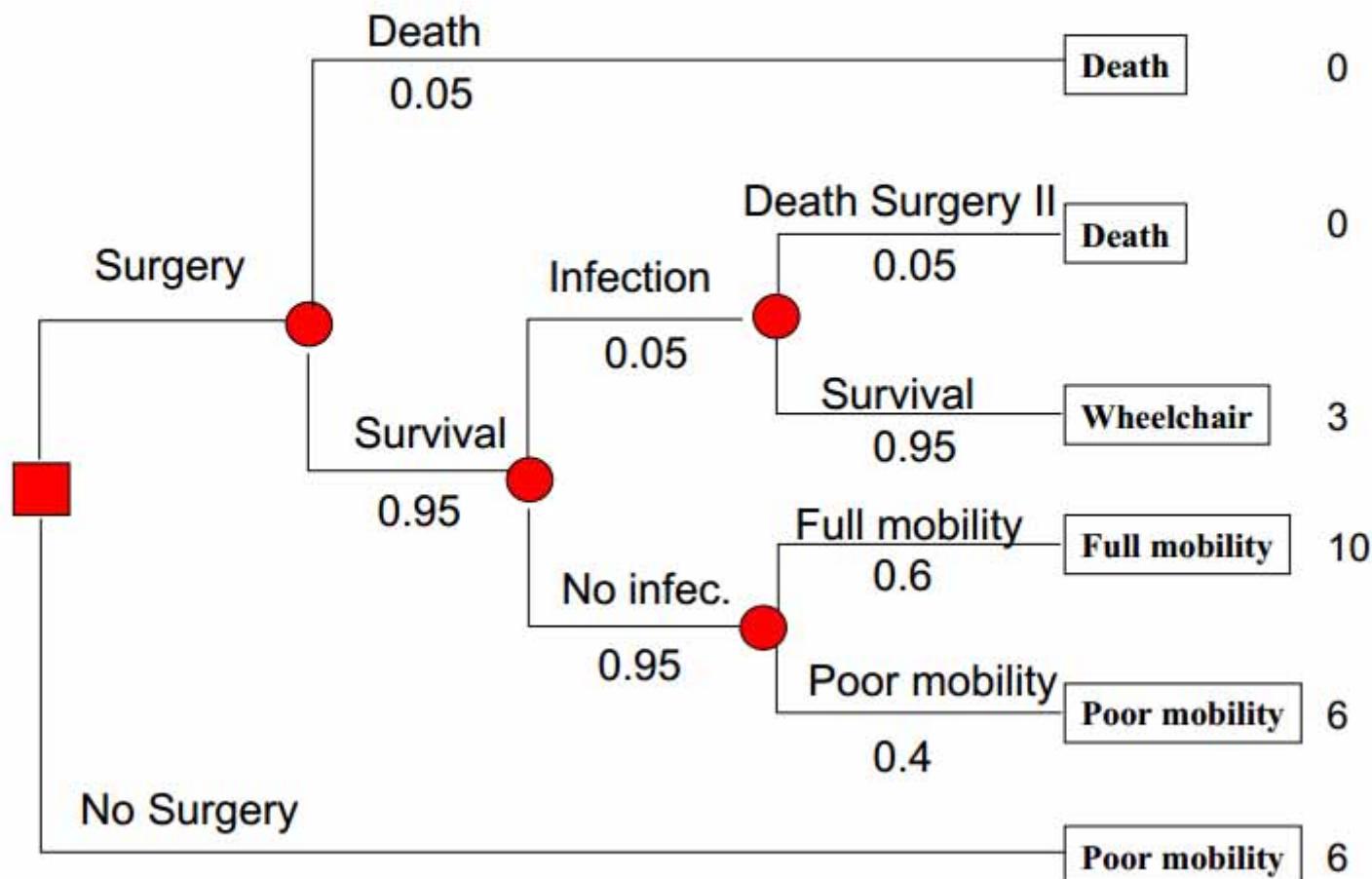
- Example of a Decision Problem
- Soccer player considering knee surgery
- Uncertainties:
- Success: recovering full mobility
- Risks: infection in surgery (if so, needs another surgery and may lose more mobility)
- Survival chances of surgery

Harvard-MIT Division of Health Sciences and Technology

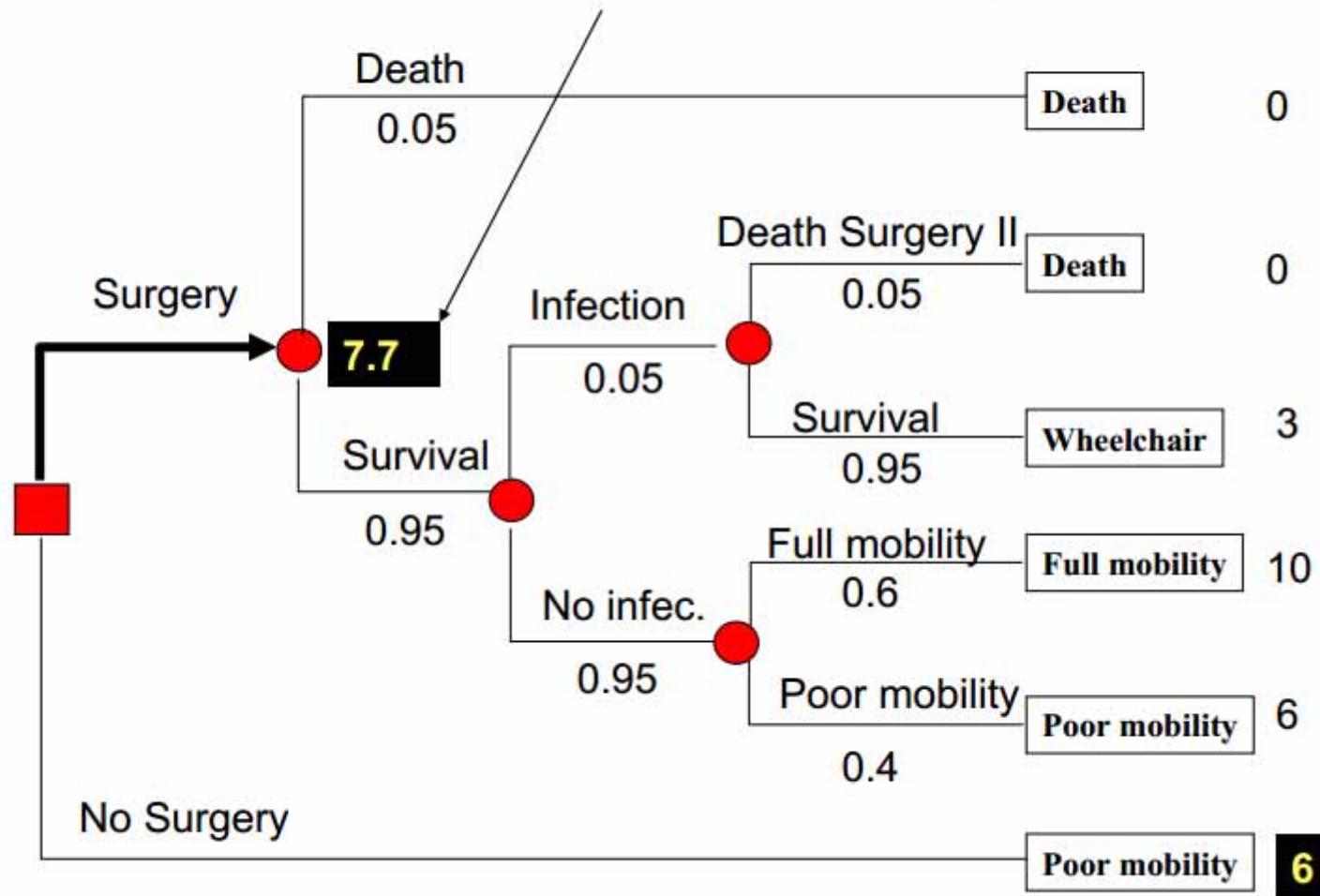
HST.951J: Medical Decision Support, Fall 2005

Instructors: Professor Lucila Ohno-Machado and Professor Staal Vinterbo

Knee Surgery



Expected Value of Surgery



For a single decision variable an agent can select
 $D = d$ for any $d \in \text{dom}(D)$.

The expected utility of decision $D = d$ is



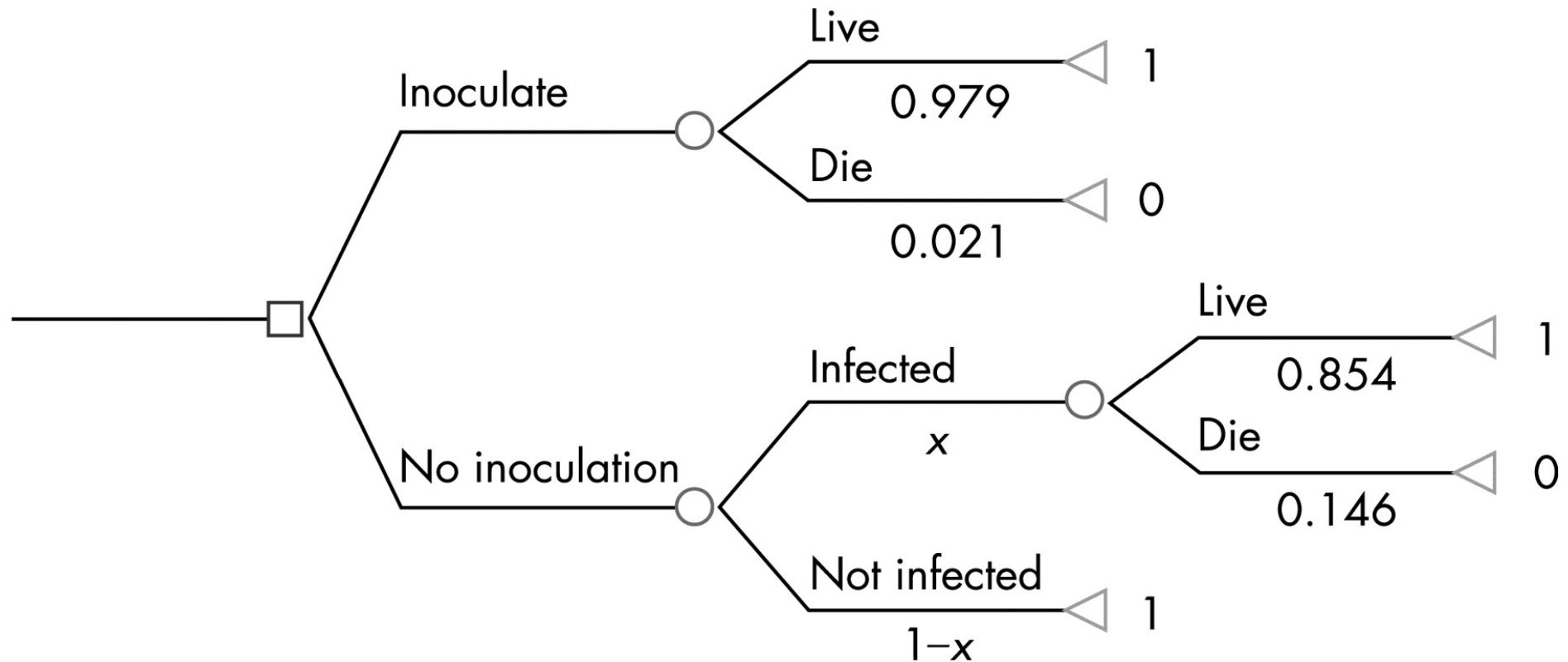
<http://www.eoht.info/page/Oskar+Morgenstern>

$$E(U | d) = \sum_{x_1, \dots, x_n} P(x_1, \dots, x_n | d) U(x_1, \dots, x_n, d)$$

An optimal single decision is the decision $D = d_{\max}$
whose expected utility is maximal:

$$d_{\max} = \arg \max_{d \in \text{dom}(D)} E(U | d)$$

Von Neumann, J. & Morgenstern, O. 1947. Theory of games and economic behavior, Princeton university press.



Ferrando, A., Pagano, E., Scaglione, L., Petrinco, M., Gregori, D. & Ciccone, G. (2009) A decision-tree model to estimate the impact on cost-effectiveness of a venous thromboembolism prophylaxis guideline. *Quality and Safety in Health Care*, 18, 4, 309-313.

Decision Model

Quantitative (statistical)

supervised

Bayesian

unsupervised

Fuzzy sets

Neural
network

Logistic

Qualitative (heuristic)

Decision
trees

Truth tables

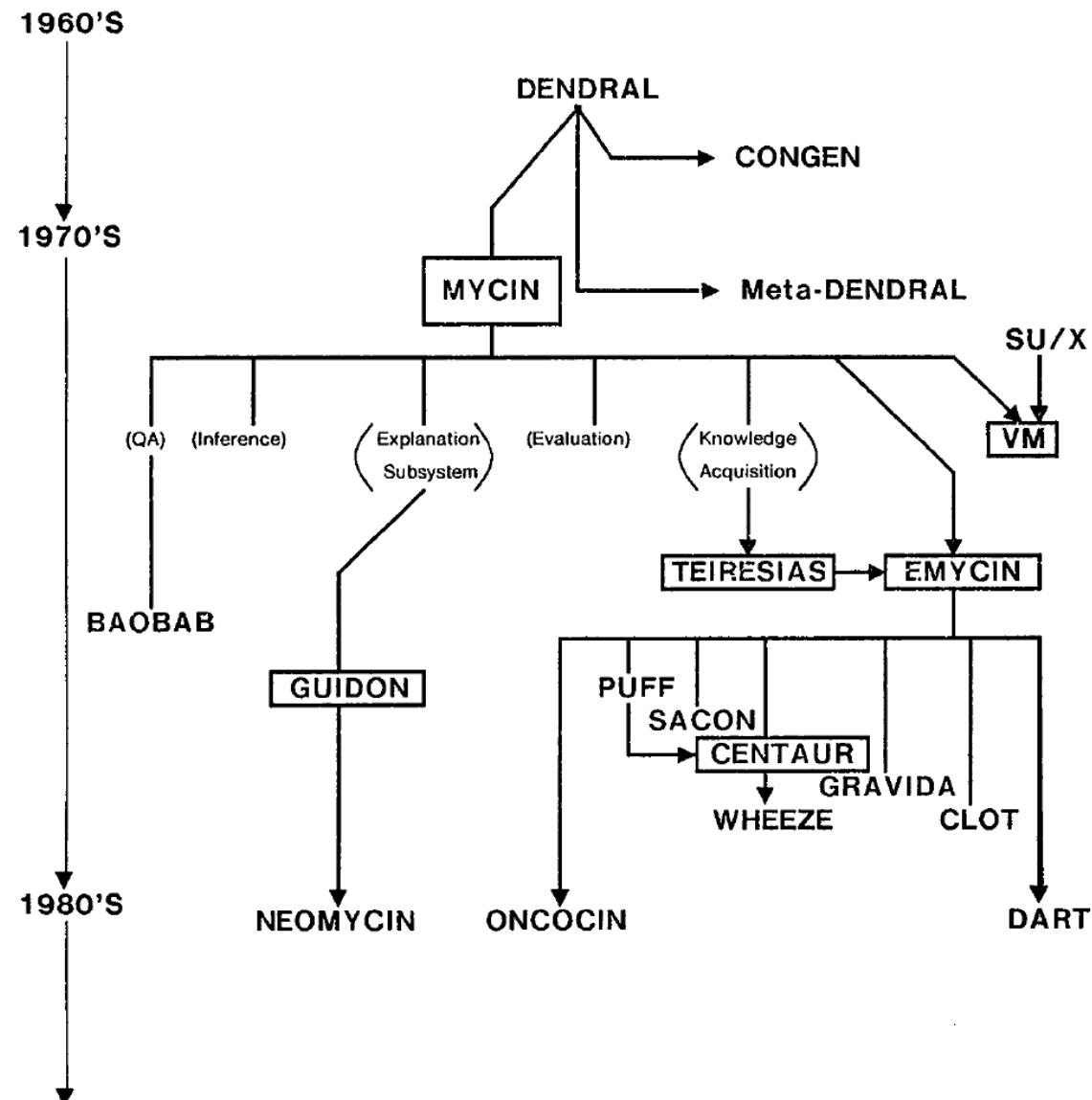
Boolean
LogicNon-
parametric
PartitioningReasoning
modelsExpert
systemsCritiquing
systems

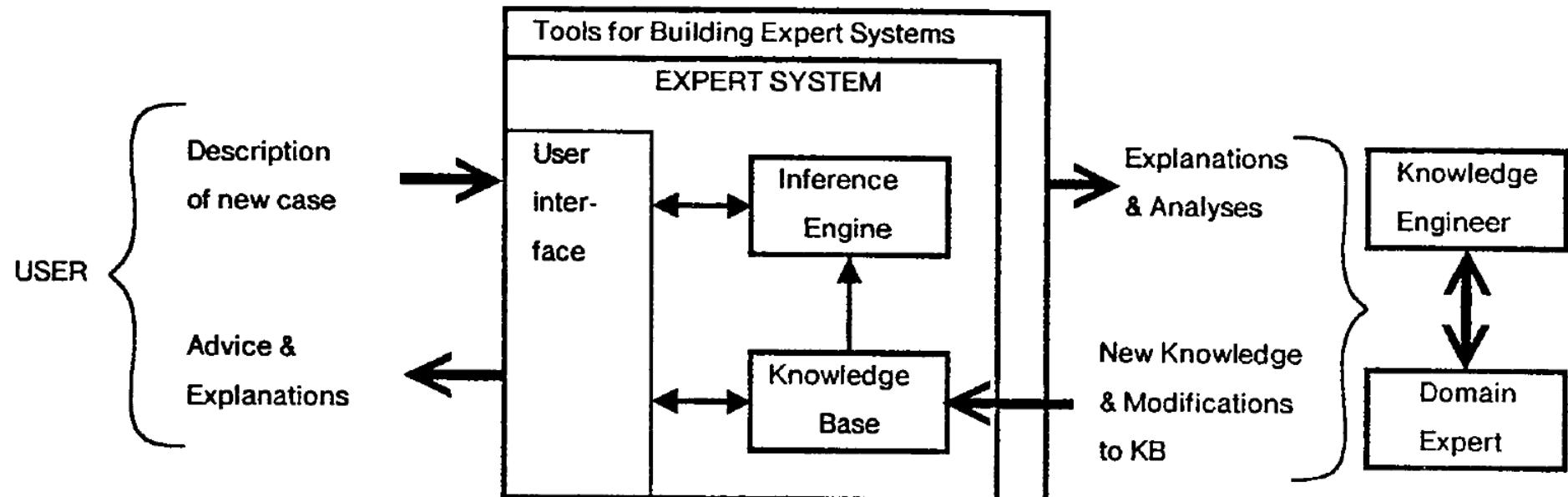
Extended by A. Holzinger after: Bemmel, J. H. v. & Musen, M. A. (1997) *Handbook of Medical Informatics*. Heidelberg, Springer.

02 History of DSS = History of AI

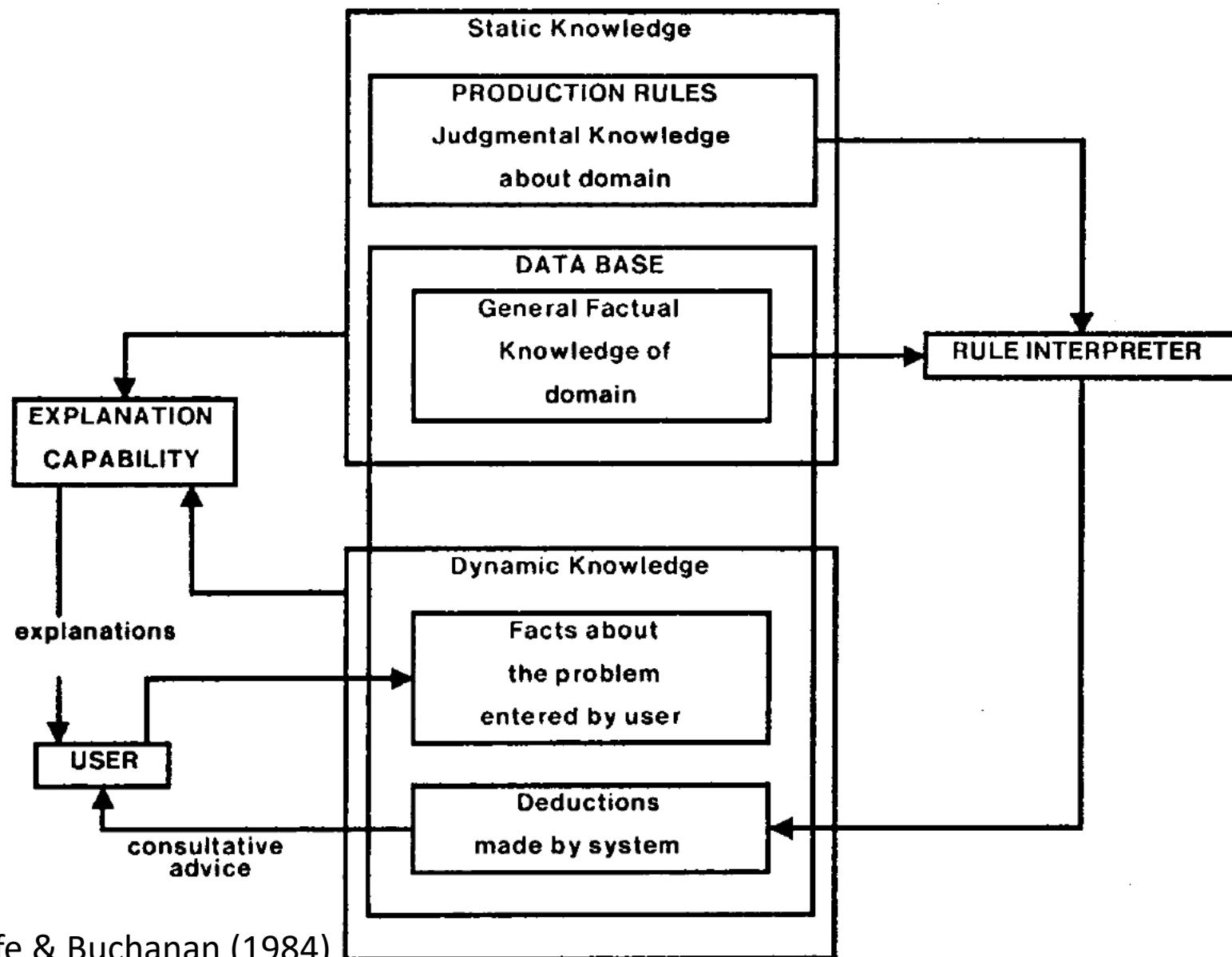
- **1943** McCulloch, W.S. & Pitts, W. A logical calculus of the ideas immanent in nervous activity. *Bulletin of Mathematical Biology*, 5, (4), 115-133, doi:10.1007/BF02459570.
- **1950** Turing, A.M. Computing machinery and intelligence. *Mind*, 59, (236), 433-460.
- **1958** John McCarthy Advice Taker: programs with common sense
- **1959** Samuel, A.L. Some studies in machine learning using the game of checkers. *IBM Journal of research and development*, 3, (3), 210-229, doi:10.1147/rd.33.0210.
- **1975** Shortliffe, E.H. & Buchanan, B.G. 1975. A model of inexact reasoning in medicine. *Mathematical biosciences*, 23, (3-4), 351-379, doi:10.1016/0025-5564(75)90047-4.
- **1978** Bellman, R. Can Computers Think? Automation of Thinking, problem solving, decision-making ...

Shortliffe, E. H. &
Buchanan, B. G. (1984)
*Rule-based expert
systems: the MYCIN
experiments of the
Stanford Heuristic
Programming Project.*
Addison-Wesley.





Shortliffe, T. & Davis, R. (1975) Some considerations for the implementation of knowledge-based expert systems *ACM SIGART Bulletin*, 55, 9-12.



- The information available to humans is often imperfect – imprecise - uncertain.
- This is especially in the medical domain the case.
- An **human agent** can cope with deficiencies.
- Classical logic permits only **exact reasoning**:
- IF A is true THEN A is non-false and
IF B is false THEN B is non-true
- Most real-world problems do not provide this exact information, mostly it is **inexact, incomplete, uncertain and/or un-measurable!**

- MYCIN is a rule-based Expert System, which is used for therapy planning for patients with bacterial infections
- Goal oriented strategy (“Rückwärtsverkettung”)
- To every rule and every entry a certainty factor (CF) is assigned, which is between 0 und 1
- Two measures are derived:
- MB: measure of belief
- MD: measure of disbelief
- Certainty factor – CF of an element is calculated by:
$$CF[h] = MB[h] - MD[h]$$
- CF is positive, if more evidence is given for a hypothesis, otherwise CF is negative
- $CF[h] = +1 \rightarrow h$ is 100 % true
- $CF[h] = -1 \rightarrow h$ is 100% false

h_1 = The identity of ORGANISM-1 is streptococcus

h_2 = PATIENT-1 is febrile

h_3 = The name of PATIENT-1 is John Jones

$CF[h_1, E] = .8$: There is strongly suggestive evidence (.8) that the identity of ORGANISM-1 is streptococcus

$CF[h_2, E] = -.3$: There is weakly suggestive evidence (.3) that PATIENT-1 is not febrile

$CF[h_3, E] = +1$: It is definite (1) that the name of PATIENT-1 is John Jones

Shortliffe, E. H. & Buchanan, B. G. (1984) *Rule-based expert systems: the MYCIN experiments of the Stanford Heuristic Programming Project*. Addison-Wesley.



The image consists of three separate panels. The left panel is a magazine cover for 'SIBYLLE' from May 1968, showing a woman wearing a beret and a plaid scarf. The middle panel is a page from a magazine titled 'Die Geheimnisse des Rechenautomaten' (The Secrets of the Computer), featuring a small illustration of a woman's head. The right panel is a small illustration of a man and a woman holding hands.

Image credit to Bernhard Schölkopf

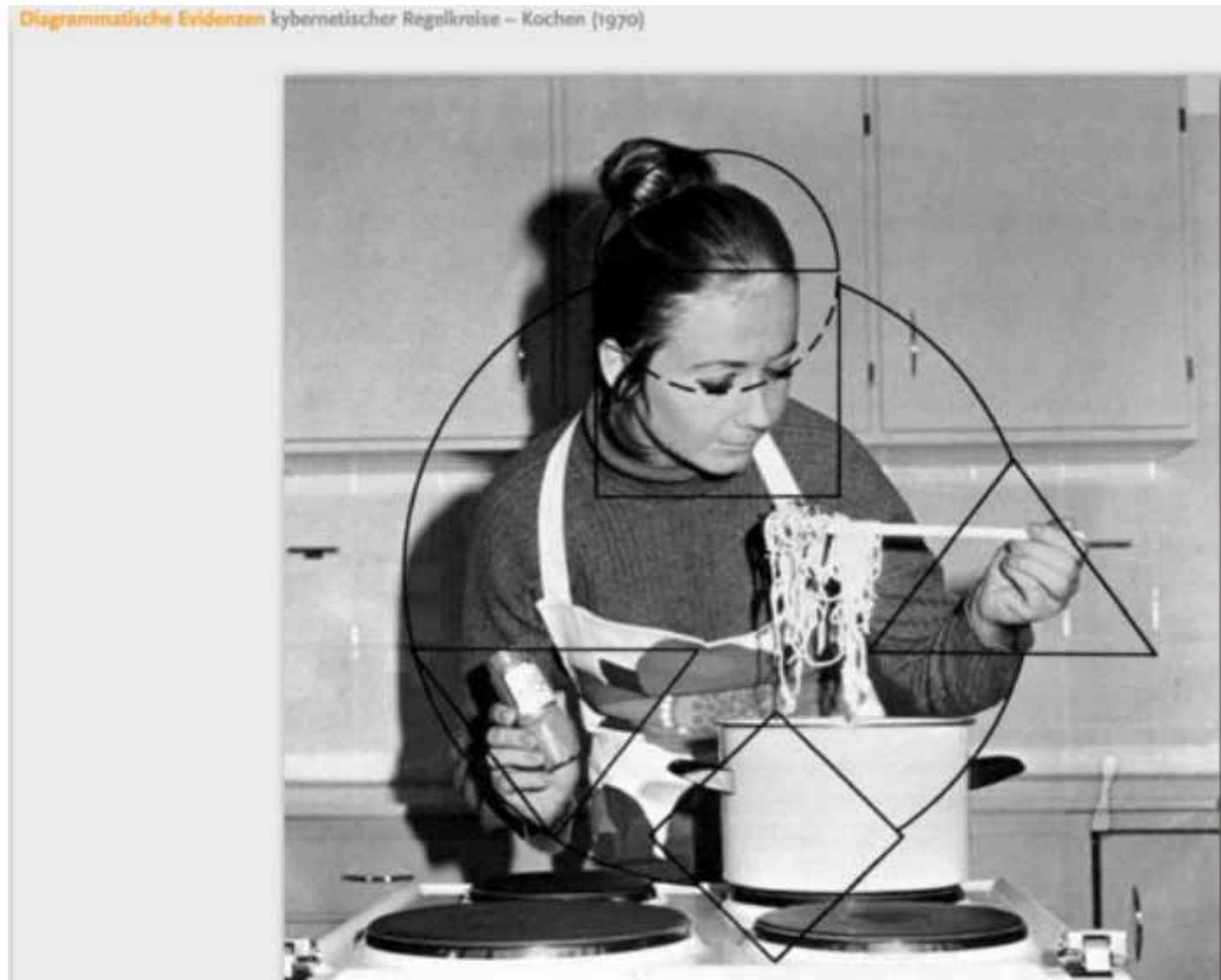
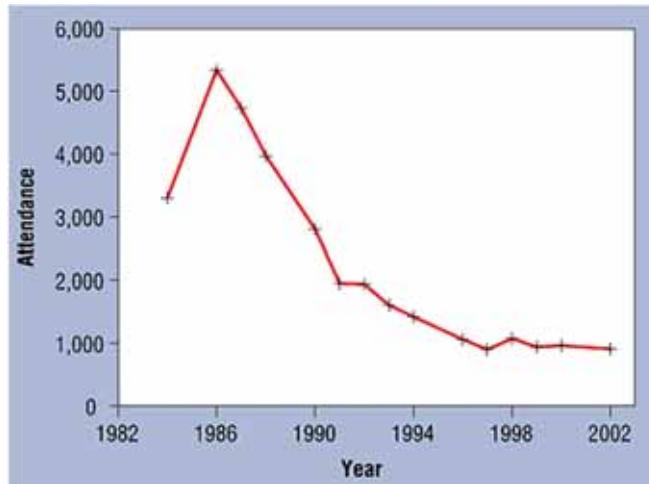


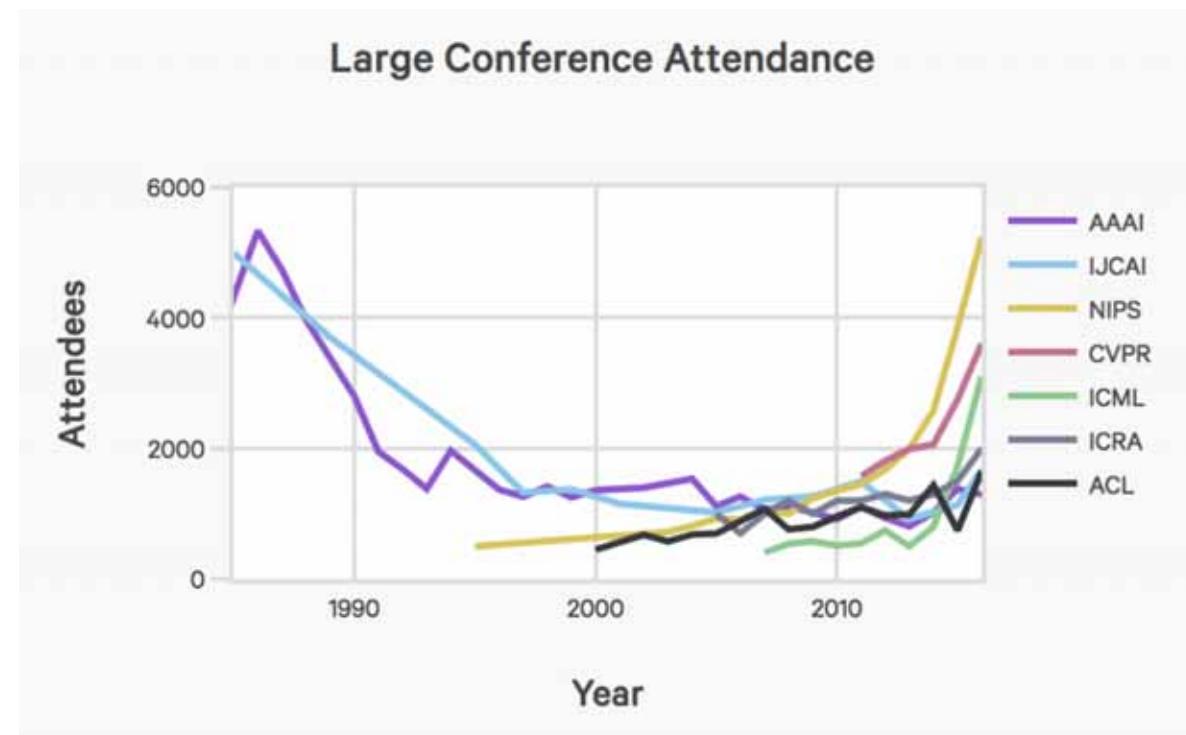
Image credit to Bernhard Schölkopf



<https://blogs.dxc.technology/2017/04/25/are-we-heading-toward-an-ai-winter/>

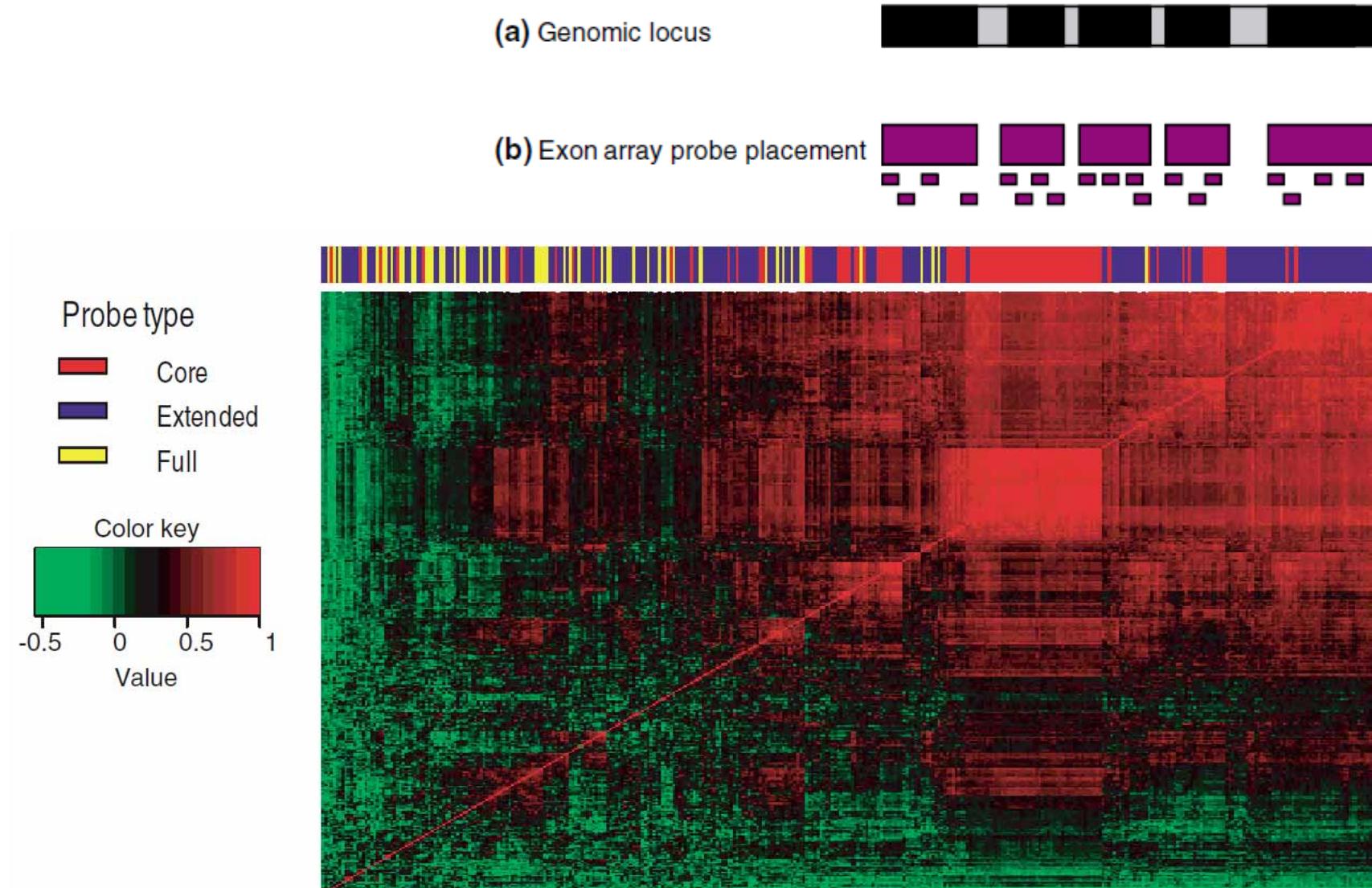


<https://www.computer.org/csl/mags/ex/2003/03/x3018.html>

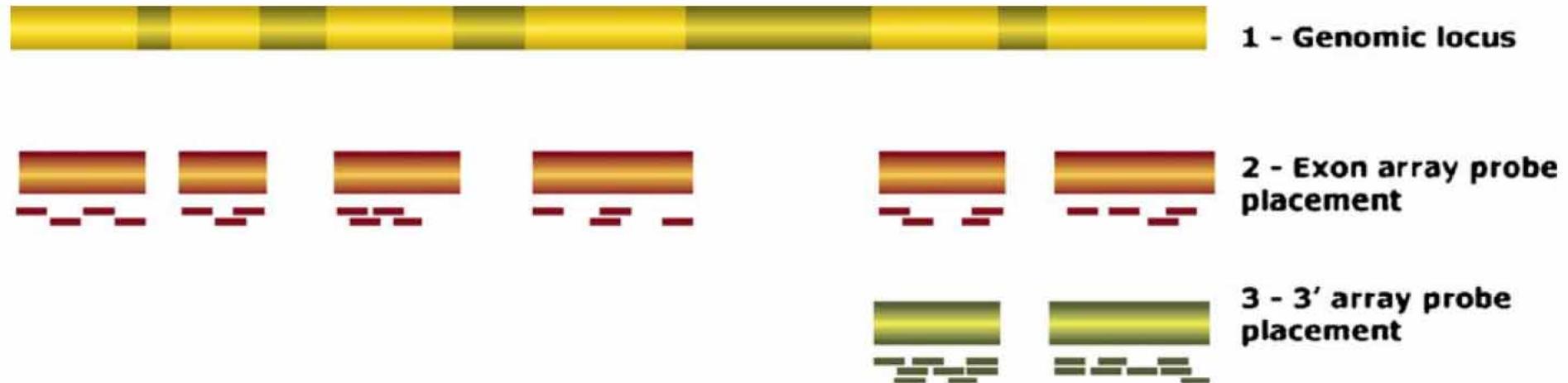


<https://medium.com/machine-learning-in-practice/nips-accepted-papers-stats-26f124843aa0>

03 Example: P4-Medicine

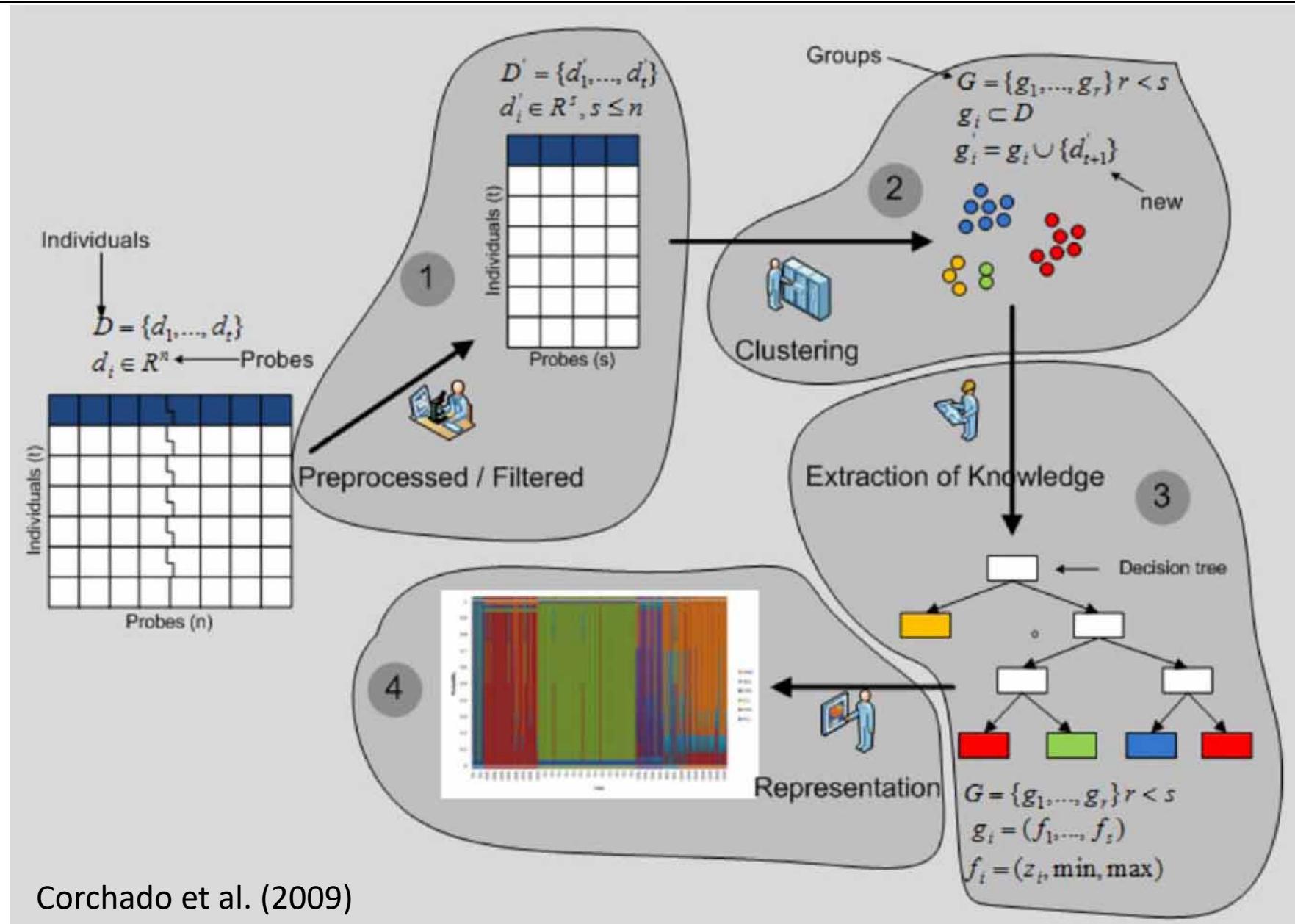


Kapur, K., Xing, Y., Ouyang, Z. & Wong, W. (2007) Exon arrays provide accurate assessments of gene expression. *Genome Biology*, 8, 5, R82.



Exon array structure. Probe design of exon arrays. (1) Exon—intron structure of a gene. Gray boxes represent introns, rest represent exons. Introns are not drawn to scale. (2) Probe design of exon arrays. Four probes target each putative exon. (3) Probe design of 30expression arrays. Probe target the 30end of mRNA sequence.

Corchado, J. M., De Paz, J. F., Rodriguez, S. & Bajo, J. (2009) Model of experts for decision support in the diagnosis of leukemia patients. *Artificial Intelligence in Medicine*, 46, 3, 179-200.

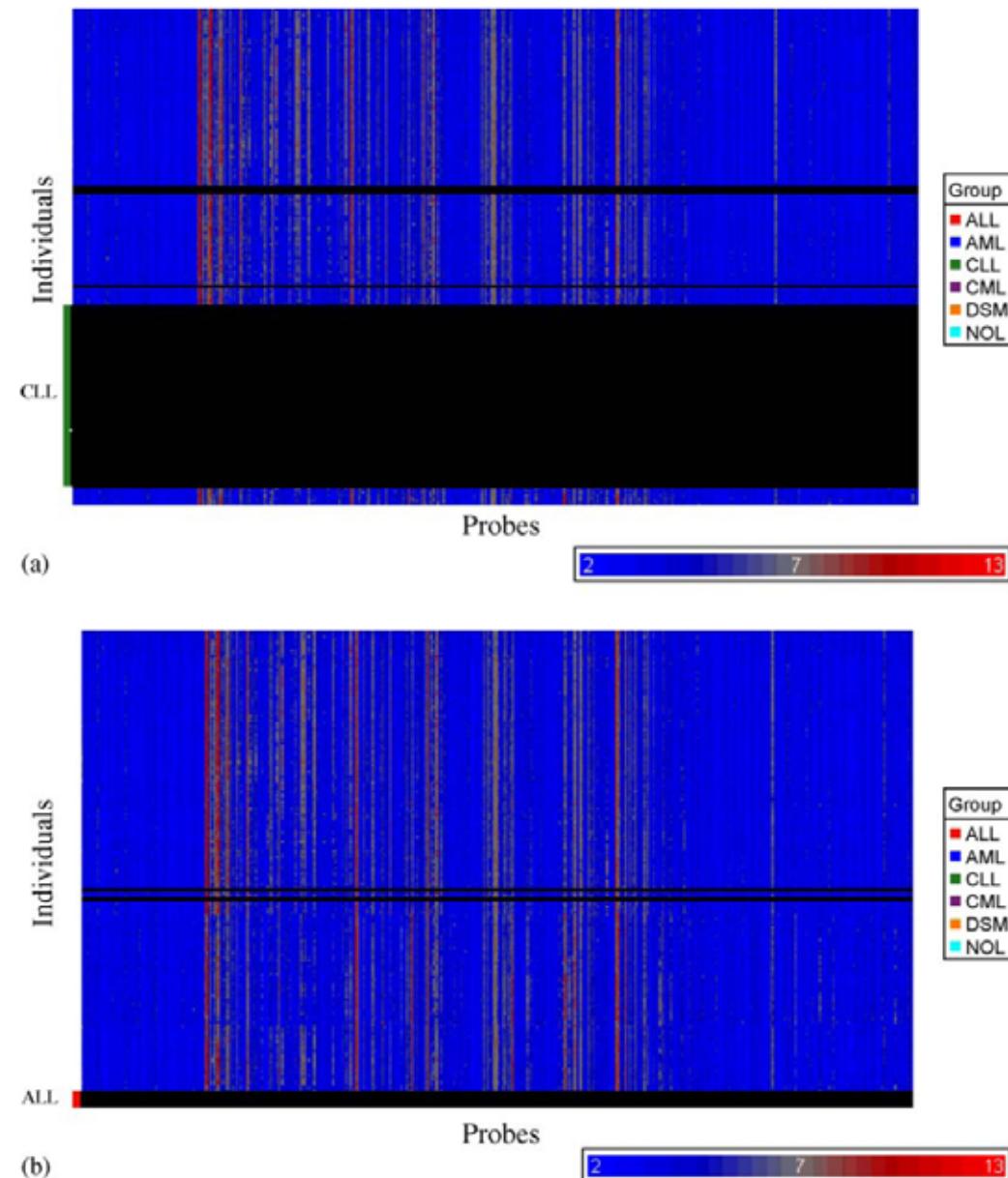


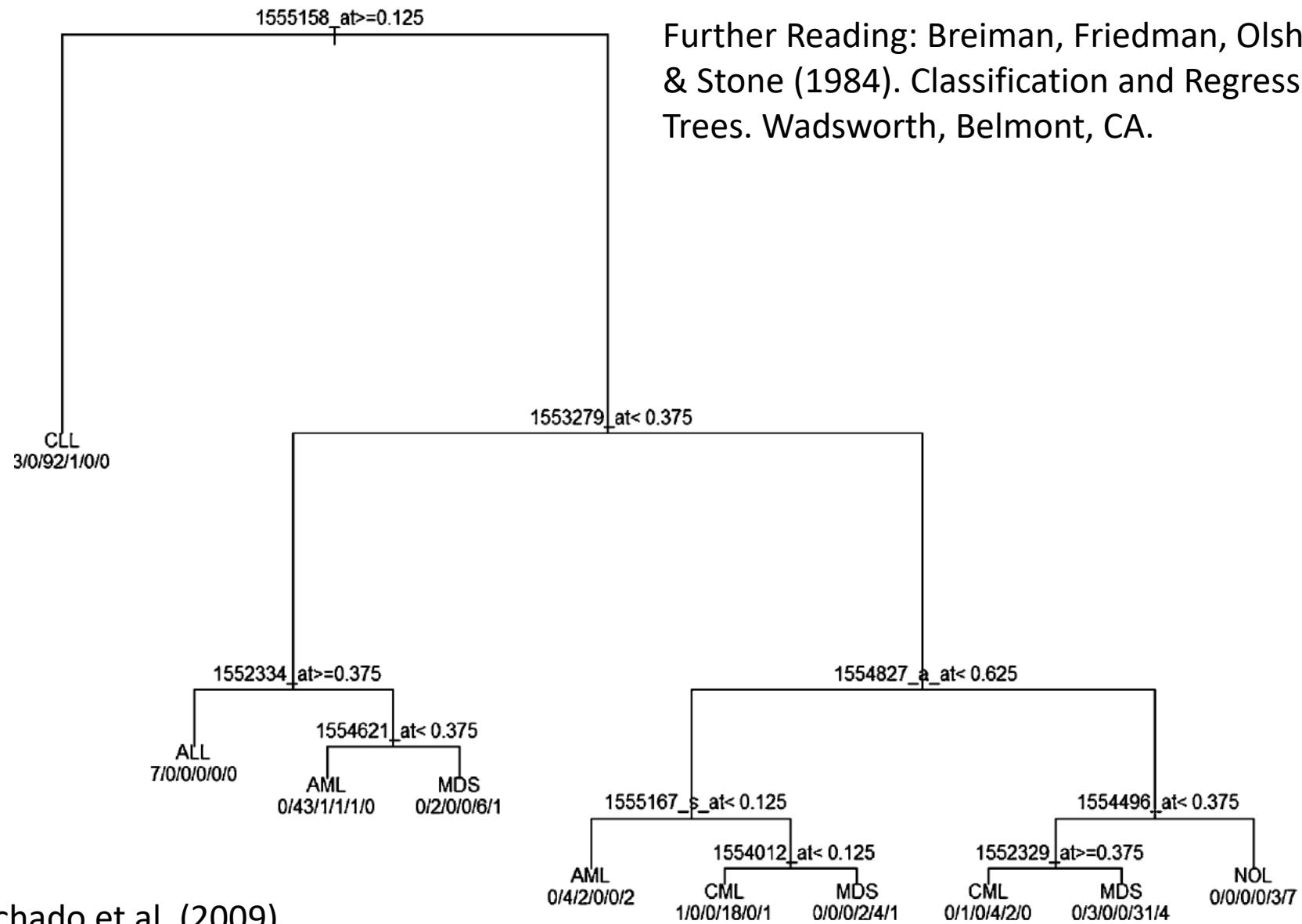
Corchado et al. (2009)

A = acute, C = chronic,
L = lymphocytic, M = myeloid

- **ALL** = cancer of the blood AND bone marrow caused by an abnormal proliferation of lymphocytes.
- **AML** = cancer in the bone marrow characterized by the proliferation of myeloblasts, red blood cells or abnormal platelets.
- **CLL** = cancer characterized by a proliferation of lymphocytes in the bone marrow.
- **CML** = caused by a proliferation of white blood cells in the bone marrow.
- **MDS (Myelodysplastic Syndromes)** = a group of diseases of the blood and bone marrow in which the bone marrow does not produce a sufficient amount of healthy cells.
- **NOL (Normal) = No leukemias**

Corchado et al. (2009)

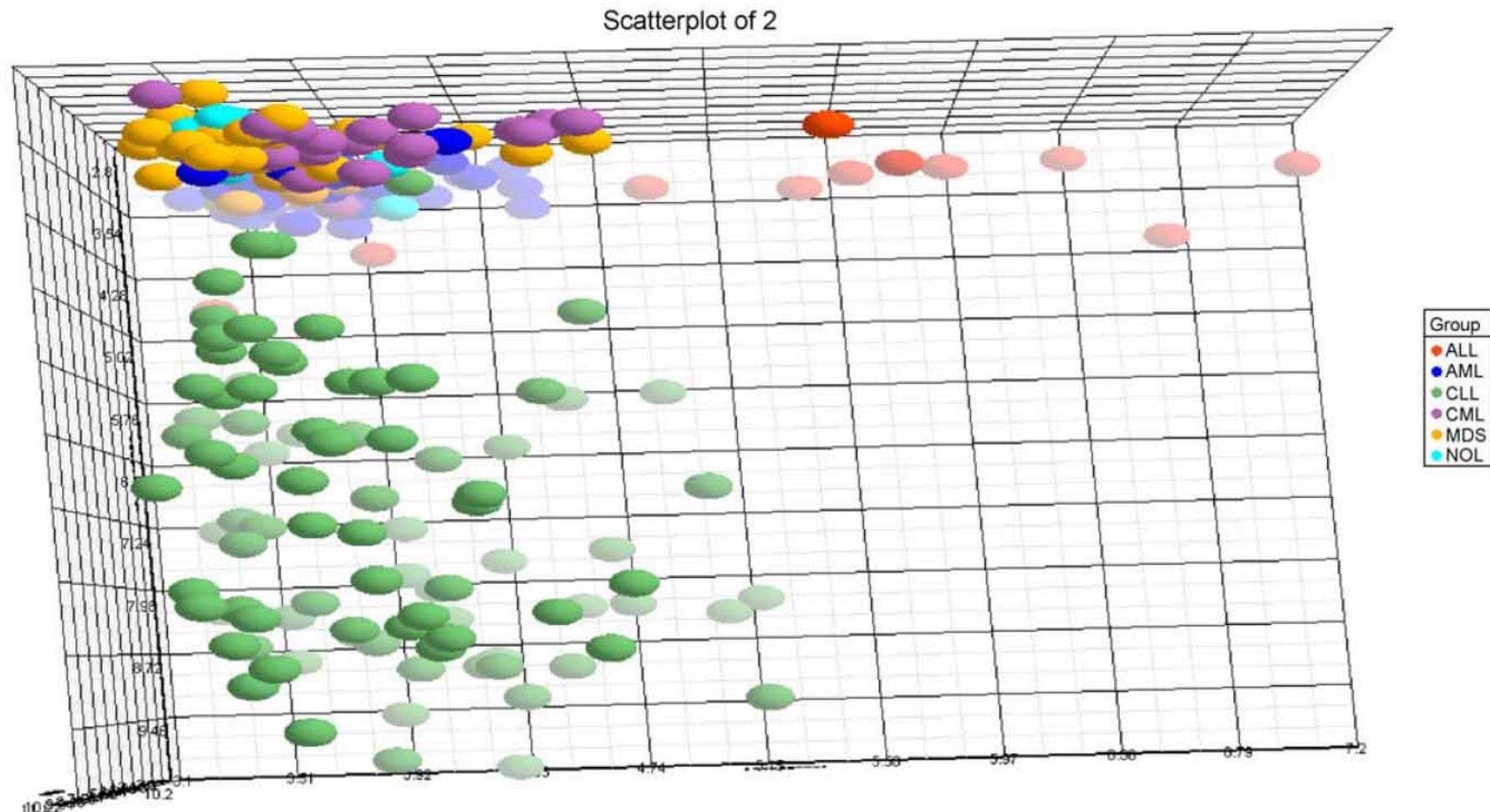




Corchado et al. (2009)

Further Reading: Breiman, Friedman, Olshen, & Stone (1984). Classification and Regression Trees. Wadsworth, Belmont, CA.

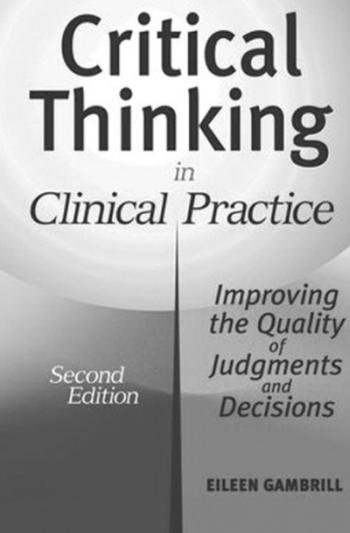
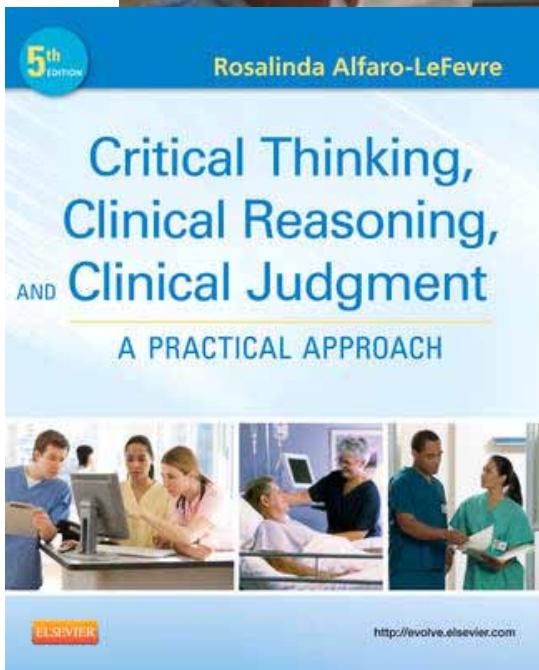
Classification CLL—ALL. Representation of the probes of the decision tree which classify the CLL and ALL to 1555158_at, 1553279_at and 1552334_at

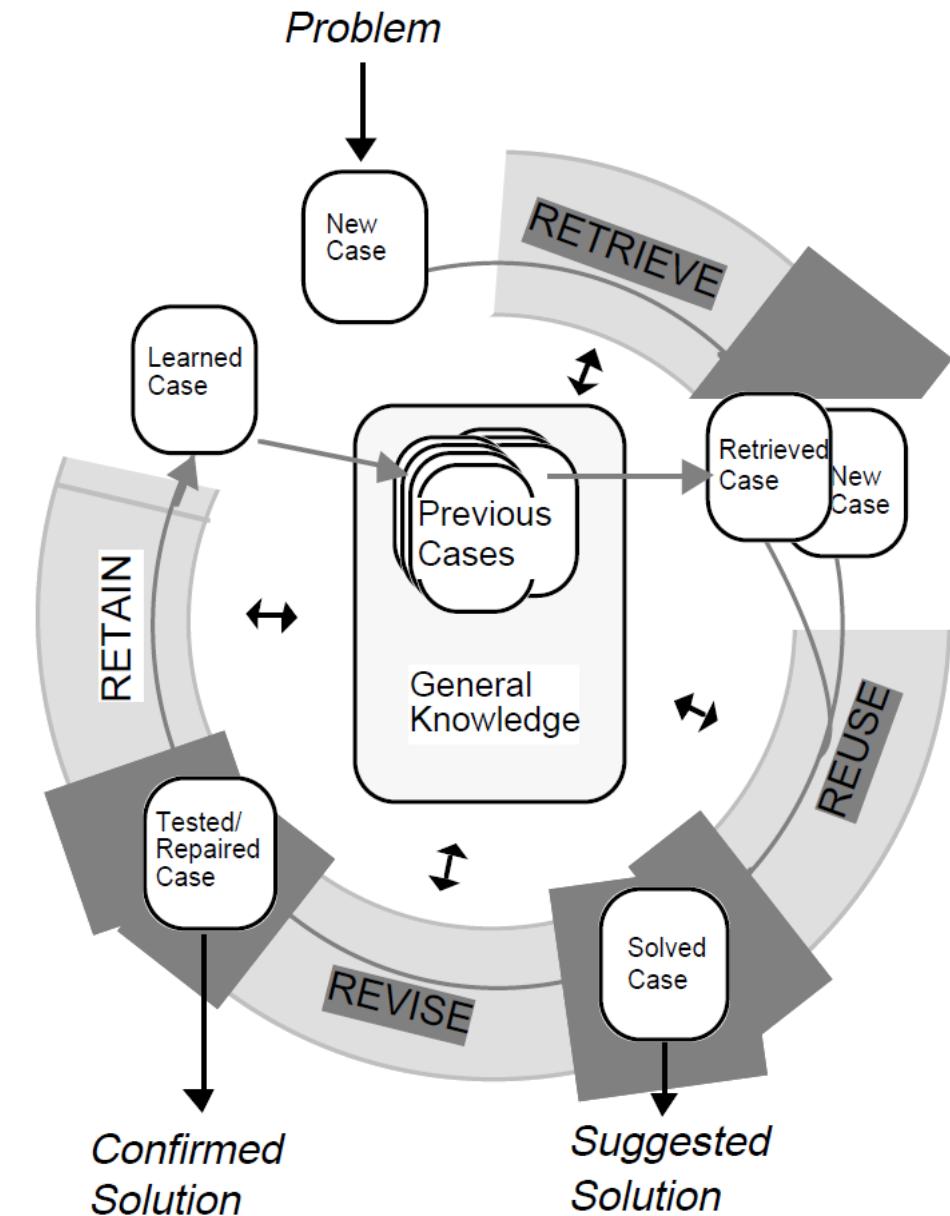


Corchado et al. (2009)

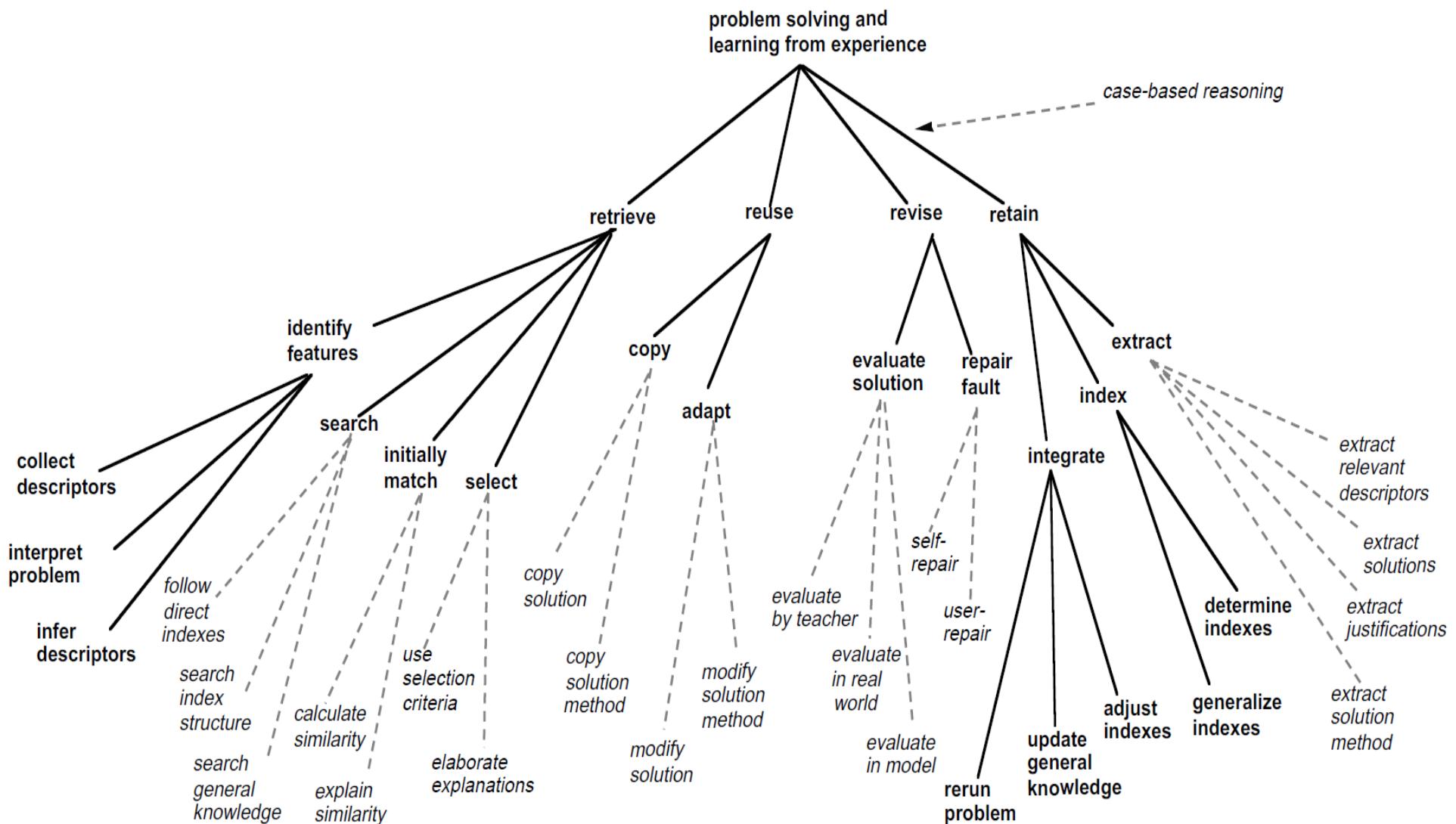
- The model of Corchado et al. (2009) combines:
- 1) methods to **reduce the dimensionality** of the original data set;
- 2) pre-processing and data filtering techniques;
- 3) a clustering method to classify patients; and
- 4) extraction of knowledge techniques
- The system reflects how human experts work in a lab, but
 - 1) **reduces the time** for making predictions;
 - 2) **reduces the rate of human error**; and
 - 3) **works with high-dimensional data** from exon arrays

04 Example: Case Based Reasoning (CBR)





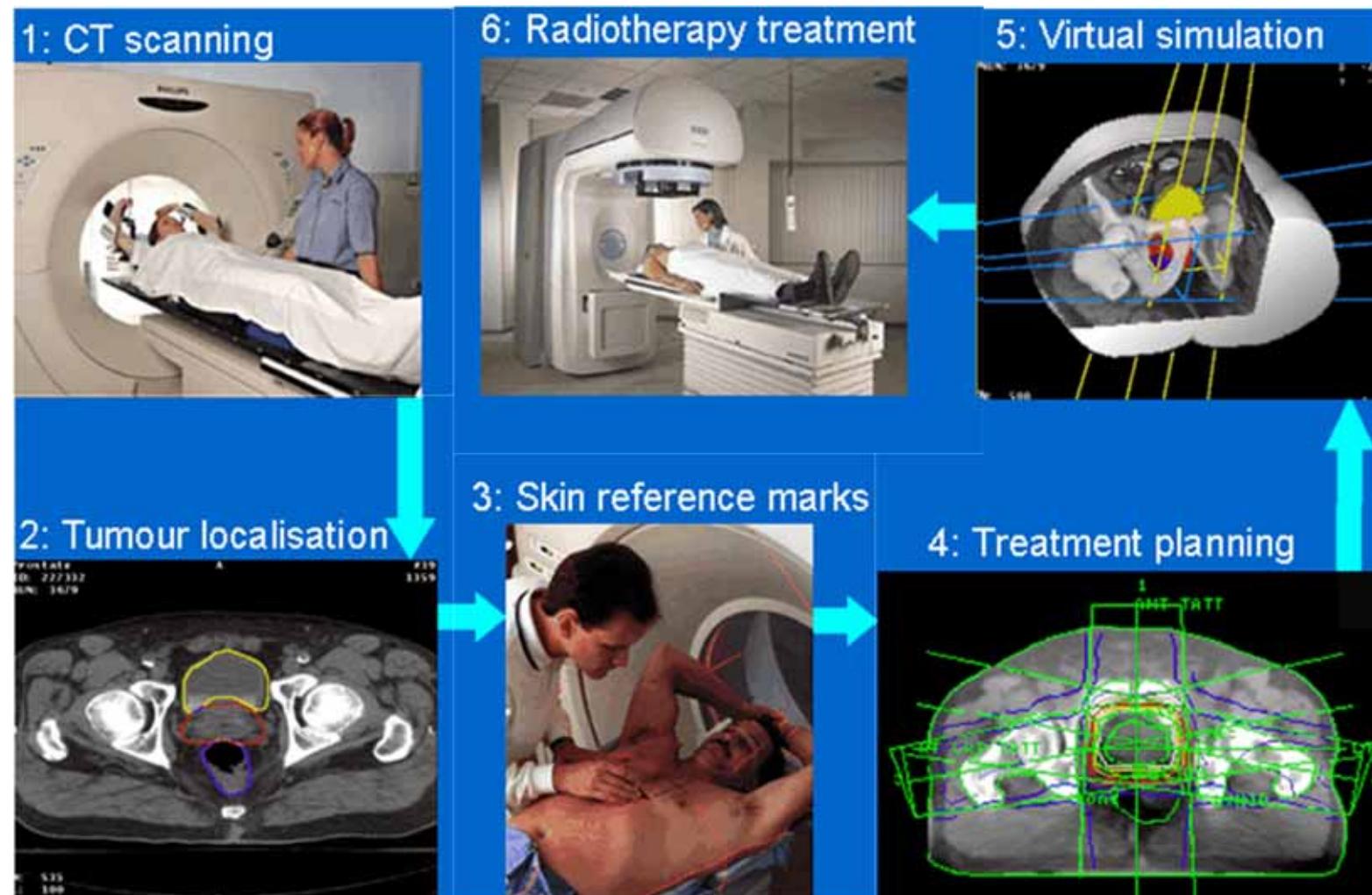
Aamodt, A. & Plaza, E. (1994) Case-based reasoning: Foundational issues, methodological variations, and system approaches. *AI Communications*, 7, 1, 39-59.



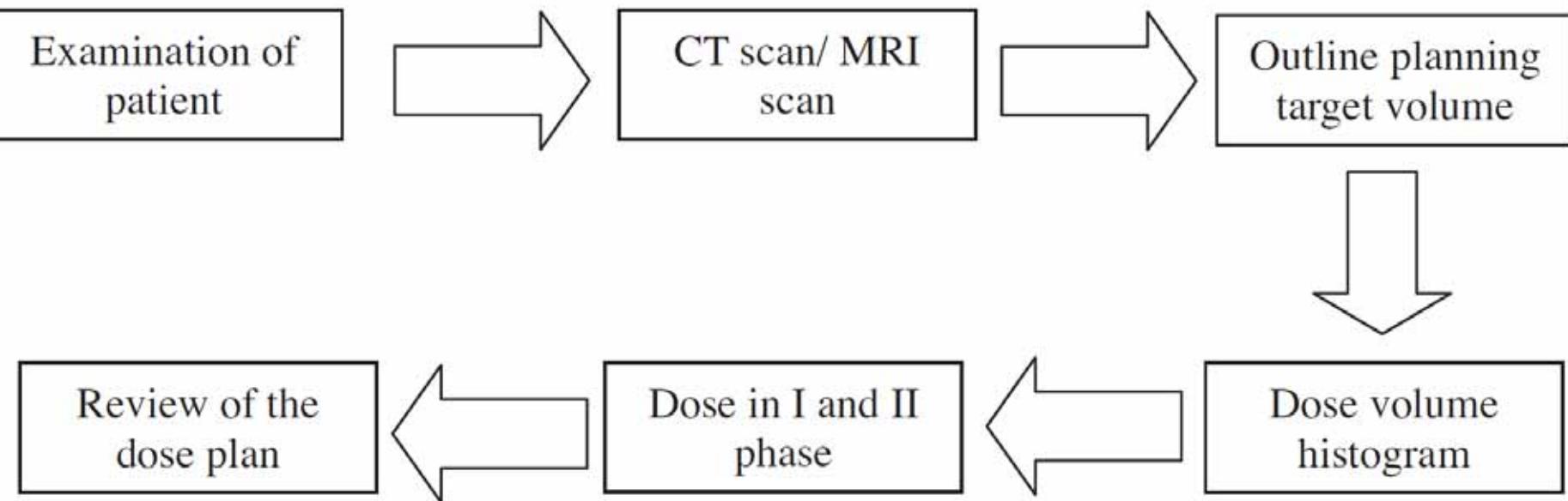
Aamodt & Plaza (1994)



Source: <http://www.teachingmedicalphysics.org.uk>



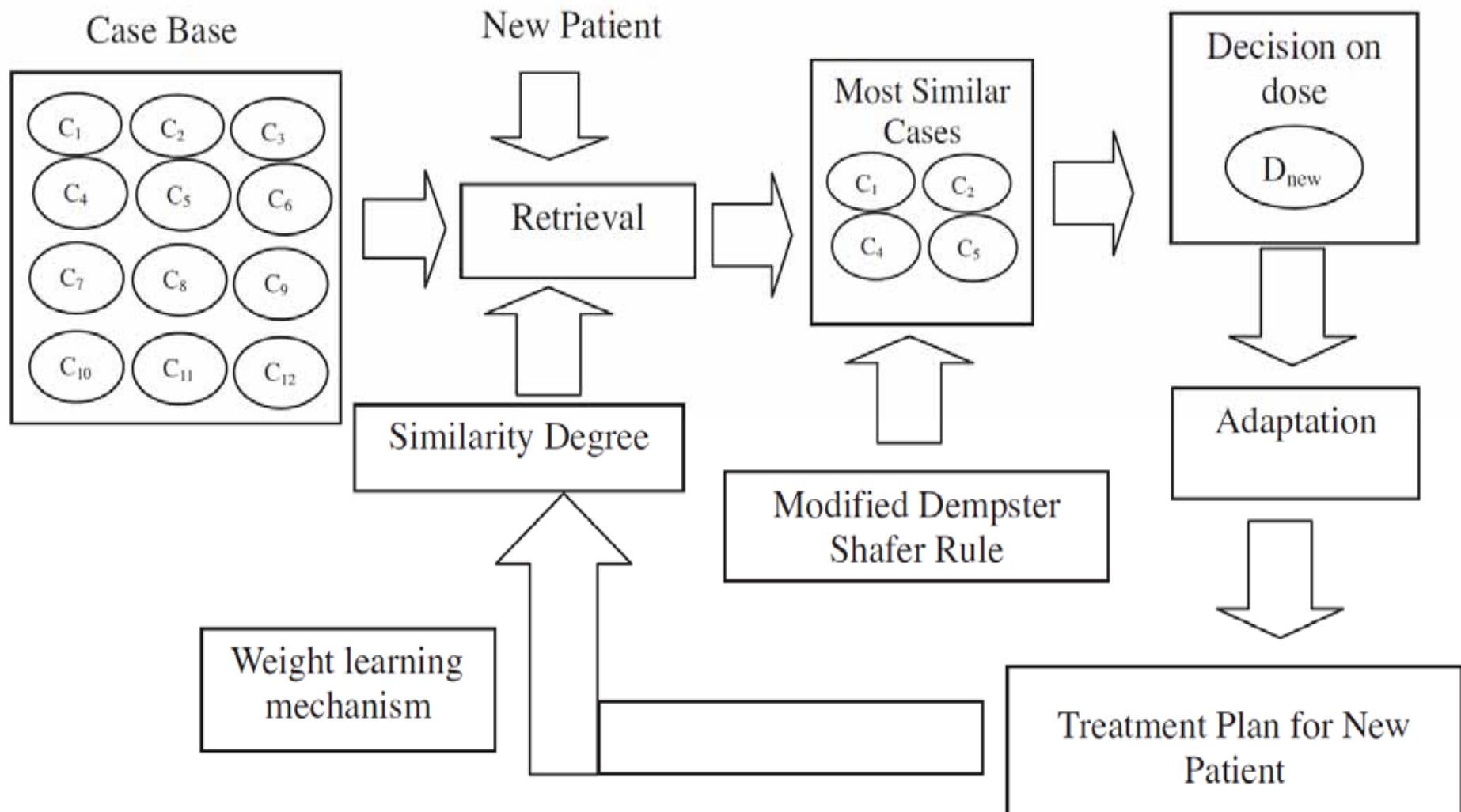
Source: Imaging Performance Assessment of CT Scanners Group, <http://www.impactscan.org>



Measures:

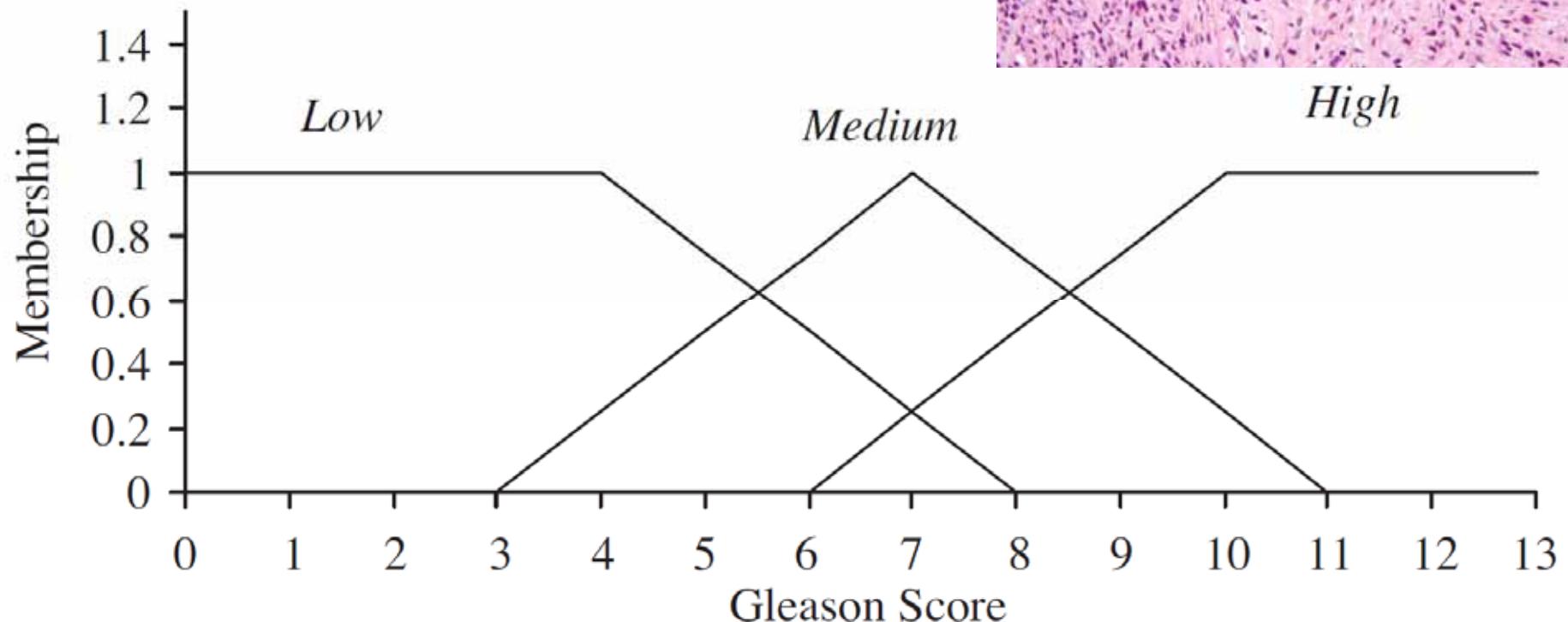
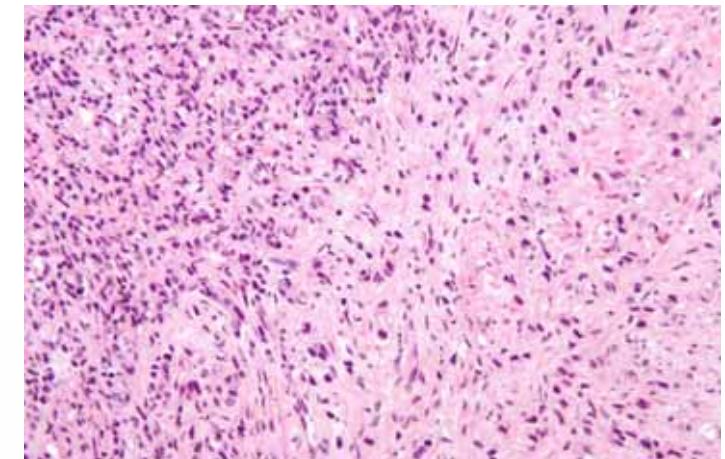
- 1) Clinical Stage = a labelling system
- 2) Gleason Score = grade of prostate cancer = integer between 1 to 10; and
- 3) Prostate Specific Antigen (PSA) value between 1 to 40
- 4) Dose Volume Histogram (DVH) = pot. risk to the rectum (66, 50, 25, 10 %)

Petrovic, S., Mishra, N. & Sundar, S. (2011) A novel case based reasoning approach to radiotherapy planning. *Expert Systems With Applications*, 38, 9, 10759-10769.



Petrovic, S., Mishra, N. & Sundar, S. (2011) A novel case based reasoning approach to radiotherapy planning. *Expert Systems With Applications*, 38, 9, 10759-10769.

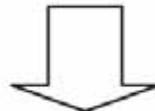
Gleason score evaluates the grade of prostate cancer. Values: integer within the range



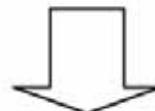
Petrovic, S., Mishra, N. & Sundar, S. (2011) A novel case based reasoning approach to radiotherapy planning. *Expert Systems With Applications*, 38, 9, 10759-10769.

Petrovic et al. (2011)

Dose plan suggested by Dempster-Shafer rule (62Gy+10Gy}



Dose received by 10% of rectum is 56.02 Gy (maximum dose limit =55 Gy)



Proposed dose plan

Yes

Feasible dose plan

No

Modification



Modification of dose plan:

New dose plan: 62Gy +8 Gy

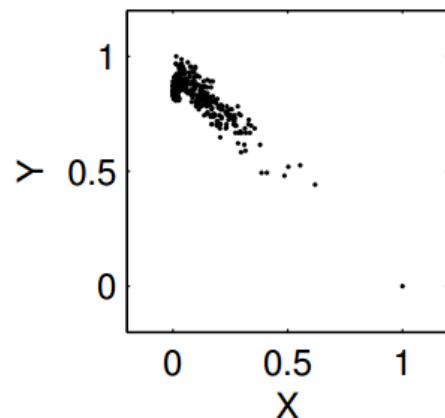
Dose received by 10% of rectum is: 54.26 Gy (feasible dose plan)

05 Causal Reasoning

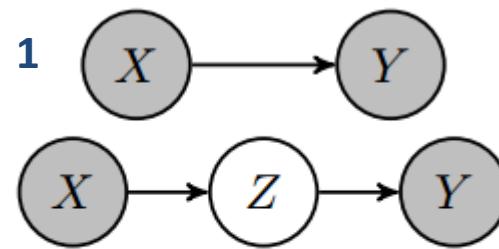
- “How do humans generalize from few examples?”
 - Learning relevant representations
 - Disentangling the explanatory factors
 - Finding the shared underlying explanatory factors, in particular between $P(x)$ and $P(Y|X)$, with a causal link between $Y \rightarrow X$

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.



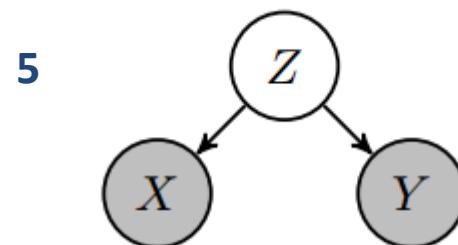
Joris M. Mooij, Jonas Peters, Dominik Janzing, Jakob Zscheischler & Bernhard Schölkopf
 2016. Distinguishing cause from effect using observational data: methods and benchmarks. The Journal of Machine Learning Research, 17, (1), 1103-1204.



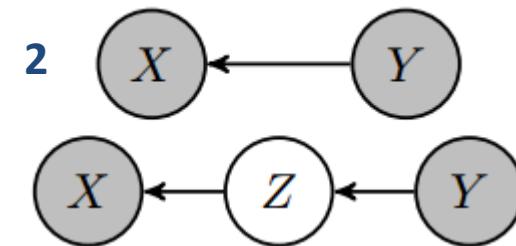
$$\begin{aligned}\mathbb{P}_Y &\neq \mathbb{P}_{Y \mid \text{do}(x)} = \mathbb{P}_{Y \mid x} \\ \mathbb{P}_X &= \mathbb{P}_{X \mid \text{do}(y)} \neq \mathbb{P}_{X \mid y}\end{aligned}$$



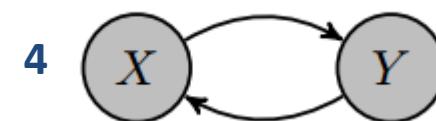
$$\begin{aligned}\mathbb{P}_Y &= \mathbb{P}_{Y \mid \text{do}(x)} = \mathbb{P}_{Y \mid x} \\ \mathbb{P}_X &= \mathbb{P}_{X \mid \text{do}(y)} = \mathbb{P}_{X \mid y}\end{aligned}$$



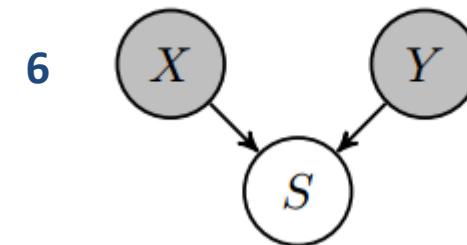
$$\begin{aligned}\mathbb{P}_Y &= \mathbb{P}_{Y \mid \text{do}(x)} \neq \mathbb{P}_{Y \mid x} \\ \mathbb{P}_X &= \mathbb{P}_{X \mid \text{do}(y)} \neq \mathbb{P}_{X \mid y}\end{aligned}$$



$$\begin{aligned}\mathbb{P}_Y &= \mathbb{P}_{Y \mid \text{do}(x)} \neq \mathbb{P}_{Y \mid x} \\ \mathbb{P}_X &\neq \mathbb{P}_{X \mid \text{do}(y)} = \mathbb{P}_{X \mid y}\end{aligned}$$



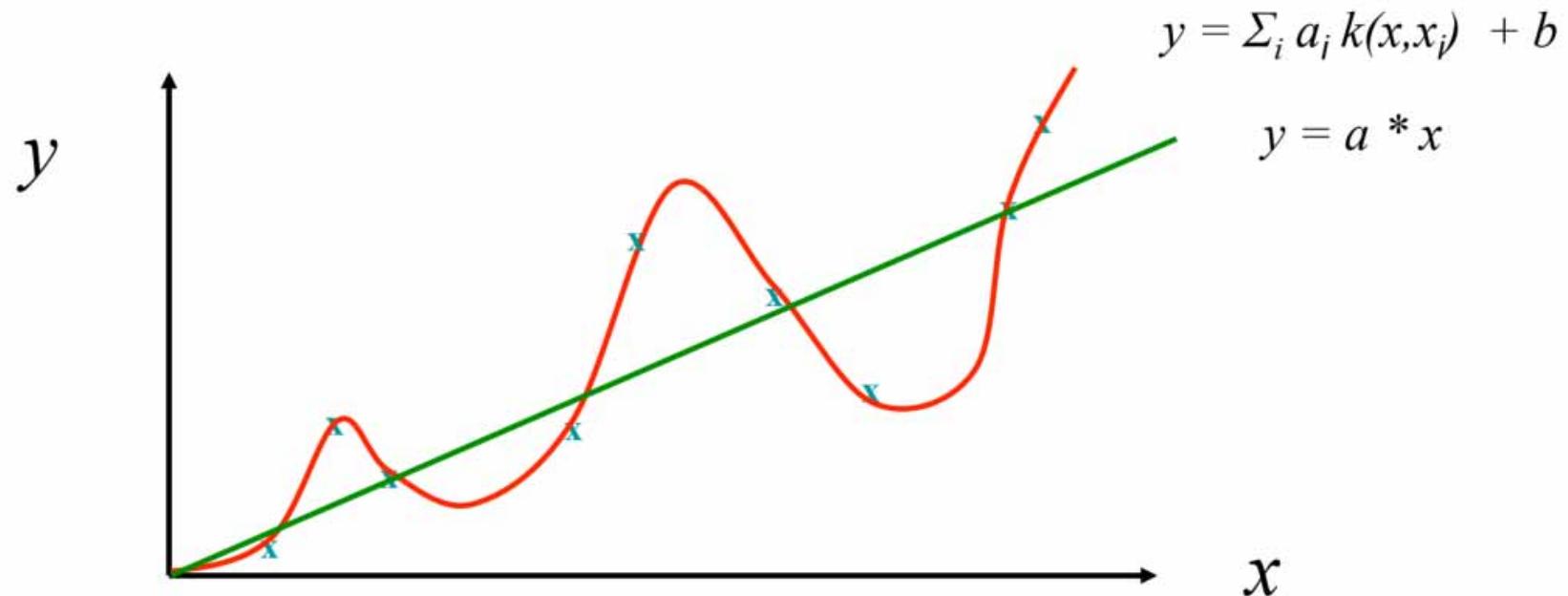
$$\begin{aligned}\mathbb{P}_Y &\neq \mathbb{P}_{Y \mid \text{do}(x)} \neq \mathbb{P}_{Y \mid x} \\ \mathbb{P}_X &\neq \mathbb{P}_{X \mid \text{do}(y)} \neq \mathbb{P}_{X \mid y}\end{aligned}$$



$$\begin{aligned}\mathbb{P}_{Y \mid s} &\neq \mathbb{P}_{Y \mid \text{do}(x), s} = \mathbb{P}_{Y \mid x, s} \\ \mathbb{P}_{X \mid s} &\neq \mathbb{P}_{X \mid \text{do}(y), s} = \mathbb{P}_{X \mid y, s}\end{aligned}$$

- **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.



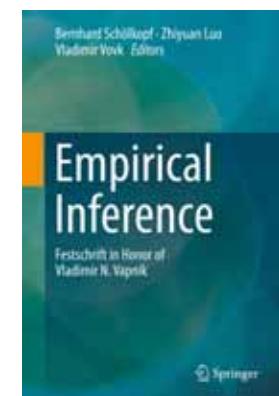
Gottfried W. Leibniz (1646-1716)

Hermann Weyl (1885-1955)

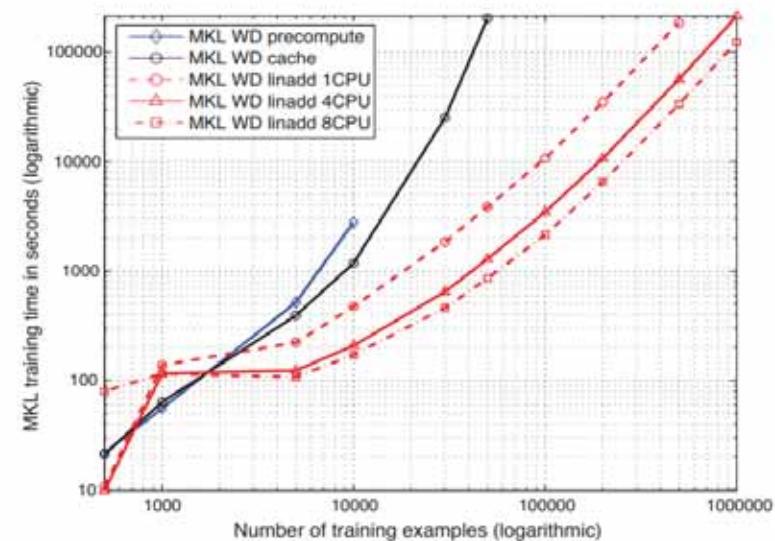
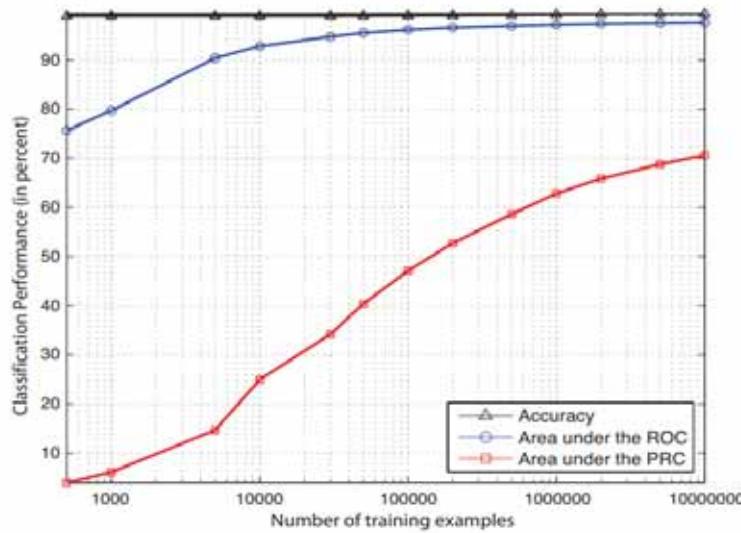
Vladimir Vapnik (1936-)

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.

06 Explainability

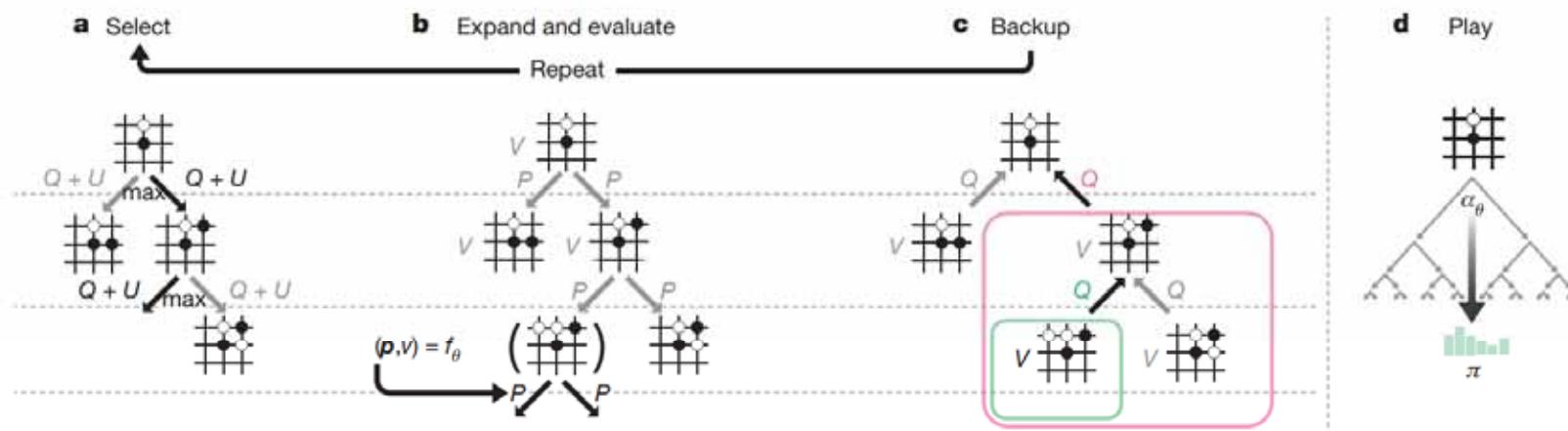


Figure 2 | MCTS in AlphaGo Zero. **a**, Each simulation traverses the tree by selecting the edge with maximum action value Q , plus an upper confidence bound U that depends on a stored prior probability P and visit count N for that edge (which is incremented once traversed). **b**, The leaf node is expanded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

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

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.



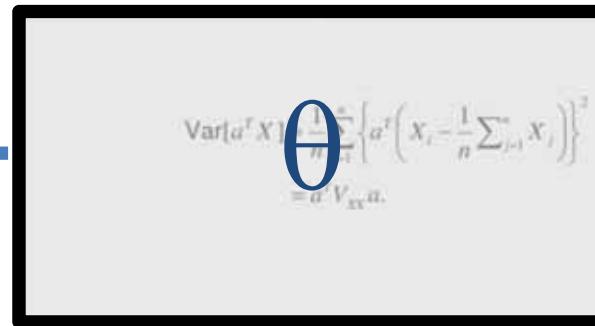
David Silver, Aja Huang, Chris J. Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, Sander Dieleman, Dominik Grewe, John Nham, Nal Kalchbrenner, Ilya Sutskever, Timothy Lillicrap, Madeleine Leach, Koray Kavukcuoglu, 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.

Causability :=
a property of a person
(Human Intelligence)

Explainability :=
a property of a system
(Artificial Intelligence)

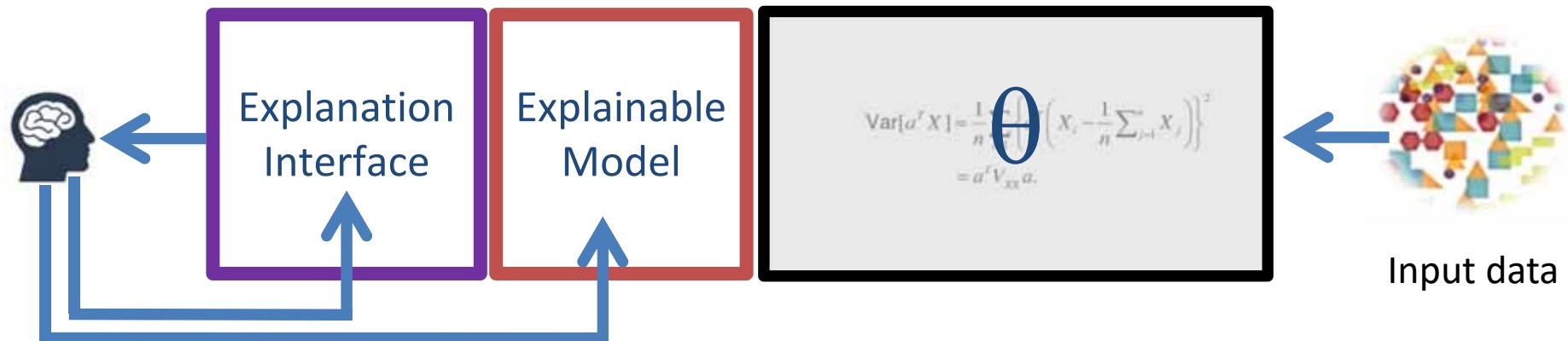
Andreas Holzinger et al. 2019. Causability and Explainability of AI in Medicine. Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery, doi:10.1002/widm.1312.

*Why did the algorithm do that?
Can I trust these results?
How can I correct an error?*



Input data

A possible solution



The domain expert can understand why ...

The domain expert can learn and correct errors ...

The domain expert can re-enact on demand ...

1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	0.953	0.894	0.620	0.699	0.629	0.546	0.540	1.000	0.526	1.000	0.522	0.483	0.471	1.000	0.522	0.576	0.658					
1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	0.722	0.638	1.000	0.785	0.743	0.792	0.801	0.875	0.712	1.000	0.444	0.947	0.431	1.000	0.793	1.000	0.635					
1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	0.658	0.633	0.569	0.561	0.589	0.640	0.659	0.845	0.932	0.512	0.575	0.941	1.000	0.991	1.000	0.892						
1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	0.932	0.639	0.575	0.544	0.501	0.489	0.470	0.454	0.576	0.581	0.707	0.992	1.000	1.000	1.000	1.000	1.000					
1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	0.711	0.644	0.569	0.541	0.461	0.430	0.425	0.381	0.364	0.437	0.562	0.509	0.528	0.678	1.000	0.991	1.000	1.000				
1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	0.680	0.594	0.579	0.513	0.490	0.429	0.405	0.425	0.381	0.401	0.387	0.367	0.484	0.428	0.483	0.659	0.936	1.000	1.000			
1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	0.761	0.677	0.610	0.565	0.511	0.498	0.457	0.416	0.396	0.388	0.369	0.355	0.359	0.468	0.392	0.380	0.487	0.4	0.500	0.744	1.000	0.485
1.000	1.000	1.000	1.000	0.861	0.640	0.579	0.560	0.542	0.476	0.470	0.441	0.405	0.389	0.392	0.396	0.436	0.355	0.327	0.394	0.407	0.390	0.330	0.300	0.266	0.766	0.676	0.437			
1.000	1.000	1.000	0.827	0.646	0.579	0.556	0.545	0.489	0.505	0.489	0.478	0.411	0.387	0.404	0.401	0.391	0.452	0.352	0.330	0.300	0.280	0.263	0.330	0.354	0.318	0.462	0.491	0.426		
0.909	1.000	0.860	0.675	0.598	0.528	0.535	0.500	0.497	0.517	0.468	0.520	0.623	0.619	0.507	0.472	0.385	0.310	0.250	0.200	0.180	0.160	0.140	0.120	0.100	0.080	0.060	0.040	0.020	0.000	
1.000	0.989	0.693	0.561	0.546	0.523	0.532	0.452	0.441	0.461	0.649	0.659	0.695	0.686	0.632	0.520	0.512	0.501	0.496	0.485	0.475	0.465	0.455	0.445	0.435	0.425	0.415	0.405	0.395	0.385	
0.969	0.849	0.606	0.530	0.521	0.494	0.437	0.396	0.421	0.626	0.698	0.741	0.737	0.730	0.729	0.728	0.727	0.726	0.725	0.724	0.723	0.722	0.721	0.720	0.719	0.718	0.717	0.716	0.715	0.714	
1.000	1.000	0.590	0.509	0.486	0.445	0.411	0.372	0.569	0.675	0.732	0.744	0.750	0.756	0.750	0.743	0.741	0.740	0.738	0.737	0.736	0.735	0.734	0.733	0.732	0.731	0.730	0.729	0.728	0.727	
1.000	0.924	0.554	0.517	0.450	0.416	0.449	0.378	0.585	0.707	0.727	0.730	0.737	0.740	0.742	0.741	0.740	0.739	0.738	0.737	0.736	0.735	0.734	0.733	0.732	0.731	0.730	0.729	0.728	0.727	
1.000	1.000	0.557	0.517	0.457	0.396	0.390	0.4	0.635	0.658	0.607	0.619	0.751	0.757	0.792	0.764	0.714	0.694	0.642	0.597	0.542	0.419	0.341	0.289	0.291	0.326	0.380	0.330	0.446		
1.000	1.000	0.556	0.4	0.42	0.385	0.52	0.623	0.63	0.670	0.711	0.748	0.771	0.775	0.772	0.724	0.594	0.4	0.434	0.378	0.354	0.414	0.307	0.282	0.278	0.402	0.306	0.290			
0.763	1.000	0.617	0.59	0.59	0.59	0.59	0.59	0.59	0.59	0.646	0.687	0.718	0.724	0.748	0.717	0.559	0.45	0.40	0.360	0.361	0.394	0.483	0.499	0.472	0.273	0.234	0.279	0.306		
1.000	1.000	0.750	0.61	0.60	0.594	0.328	0.490	0.550	0.623	0.593	0.515	0.521	0.615	0.618	0.616	0.610	0.604	0.601	0.594	0.627	0.590	0.613	0.585	0.529	0.438	0.328	0.487	0.200		
0.754	0.830	1.000	0.471	0.435	0.326	0.327	0.489	0.474	0.421	0.388	0.4	0.34	0.5	0	0.56	0.64	0.601	0.594	0.627	0.590	0.613	0.585	0.529	0.438	0.328	0.487	0.200	0.258		
0.929	0.672	0.503	0.654	0.388	0.335	0.306	0.475	0.416	0.375	0.46	0.4	0	0.74	0	0.559	0.616	0.550	0.649	0.686	0.658	0.667	0.587	0.564	0.486	0.416	0.546	0.263			
1.000	0.758	0.639	0.726	0.931	0.330	0.299	0.398	0.54	0.5	0	0.21	0.67	0.646	0.644	0.517	0.605	0.517	0.546	0.616	0.714	0.683	0.609	0.578	0.563	0.478	0.314	0.252			
1.000	0.790	0.907	0.701	0.897	0.382	0.296	0.358	0	0.63	0.628	0.674	0.683	0.666	0.605	0.526	0.620	0.527	0.514	0.616	0.666	0.670	0.628	0.549	0.512	0.262	0.321	0.254			
0.760	0.587	0.639	0.557	0.681	0.593	0.397	0.340	0.575	0.574	0.647	0.691	0.666	0.620	0.506	0.614	0.550	0.532	0.487	0.589	0.610	0.616	0.504	0.482	0.310	0.271	0.237				
0.577	0.599	0.443	0.561	0.657	0.363	0.914	0.626	0.482	0.553	0.631	0.678	0.722	0.561	0.523	0.639	0.634	0.510	0.481	0.558	0.533	0.597	0.570	0.509	0.342	0.263	0.243				
0.639	0.615	0.748	0.639	0.911	0.796	0.647	0.614	0.529	0.553	0.588	0.651	0.644	0.585	0.433	0.606	0.588	0.467	0.313	0.363	0.349	0.415	0.578	0.512	0.305	0.274	0.256				
0.569	0.661	0.486	0.605	0.448	0.494	0.705	0.730	0.579	0.532	0.526	0.623	0.518	0.387	0.310	0.338	0.466	0.378	0.559	0.479	0.444	0.430	0.494	0.465	0.232	0.248	0.237				
0.493	0.522	0.508	0.553	0.458	0.457	0.435	0.742	0.636	0.434	0.553	0.578	0.369	0.394	0.502	0.539	0.532	0.555	0.601	0.582	0.548	0.498	0.328	0.237	0.242	0.252	0.273				
0.891	0.817	0.441	0.445	0.473	0.452	0.720	0.423	0.700	0.492	0.525	0.509	0.463	0.614	0.466	0.477	0.603	0.615	0.509	0.517	0.563	0.405	0.224	0.258	0.234	0.211	0.228				
0.543	0.548	0.598	0.433	0.386	0.627	0.482	0.345	0.835	0.751	0.581	0.502	0.482	0.610	0.531	0.524	0.615	0.625	0.562	0.481	0.566	0.306	0.266	0.407	0.366	0.243	0.252				
0.762	0.720	0.506	0.496	0.495	0.698	0.396	0.627	0.555	0.317	0.491	0.294	0.382	0.393	0.572	0.449	0.405	0.407	0.357	0.567	0.518	0.243	0.255	0.465	0.415	0.323	0.248				
0.472	0.437	0.618	0.547	0.500	0.439	0.580	0.579	0.474	0.406	0.320	0.302	0.233	0.262	0.387	0.622	0.556	0.499	0.580	0.558	0.378	0.214	0.364	0.502	0.413	0.311	0.269				
0.461	0.503	0.513	0.432	0.537	0.537	0.467	0.530	0.387	0.504	0.353	0.362	0.456	0.222	0.241	0.342	0.510	0.622	0.454	0.441	0.285	0.218	0.545	0.502	0.445	0.508	0.623				
0.529	0.464	0.455	0.824	0.476	0.411	0.498	0.405	0.408	0.400	0.382	0.387	0.482	0.422	0.210	0.242	0.281	0.309	0.295	0.241	0.213	0.549	0.569	0.522	0.500	0.493	0.529				
0.383	0.458	0.482	0.370	0.384	0.361	0.400	0.391	0.320	0.319	0.425	0.377	0.433	0.528	0.497	0.285	0.247	0.198	0.226	0.410	0.570	0.597	0.576	0.588	0.531	0.493	0.546				
0.459	0.476	0.391	0.431	0.563	0.321	0.364	0.382	0.365	0.368	0.405	0.287	0.263	0.509	0.606	0.569	0.509	0.554	0.551	0.591	0.622	0.647	0.612	0.648	0.594	0.537	0.546				

What is interpretable
for humans?

Message Predictor 1.0.5.28868

Move message to folder... Only show predictions that just changed → OFF Search Stanley Clear

Folders

Original order	Subject	Predicted topic	Prediction confidence
9287	Re: Playoff Predictions	Hockey	99%
9294	Re: Schedule...	Baseball	60% ▲
9306	Paul Kuryla and Canadian Worm	Hockey	99%
9308	Re: My Predictions For 1993	Baseball	64% ▲
9312	Re: NHL Team Captains	Baseball	64% ▲
9316	Re: ugliest swing	Baseball	63% ▲
9319	Re: Octopus in Detroit?	Hockey	67% ▼
9339	Sparky Anderson Gets win #2000, Tigers beat A's	Baseball	99%
9347	Re: Goalie masks	Baseball	53%
9362	Re: Young Catchers	Baseball	82% ▲
9371	Re: Winning Streaks	Baseball	53%
9379	Royals	Baseball	64% ▲
9390	Phillies Mailing List?	Baseball	65% ▲
9410	Reds snap 5-game losing streak: RedReport 4-18	Baseball	98%
9423	Re: Juggling Dodgers	Baseball	57% ▲
9424	Re: Candlestick Park experience (long)	Baseball	99%
9433	Re: Notes on Jays vs. Indians Series	Baseball	53%
9434	Re: When did Dodgers move from NY to LA?	Baseball	53%
9439	Playoff pool	Hockey	96%
9441	Re: Hockey and the Hispanic community	Hockey	99%
9449	Re: Yooi-isms	Baseball	53%

Messages in the 'Unknown' folder

A Unknown (1,180 messages) **B** **C** **D** **E** **F**

Why Hockey?

Part 1: Important words
This message has important words about Hockey

baseball **hockey**
stanley **tiger**

The difference makes the computer think this message is 2.3 times more likely to be about Hockey than Baseball.

AND

Part 2: Folder size
The Baseball folder has more messages than the Hockey folder

Hockey: 7 **Baseball:** 8

The difference makes the computer thinks each Unknown message is 1.1 times more likely to be about Baseball than Hockey.

Important words

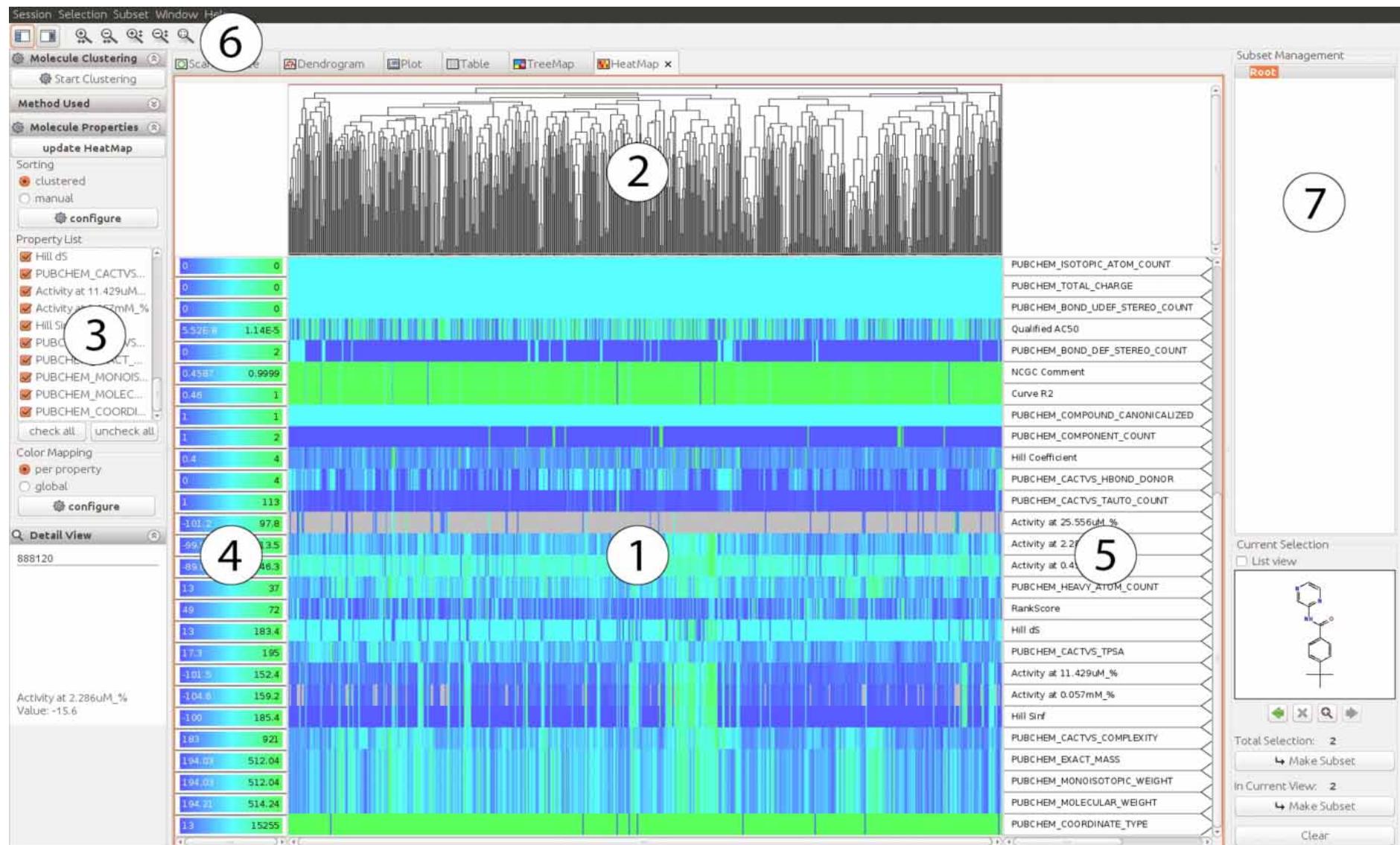
These are all of the words the computer used to make its prediction (more).

Importance

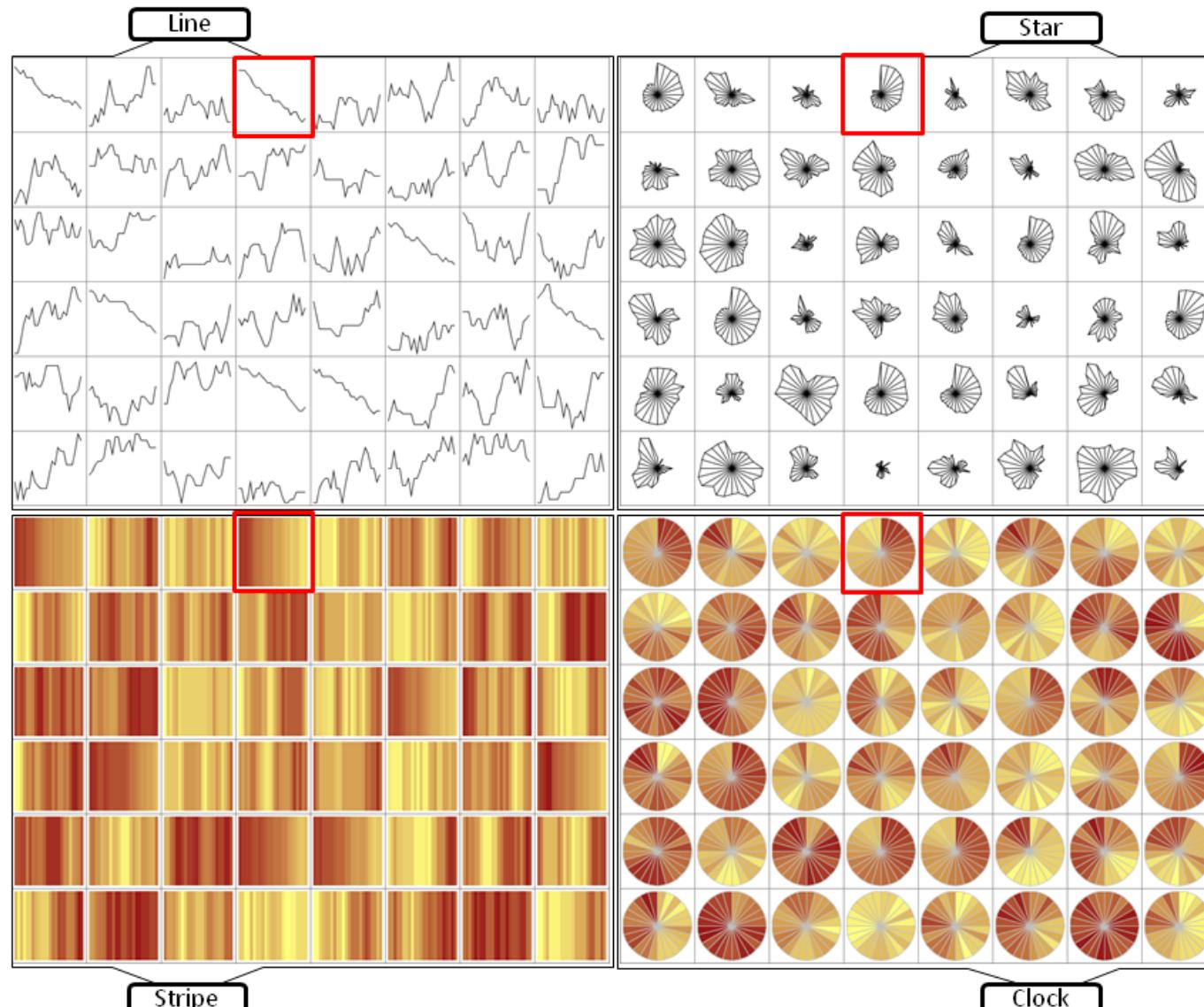
baseball, bill, canadian, dave, david, hockey, player, players, prime, stanley, stats, tiger, time

Add a new word or phrase
Remove word
Undo importance adjustment

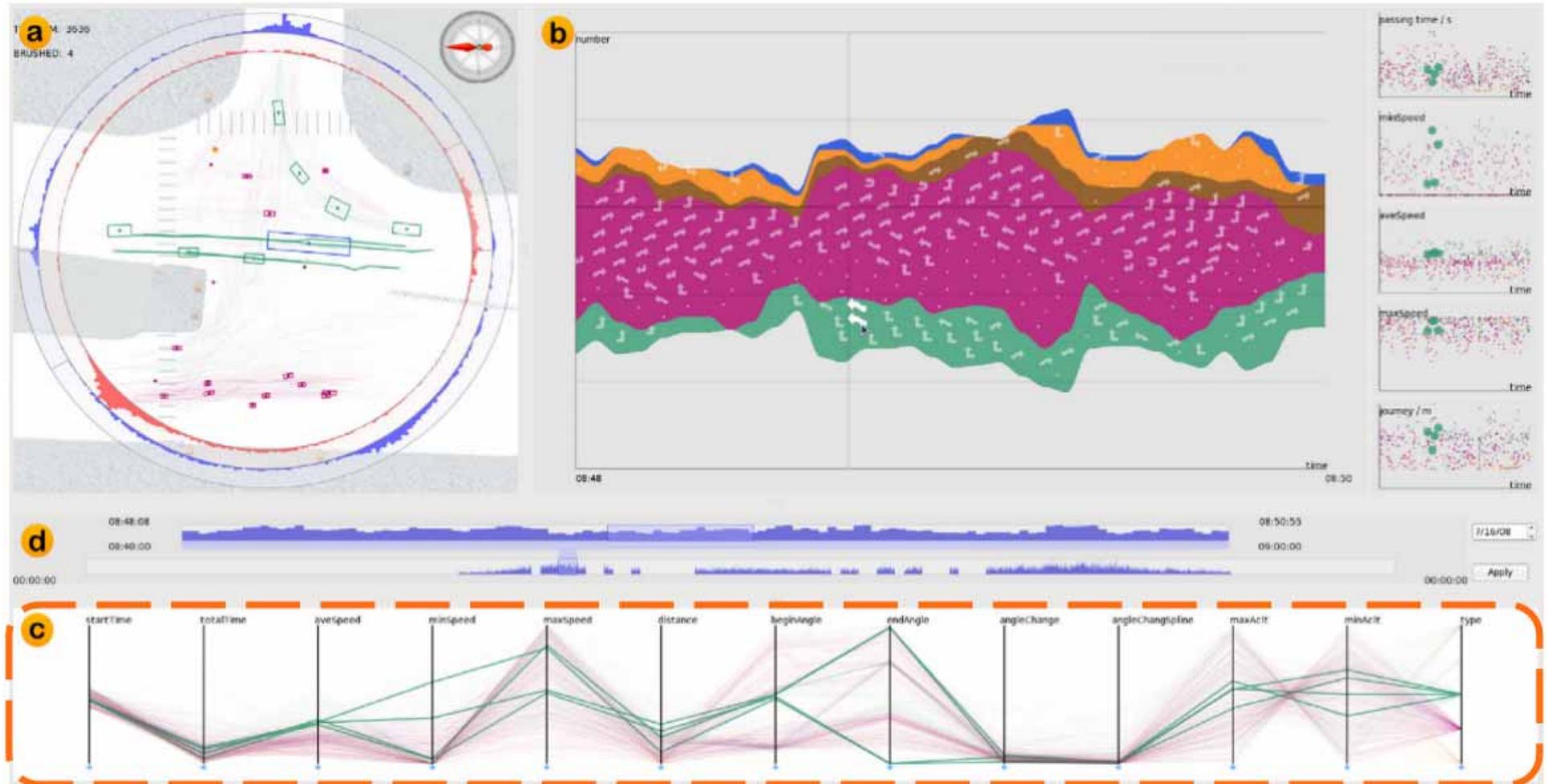
Todd Kulesza, Margaret Burnett, Weng-Keen Wong & Simone Stumpf. Principles of explanatory debugging to personalize interactive machine learning. Proceedings of the 20th International Conference on Intelligent User Interfaces (IUI 2015), 2015 Atlanta. ACM, 126-137, doi:10.1145/2678025.2701399.

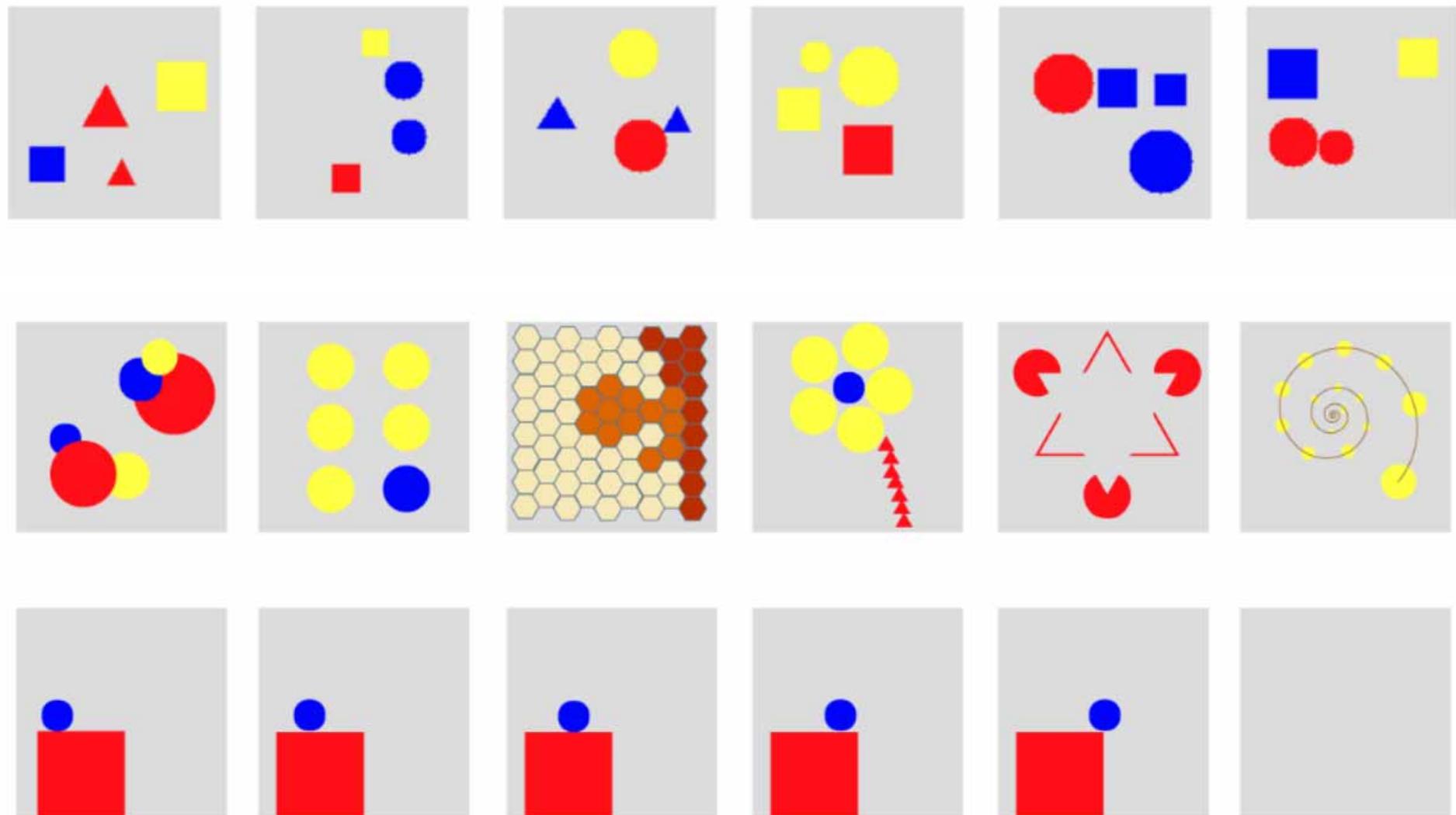


Werner Sturm, Till Schaefer, Tobias Schreck, Andreas Holzinger & Torsten Ullrich. Extending the Scaffold Hunter Visualization Toolkit with Interactive Heatmaps In: Borgo, Rita & Turkay, Cagatay, eds. EG UK Computer Graphics & Visual Computing CGVC 2015, 2015 University College London (UCL). Euro Graphics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



<https://www.vis.uni-konstanz.de/en/members/fuchs/>

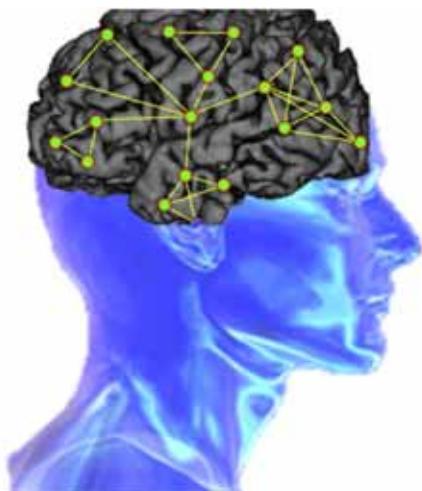




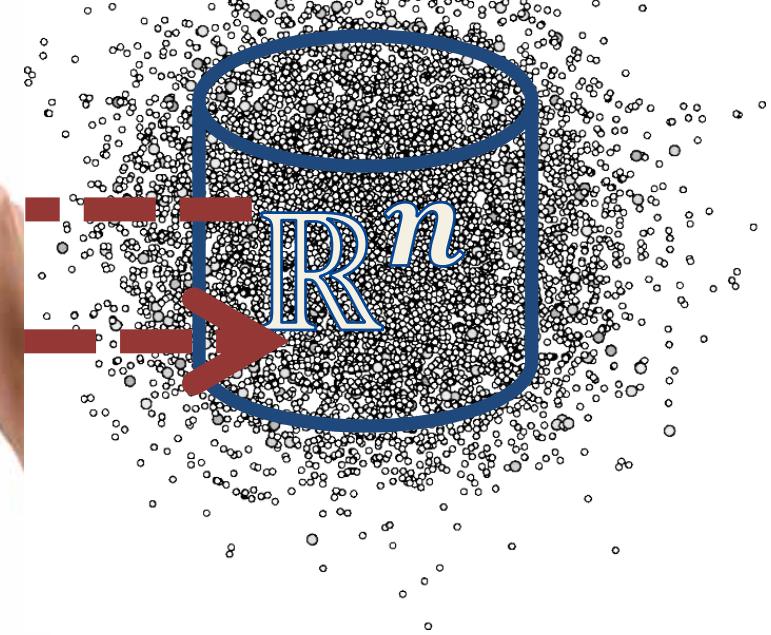
07 Methods of Explainable AI

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

Human intelligence
(Cognitive Science)

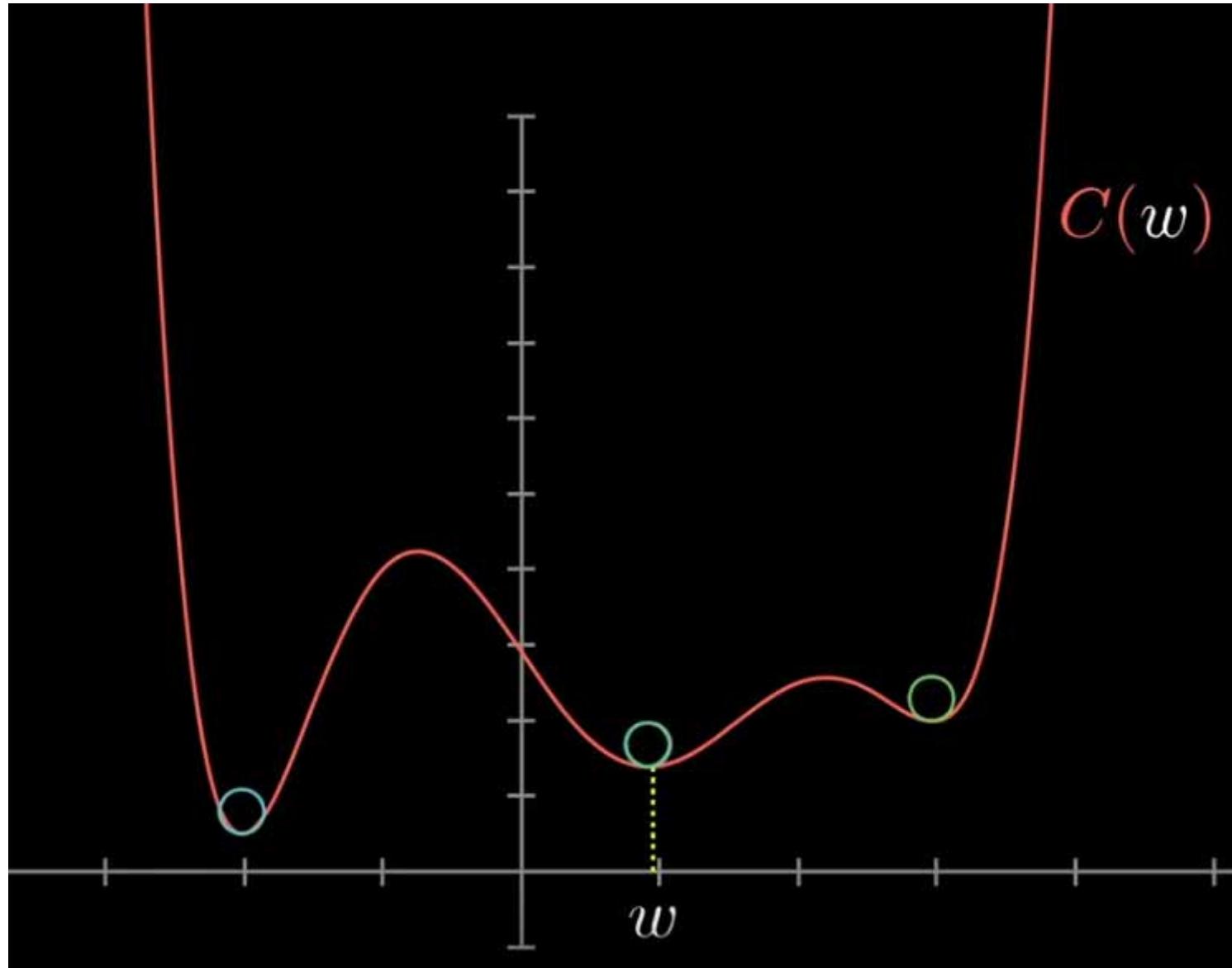


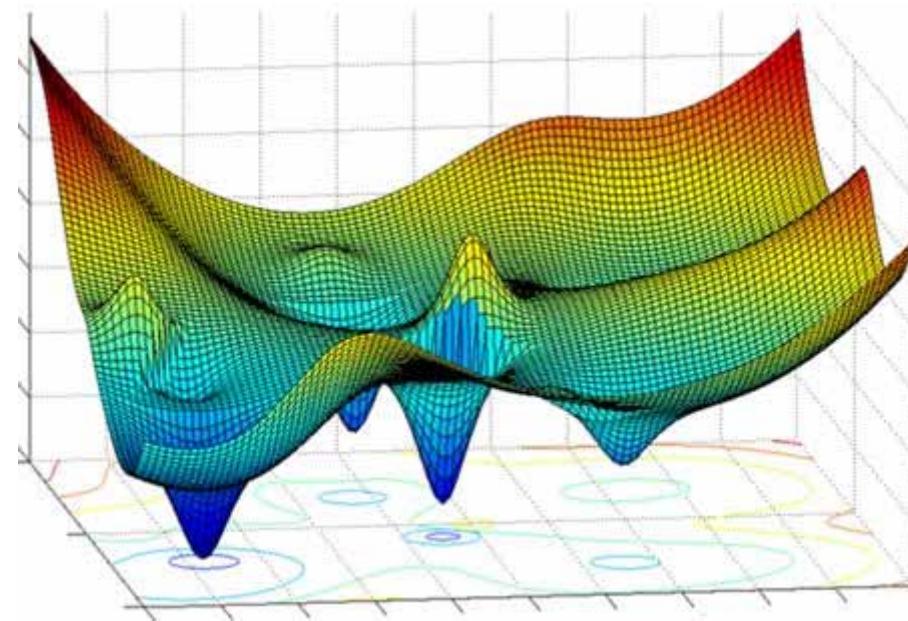
Artificial intelligence
(Computer Science)

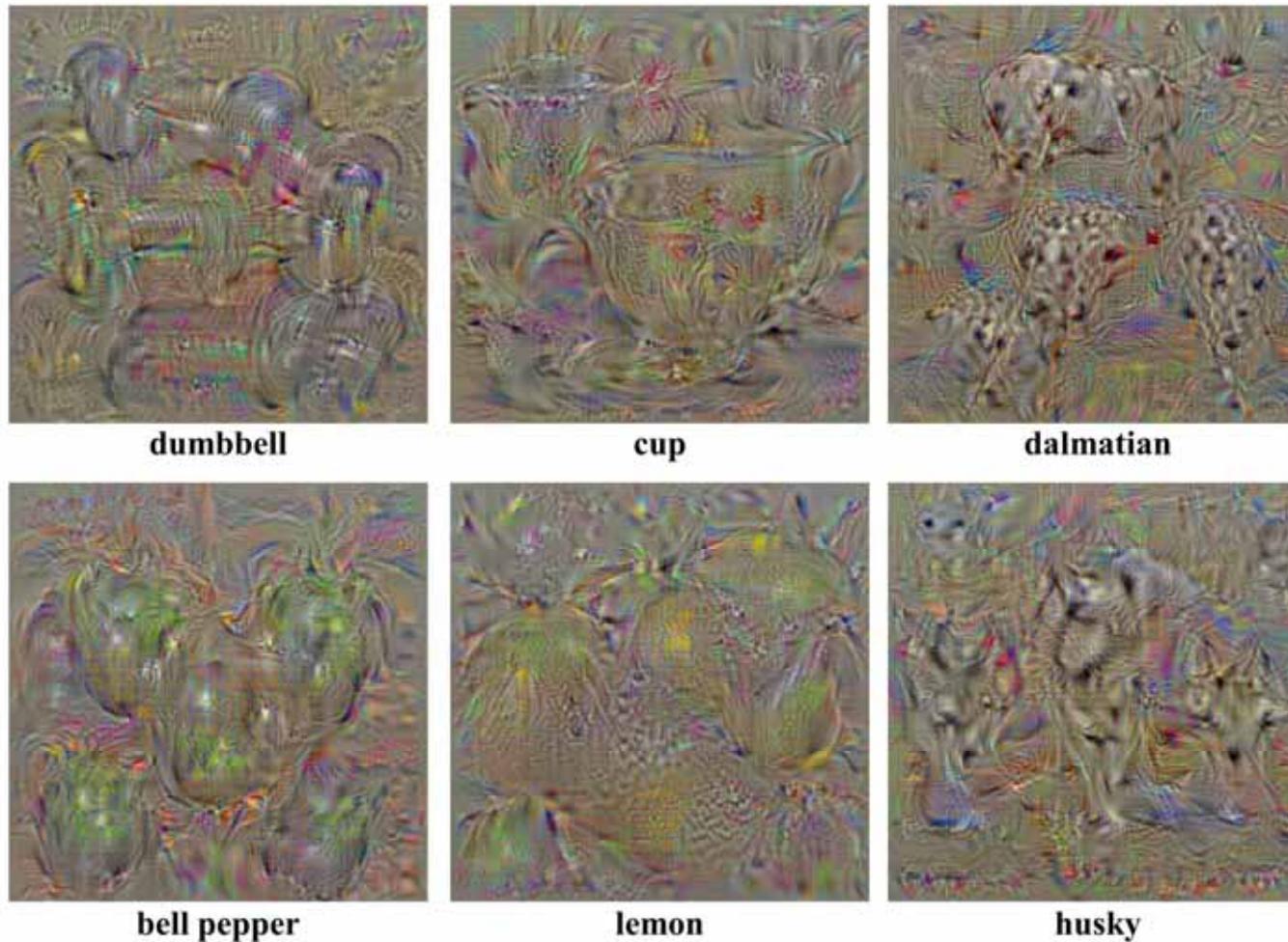


- 1) Gradients
- 2) Sensitivity Analysis
- 3) Decomposition Relevance Propagation
(Pixel-RP, Layer-RP, Deep Taylor Decomposition, ...)
- 4) Optimization (Local-IME – model agnostic,
BETA transparent approximation, ...)
- 5) Deconvolution and Guided Backpropagation
- 6) Model Understanding
 - Feature 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://hci-kdd.org/explainable-ai-causability-2019> (course given since 2016)







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

For $\text{var}_f(x_0) = k(x_0, x_0) - k_*^T(K + \Sigma)^{-1}k_*$ the derivative is given by³

$$\nabla \text{var}_f(x)|_{x=x_0} = \frac{\partial \text{var}_f}{\partial x_{0,j}} = \left(\frac{\partial}{\partial x_{0,j}} k(x_0, x_0) \right) - 2 * k_*^T(K + \Sigma)^{-1} \frac{\partial}{\partial x_{0,j}} k_* \quad \text{for } j \in \{1, \dots, d\}.$$

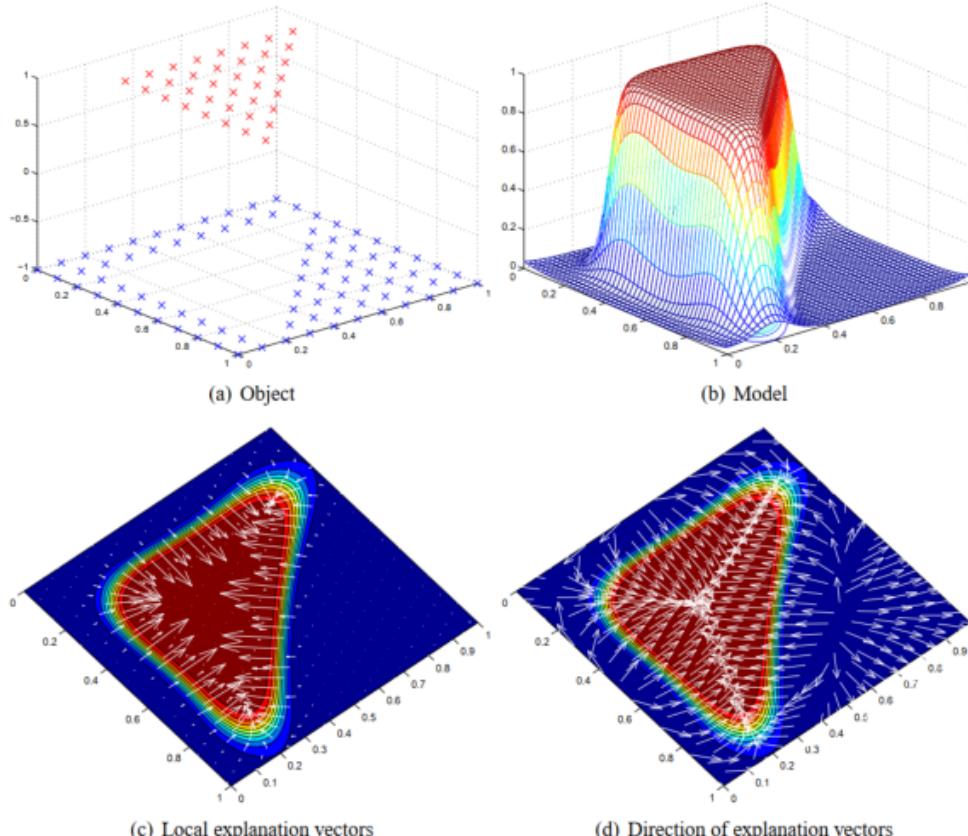


Figure 1: Explaining simple object classification with Gaussian Processes

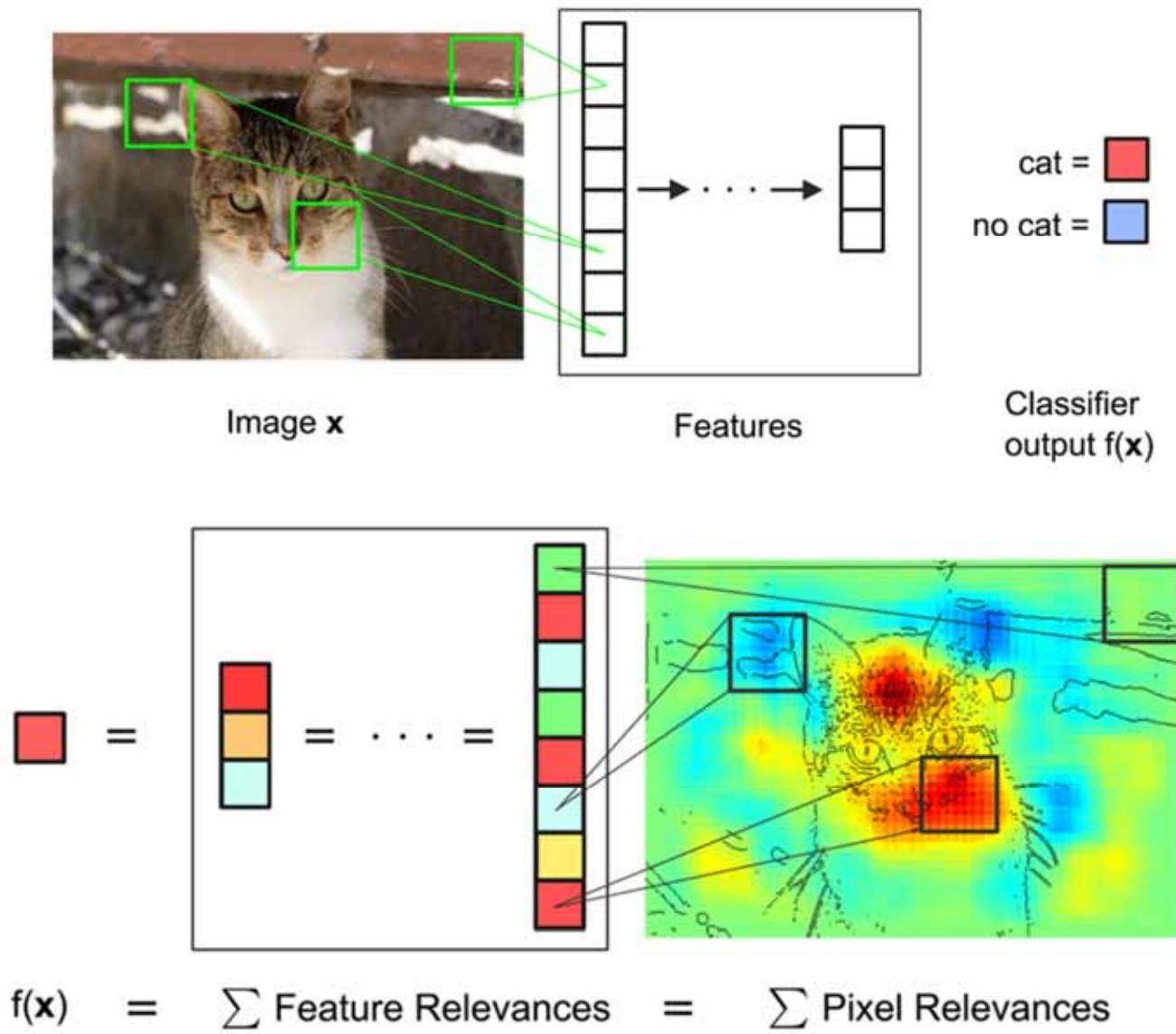
Panel (a) of Figure 1 shows the training data of a simple object classification task and panel (b) shows the model learned using GPC.⁴ The data is labeled -1 for the blue points and $+1$ for the red points. As illustrated in panel (b) the model is a probability function for the positive class which gives every data point a probability of being in this class. Panel (c) shows the probability gradient of the model together with the local gradient explanation vectors. Along the hypotenuse and at the corners of the triangle explanations from both features interact towards the triangle class while along

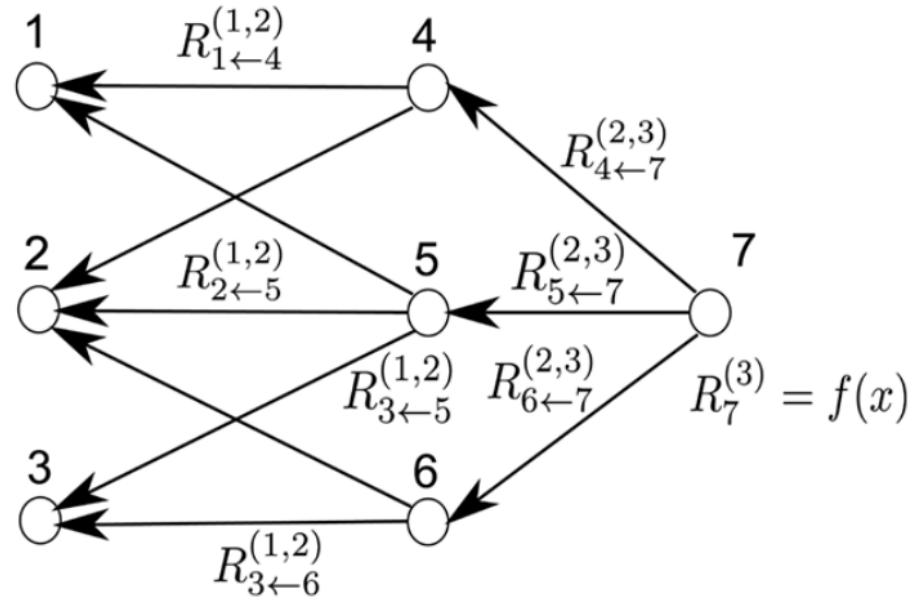
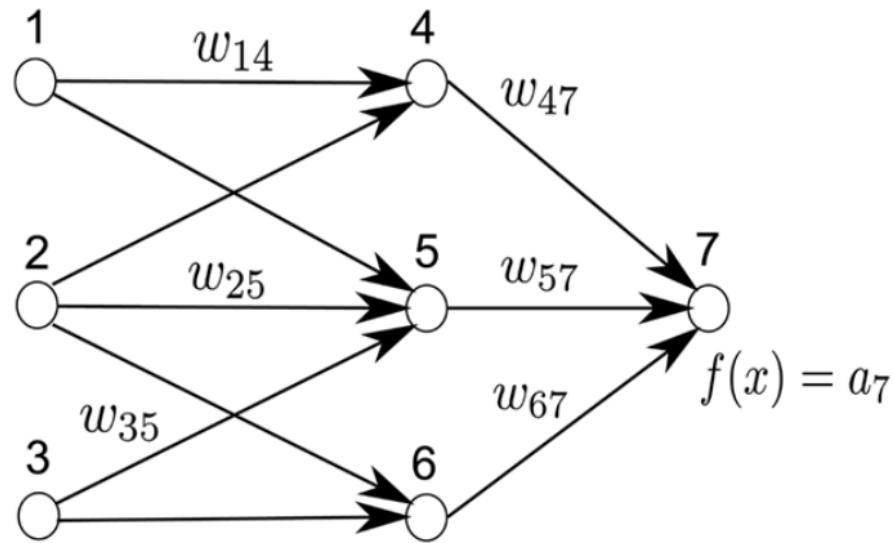
Figure 12

Robert Mueller 2010.
11, (6), 1803-1831.

Machine Learning Health 02

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.

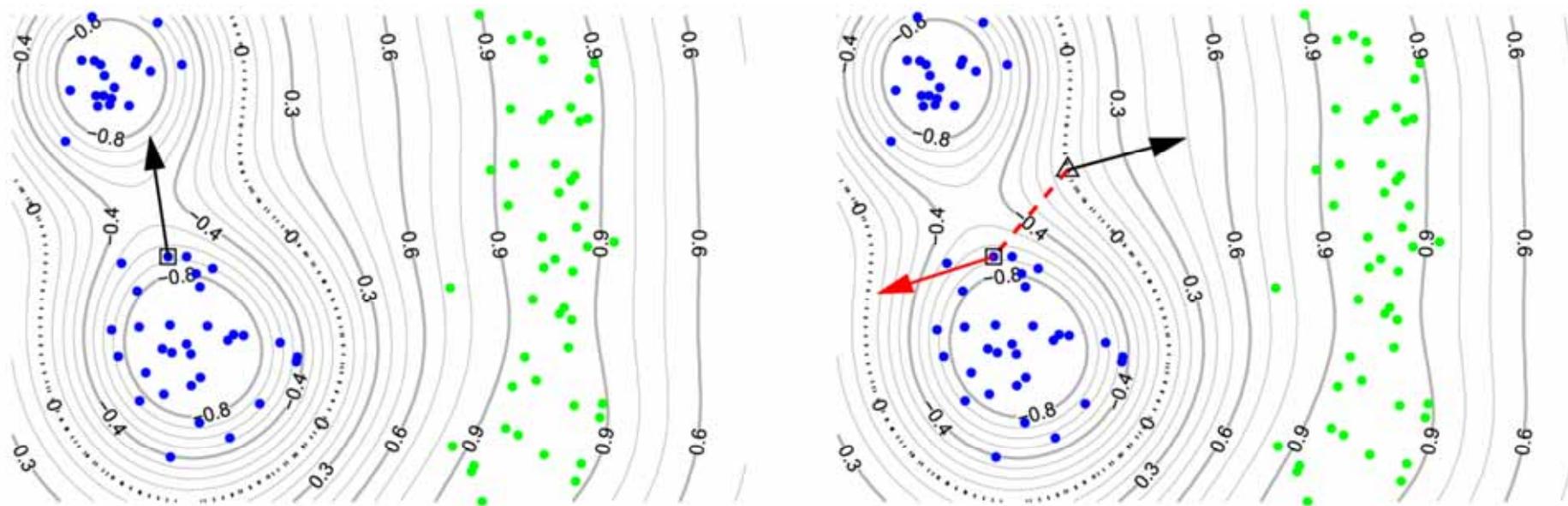


Fig 3. An exemplary real-valued prediction function for classification with the dashed black line being the decision boundary which separates the blue from the green dots. The blue dots are labeled negatively, the green dots are labeled positively. Left: Local gradient of the classification function at the prediction point. Right: Taylor approximation relative to a root point on the decision boundary. This figure depicts the intuition that a gradient at a prediction point x —here indicated by a square—does not necessarily point to a close point on the decision boundary. Instead it may point to a local optimum or to a far away point on the decision boundary. In this example the explanation vector from the local gradient at the prediction point x has a too large contribution in an irrelevant direction. The closest neighbors of the other class can be found at a very different angle. Thus, the local gradient at the prediction point x may not be a good explanation for the contributions of single dimensions to the function value $f(x)$. Local gradients at the prediction point in the left image and the Taylor root point in the right image are indicated by black arrows. The nearest root point x_0 is shown as a triangle on the decision boundary. The red arrow in the right image visualizes the approximation of $f(x)$ by Taylor expansion around the nearest root point x_0 . The approximation is given as a vector representing the dimension-wise product between $Df(x_0)$ (the black arrow in the right panel) and $x - x_0$ (the dashed red line in the right panel) which is equivalent to the diagonal of the outer product between $Df(x_0)$ and $x - x_0$.

doi:10.1371/journal.pone.0130140.g003

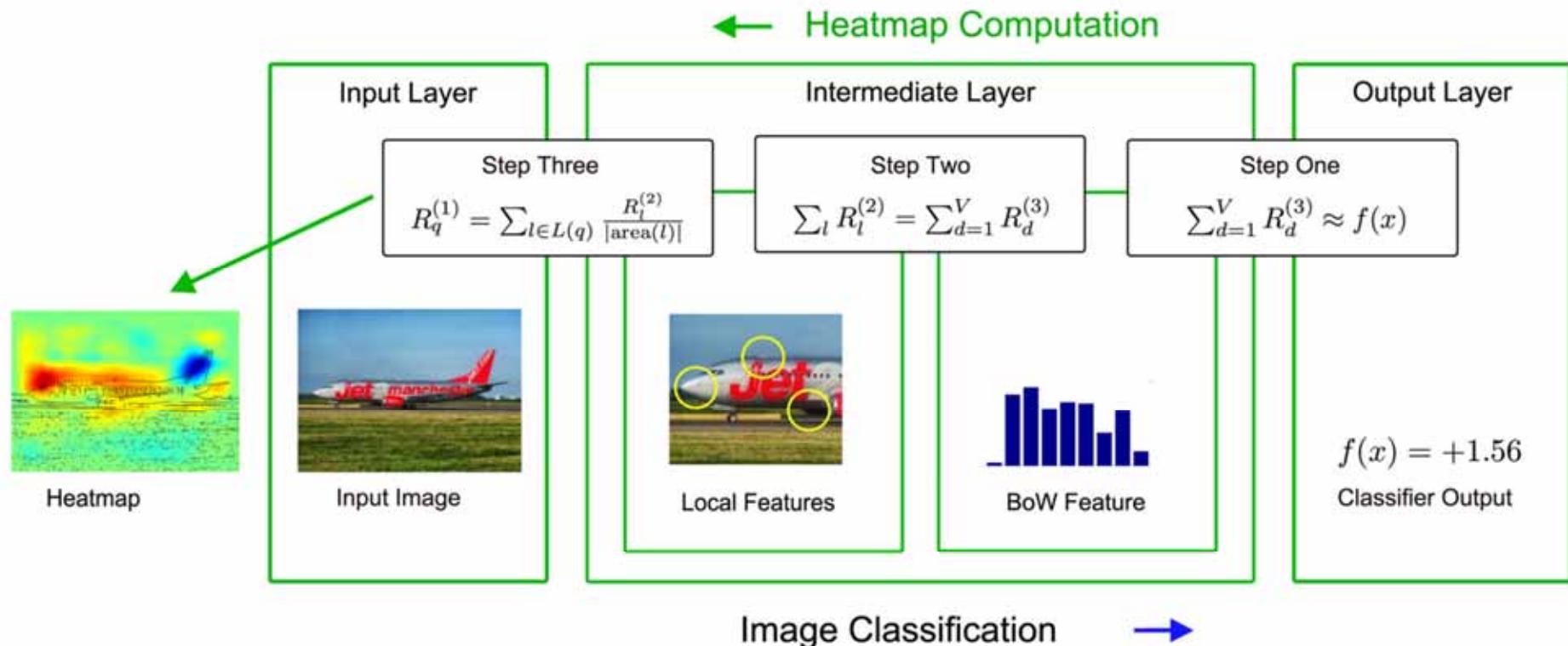
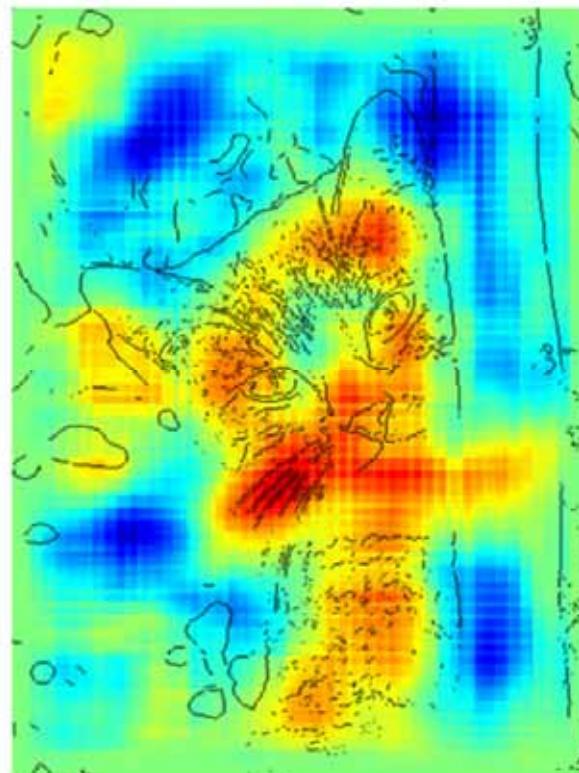


Fig 4. Local and global predictions for input images are obtained by following a series of steps through the classification- and pixel-wise decomposition pipelines. Each step taken towards the final pixel-wise decomposition has a complementing analogue within the Bag of Words classification pipeline. The calculations used during the pixel-wise decomposition process make use of information extracted by those corresponding analogues. Airplane image in the graphic by Pixabay user tpsdave.

doi:10.1371/journal.pone.0130140.g004



Definition 1. A heatmap $R(x)$ is *conservative* if the sum of assigned relevances in the pixel space corresponds to the total relevance detected by the model:

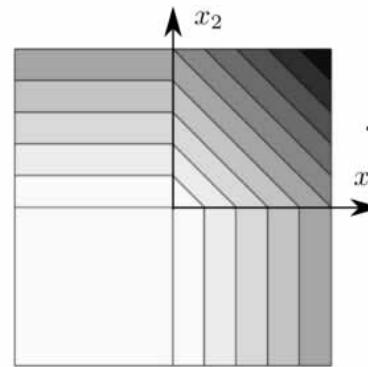
$$\forall x: f(x) = \sum_p R_p(x).$$

Definition 2. A heatmap $R(x)$ is *positive* if all values forming the heatmap are greater or equal to zero, that is:

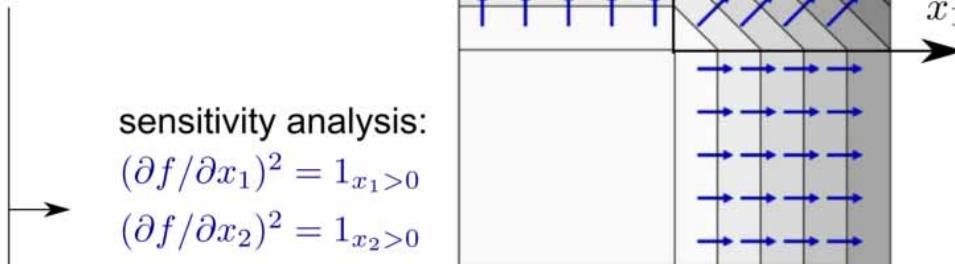
$$\forall x, p: R_p(x) \geq 0$$

Definition 3. A heatmap $R(x)$ is *consistent* if it is conservative *and* positive. That is, it is consistent if it complies with [Definitions 1 and 2](#).

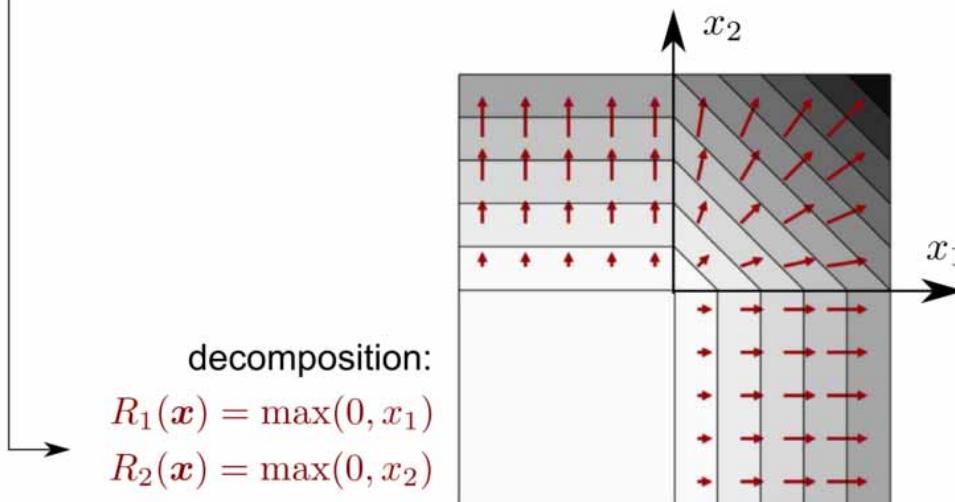
Gregoire Montavon, Sebastian Lapuschkin, Alexander Binder, Wojciech Samek & Klaus-Robert Müller 2017. Explaining nonlinear classification decisions with deep taylor decomposition. Pattern Recognition, 65, 211-222, doi:10.1016/j.patcog.2016.11.008.



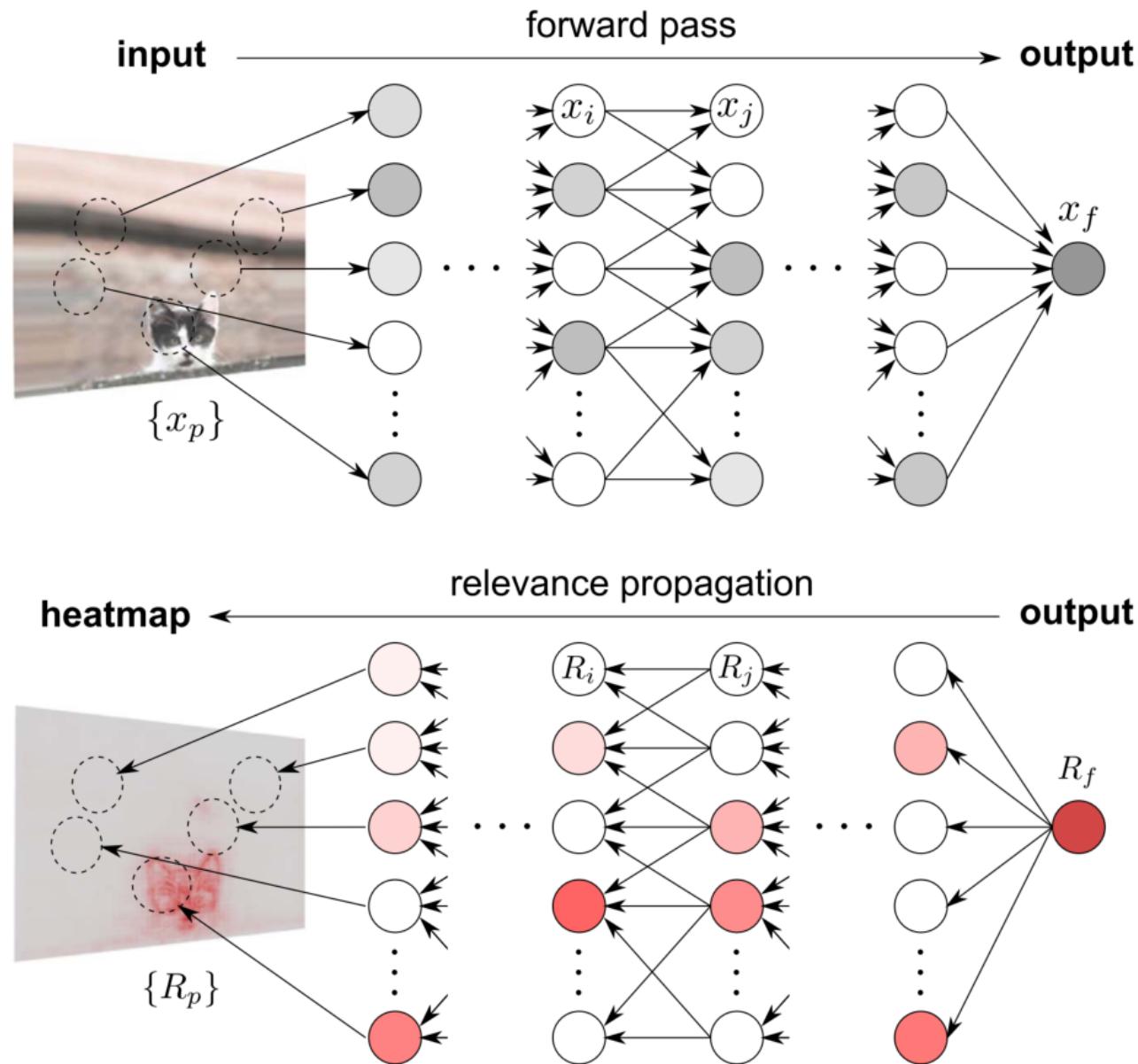
function to analyze:
 $f(\mathbf{x}) = \max(0, x_1) + \max(0, x_2)$

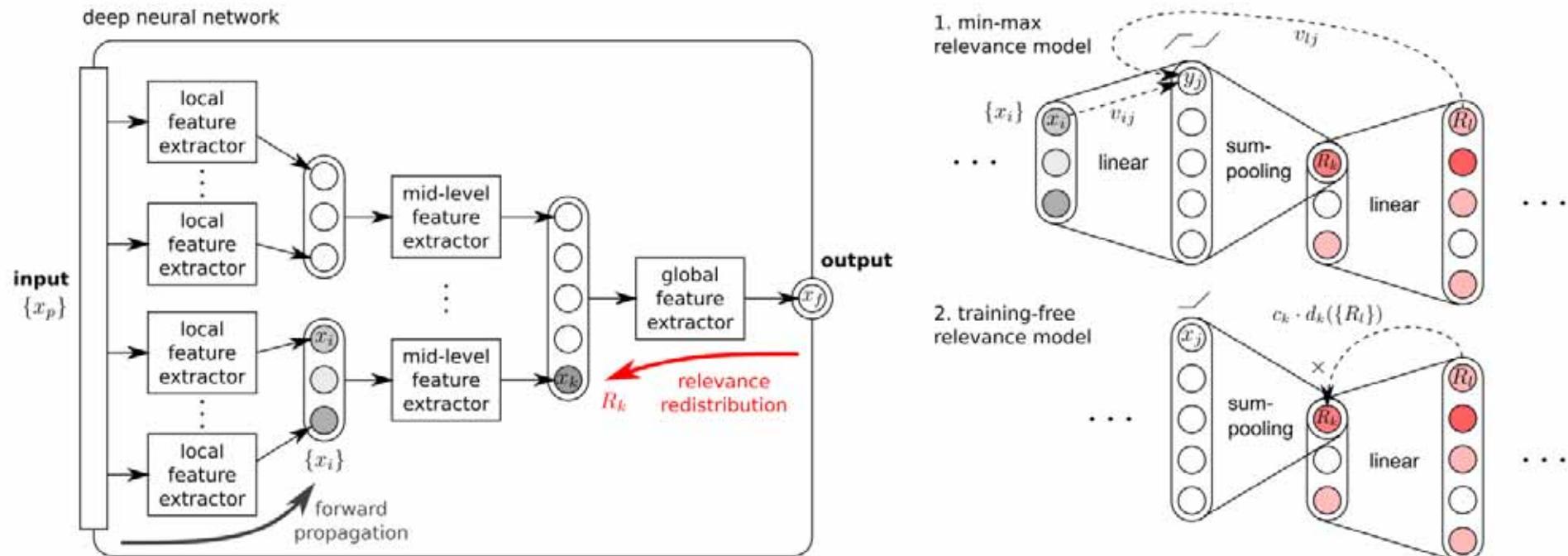


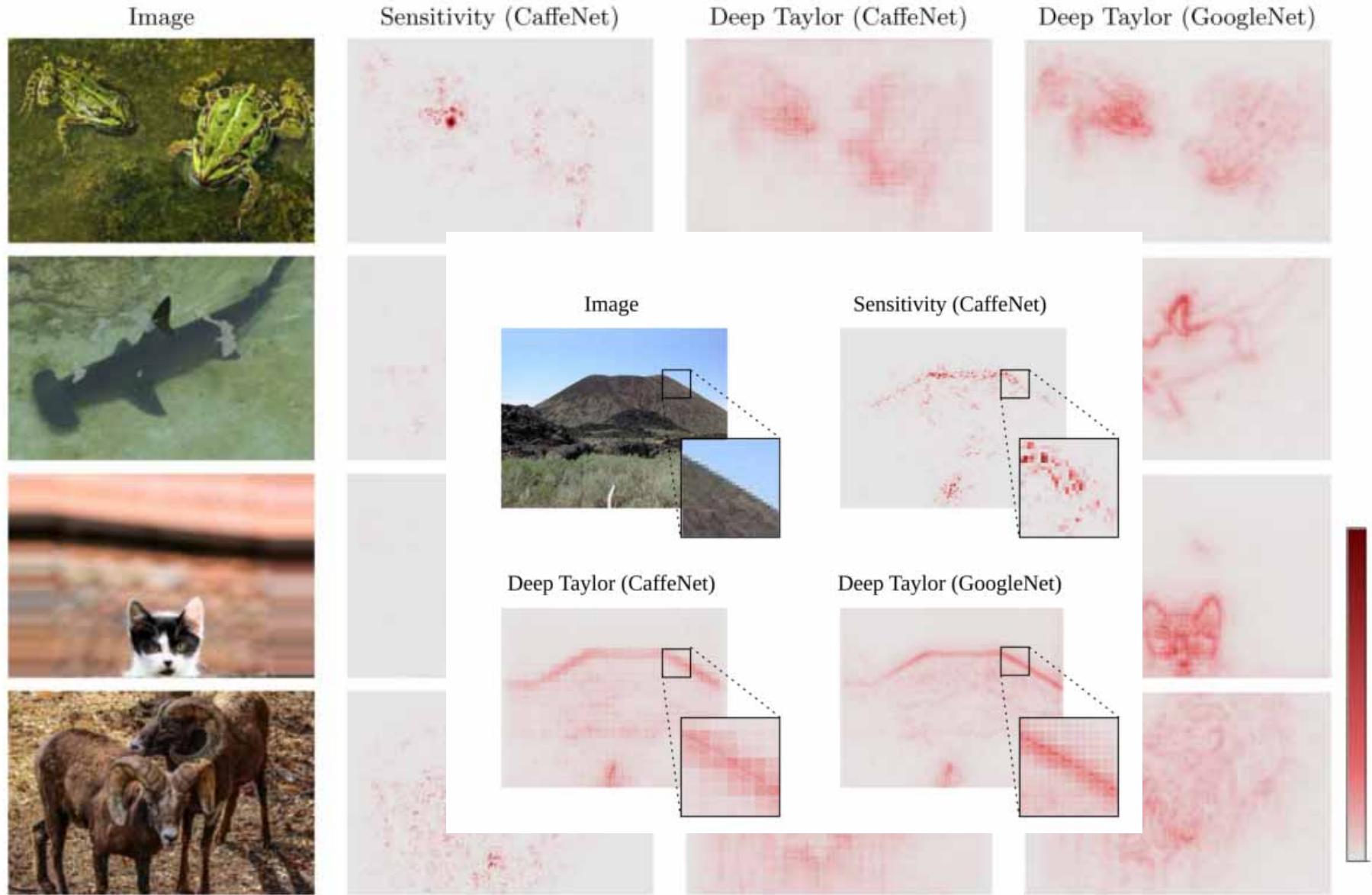
sensitivity analysis:
 $(\partial f / \partial x_1)^2 = 1_{x_1 > 0}$
 $(\partial f / \partial x_2)^2 = 1_{x_2 > 0}$

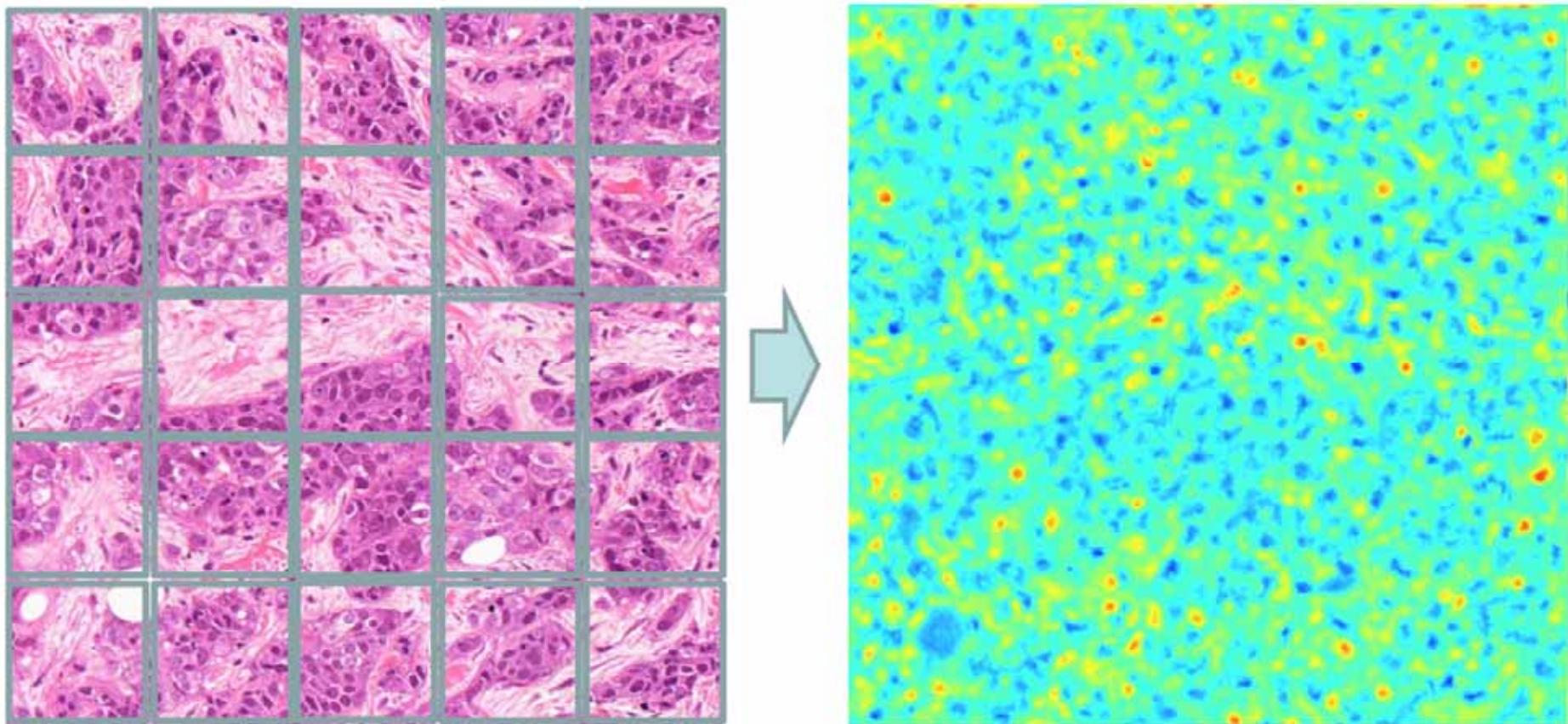


decomposition:
 $R_1(\mathbf{x}) = \max(0, x_1)$
 $R_2(\mathbf{x}) = \max(0, x_2)$

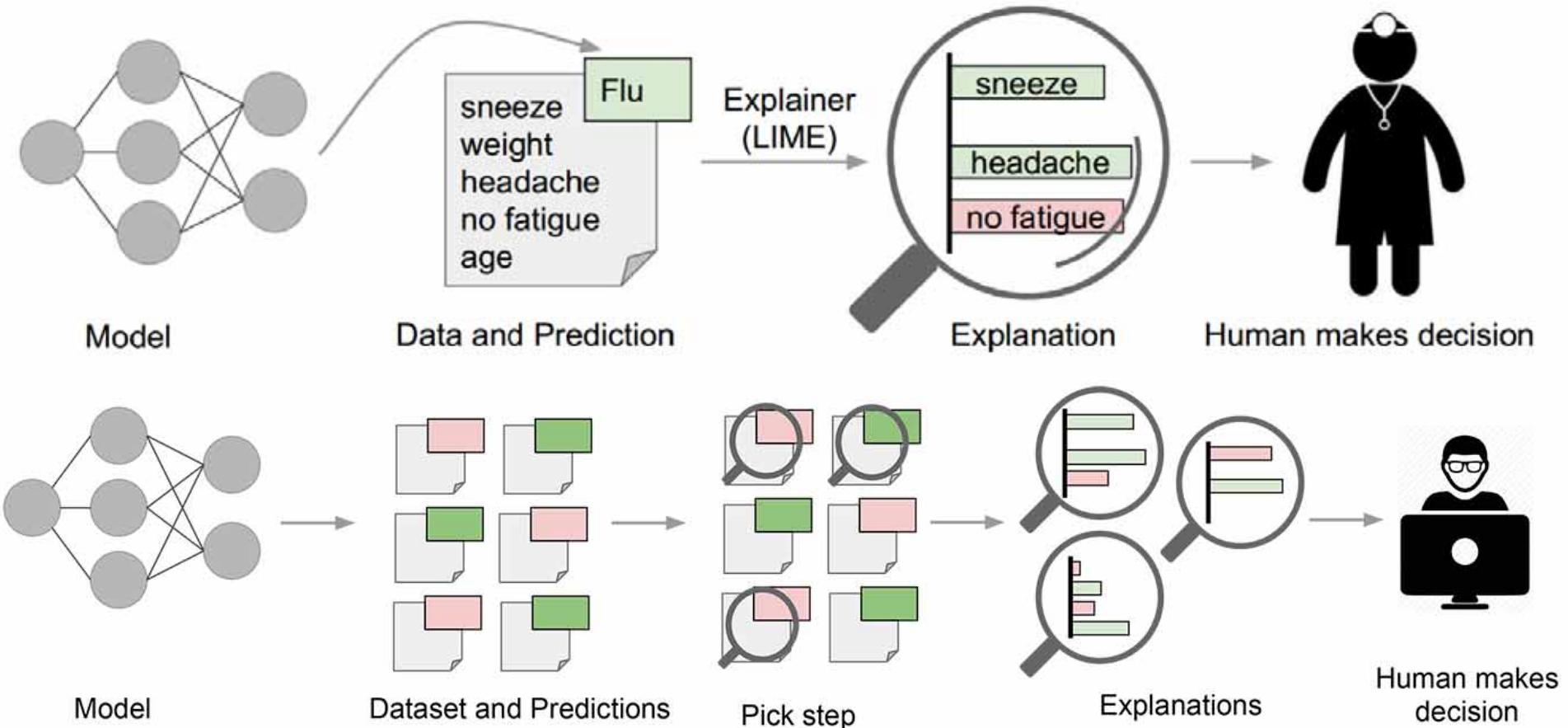








Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller & Wojciech Samek
2015. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. PloS one,
10, (7), e0130140, doi:10.1371/journal.pone.0130140.



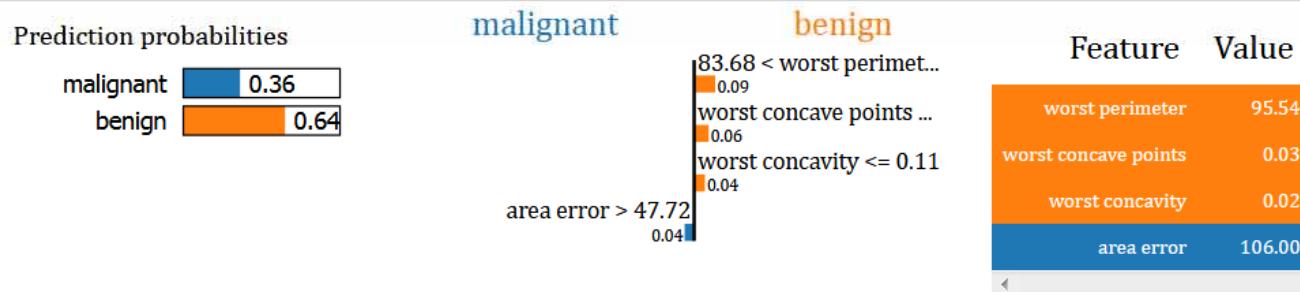
Marco Tulio Ribeiro, Sameer Singh & Carlos Guestrin. Why should i trust you?: Explaining the predictions of any classifier. 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2016. ACM, 1135-1144, doi:10.1145/2939672.2939778.

```
In [12]: explainer = lime.lime_tabular.LimeTabularExplainer(X_train, feature_names=breast.feature_names, class_names=breast.targe
```

Here we will take a sample from the test set (in this case the sample at index 76) and create an explainer instance for this sample. This will let us see why the algorithm made its prediction visually.

```
In [18]: # For this demonstration, let's take the same sample each time, in this case sample index 86
i = 76
# For a random sample uncomment out the following line
# i = np.random.randint(0, X_test.shape[0])

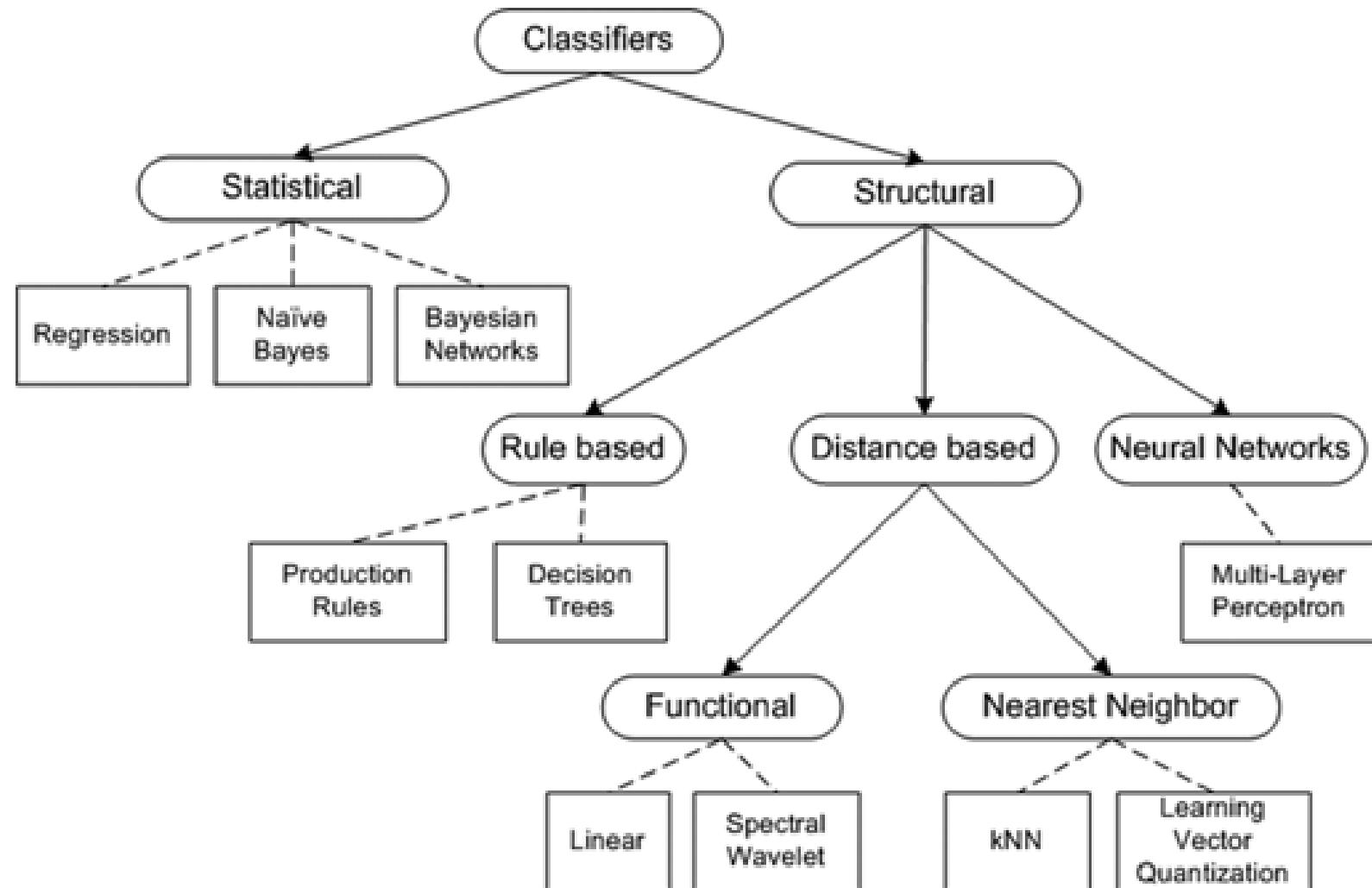
exp = explainer.explain_instance(X_test[i], random_forest.predict_proba, num_features=4)
exp.show_in_notebook(show_table=True, show_all=False)
```



As you can see, the random forest algorithm has predicted with a probability of 0.64 that the sample at index 76 in the test set is malignant.

When using the explainer, we set the `num_features` parameter to 4, meaning the explainer shows the top 4 features that contributed to the prediction probabilities.

We chose 76 as it was a borderline decision. For example sample 86 is much more clear (this will we will set the `num_features` parameter to include all features so that we see each feature's contribution to the probability):



<https://stats.stackexchange.com/questions/271247/machine-learning-statistical-vs-structural-classifiers>

If Age < 50 and Male = Yes:

If Past-Depression = Yes and Insomnia = No and Melancholy = No, then Healthy

If Past-Depression = Yes and Insomnia = Yes and Melancholy = Yes and Tiredness = Yes, then Depression

If Age \geq 50 and Male = No:

If Family-Depression = Yes and Insomnia = No and Melancholy = Yes and Tiredness = Yes, then Depression

If Family-Depression = No and Insomnia = No and Melancholy = No and Tiredness = No, then Healthy

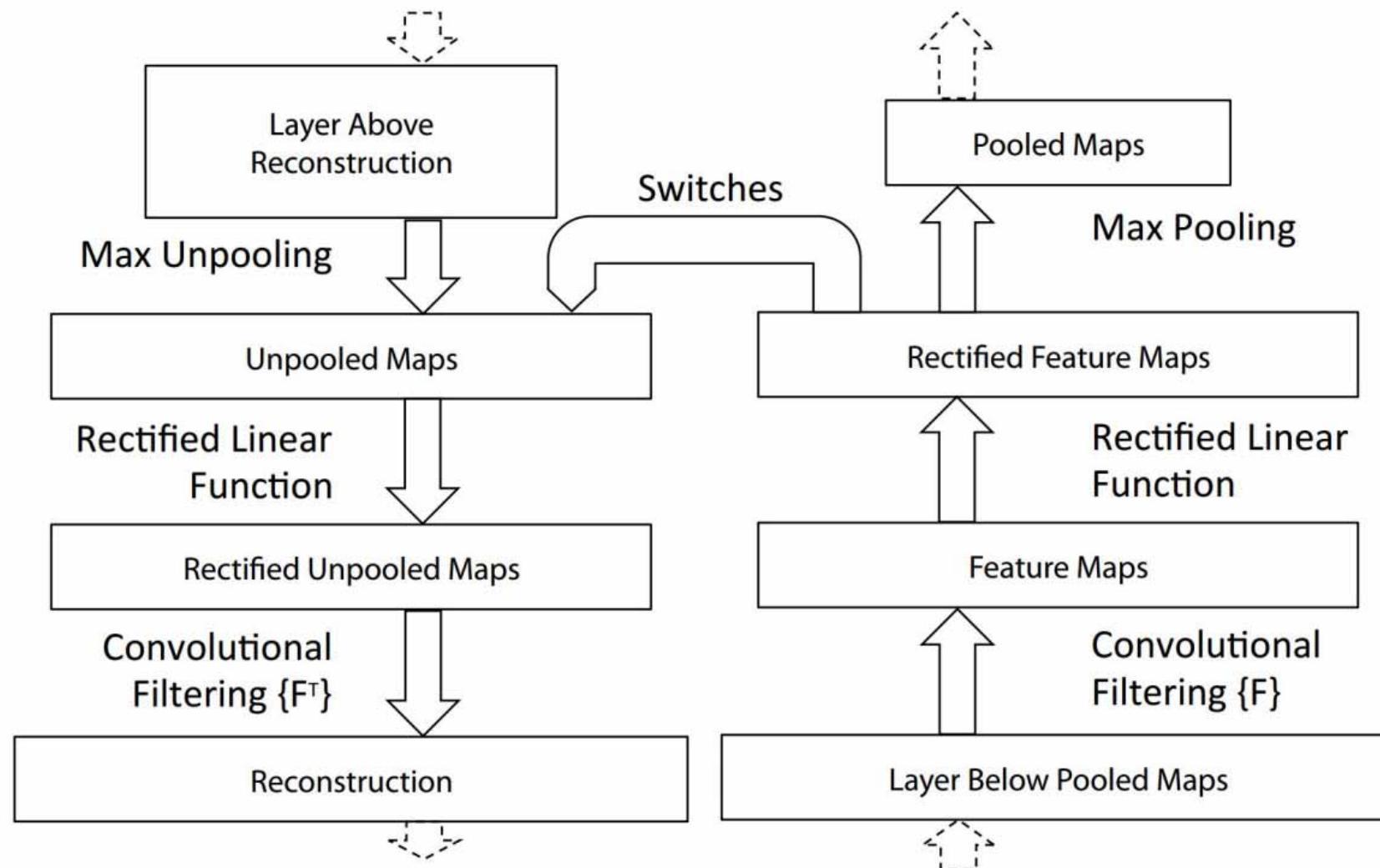
Default:

If Past-Depression = Yes and Tiredness = No and Exercise = No and Insomnia = Yes, then Depression

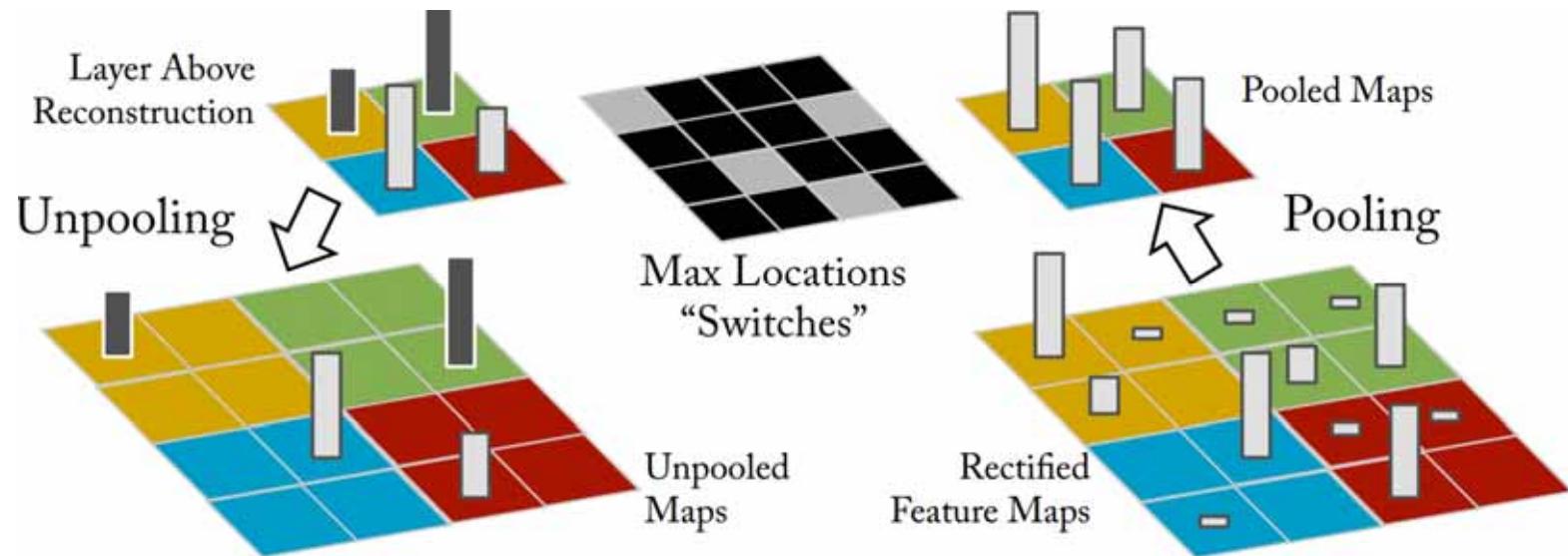
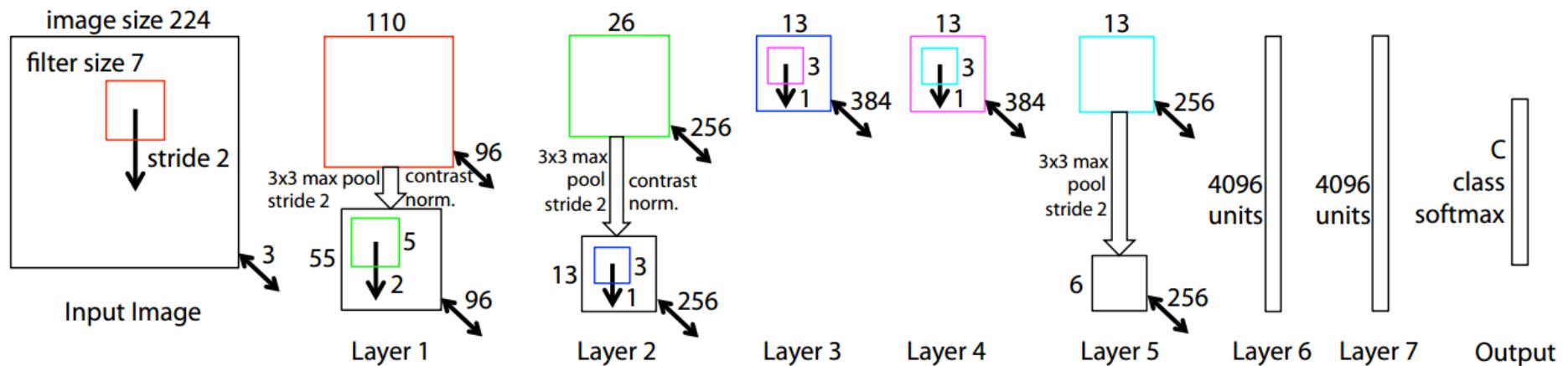
If Past-Depression = No and Weight-Gain = Yes and Tiredness = Yes and Melancholy = Yes, then Depression

If Family-Depression = Yes and Insomnia = Yes and Melancholy = Yes and Tiredness = Yes, then Depression

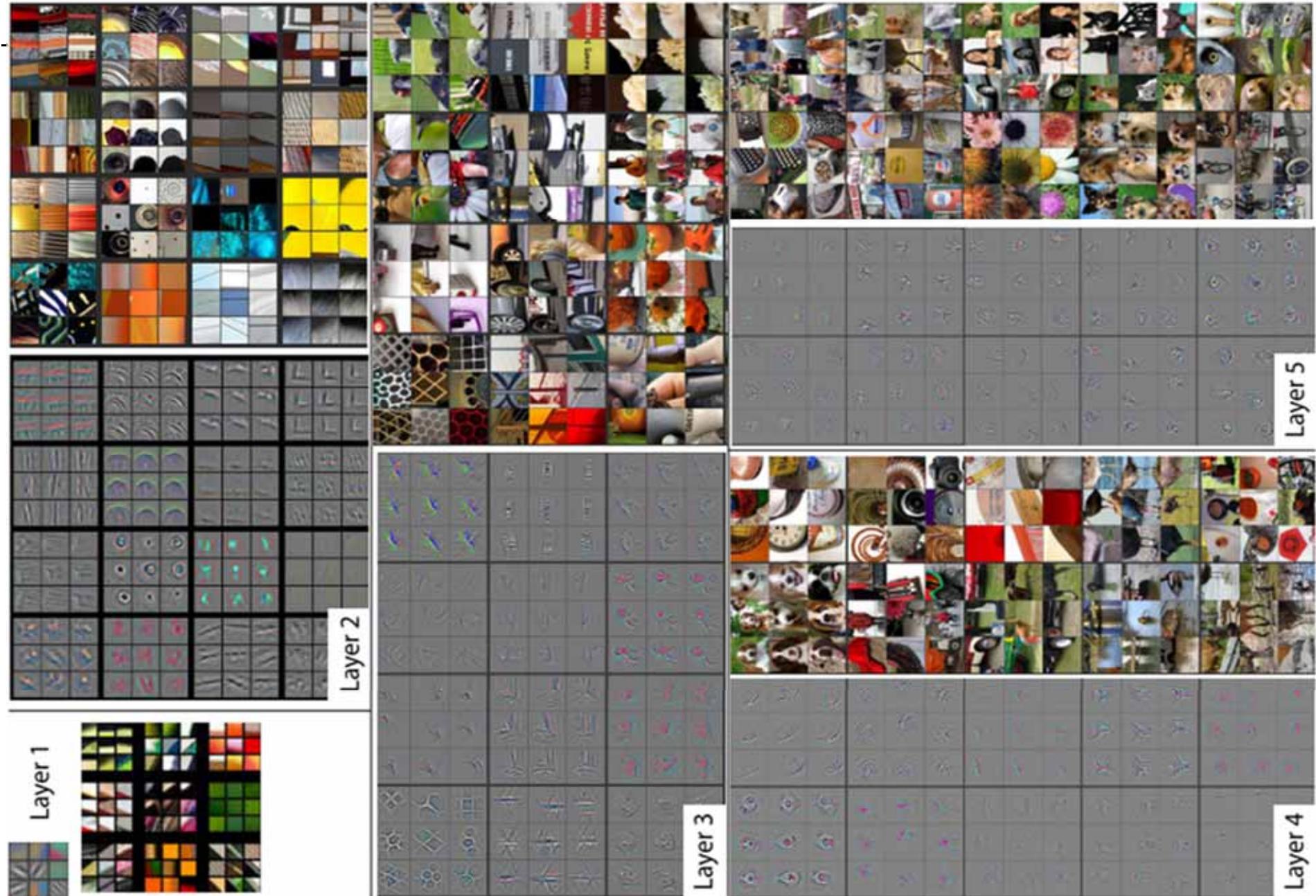
Himabindu Lakkaraju, Ece Kamar, Rich Caruana & Jure Leskovec 2017. Interpretable and Explorable Approximations of Black Box Models. arXiv:1707.01154.



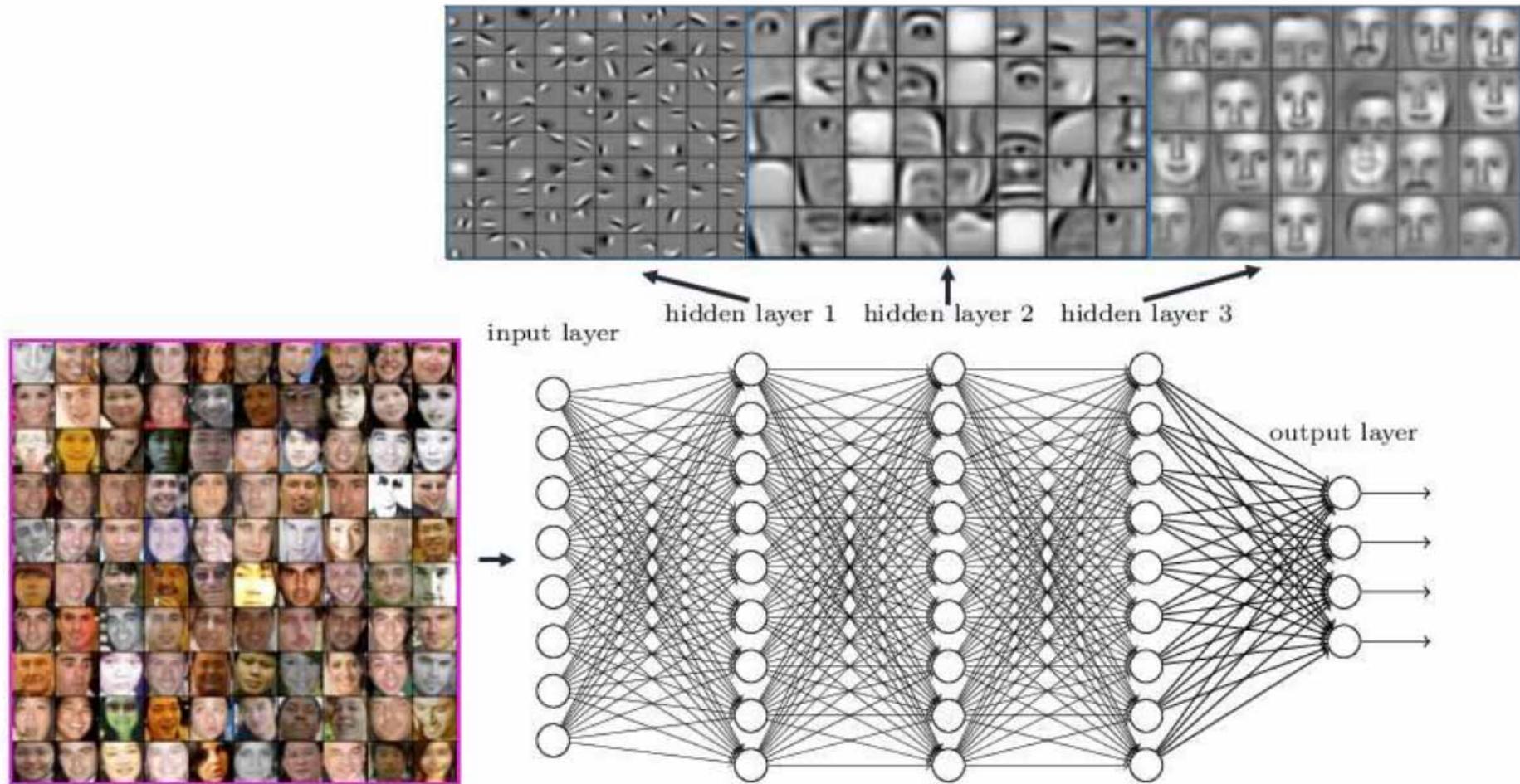
Matthew D. Zeiler & Rob Fergus 2013. Visualizing and Understanding Convolutional Networks. arXiv:1311.2901.



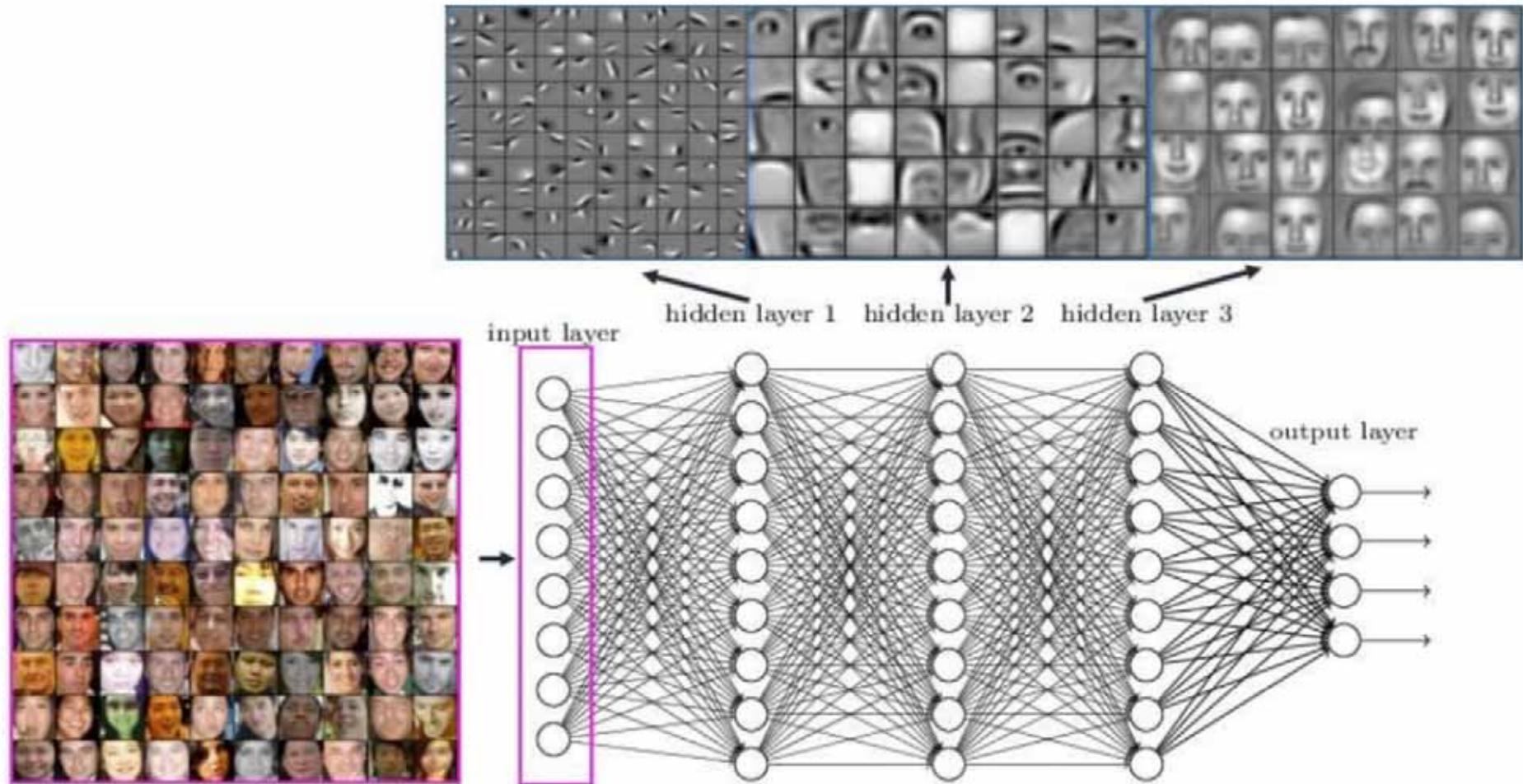
Matthew D. Zeiler & Rob Fergus 2014. Visualizing and understanding convolutional networks. In: D., Fleet, T., Pajdla, B., Schiele & T., Tuytelaars (eds.) ECCV, Lecture Notes in Computer Science LNCS 8689. Cham: Springer, pp. 818-833, doi:10.1007/978-3-319-10590-1_53.

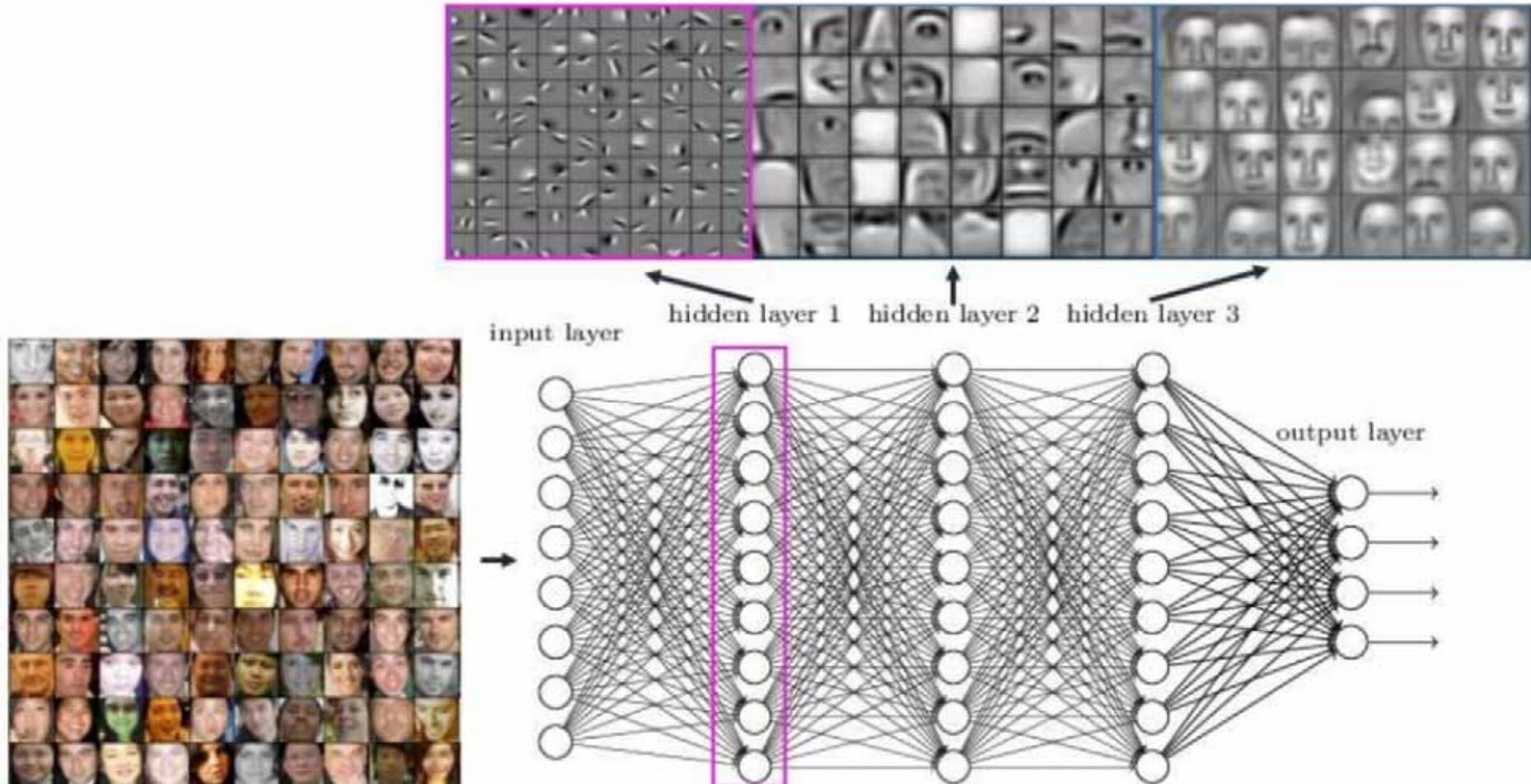


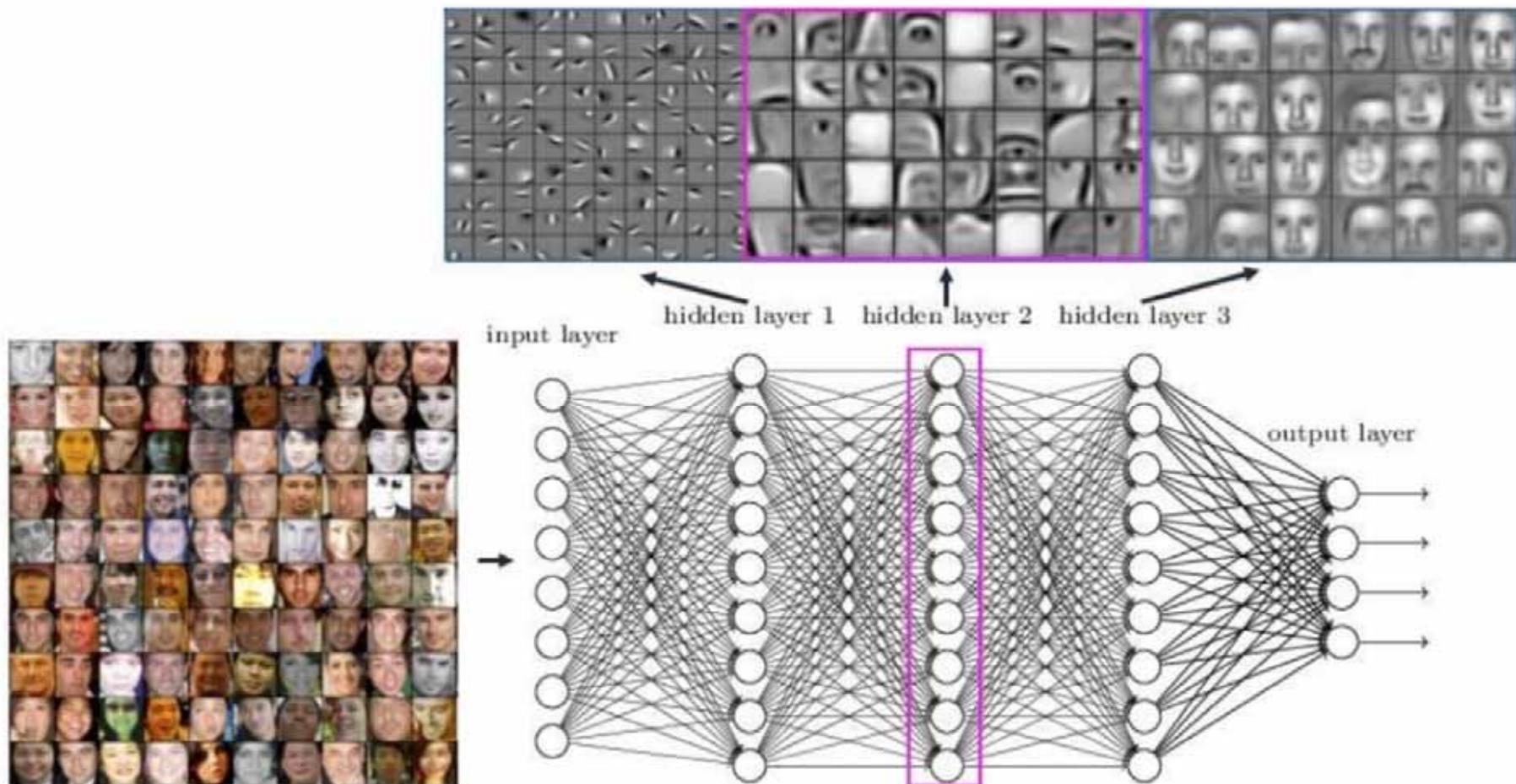
Matthew D. Zeiler & Rob Fergus 2013. Visualizing and Understanding Convolutional Networks. arXiv:1311.2901.
Holzinger Group hci-kdd.org

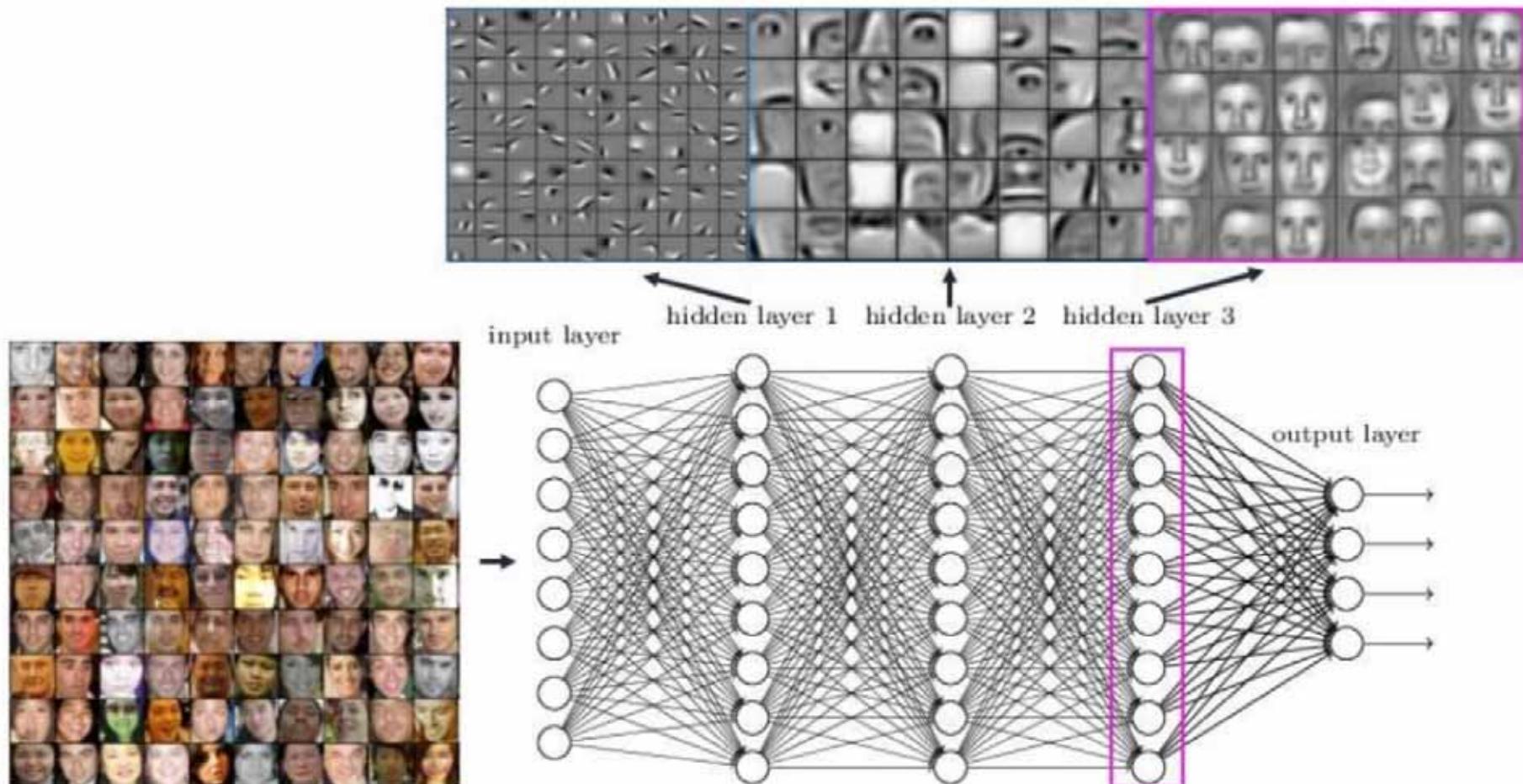


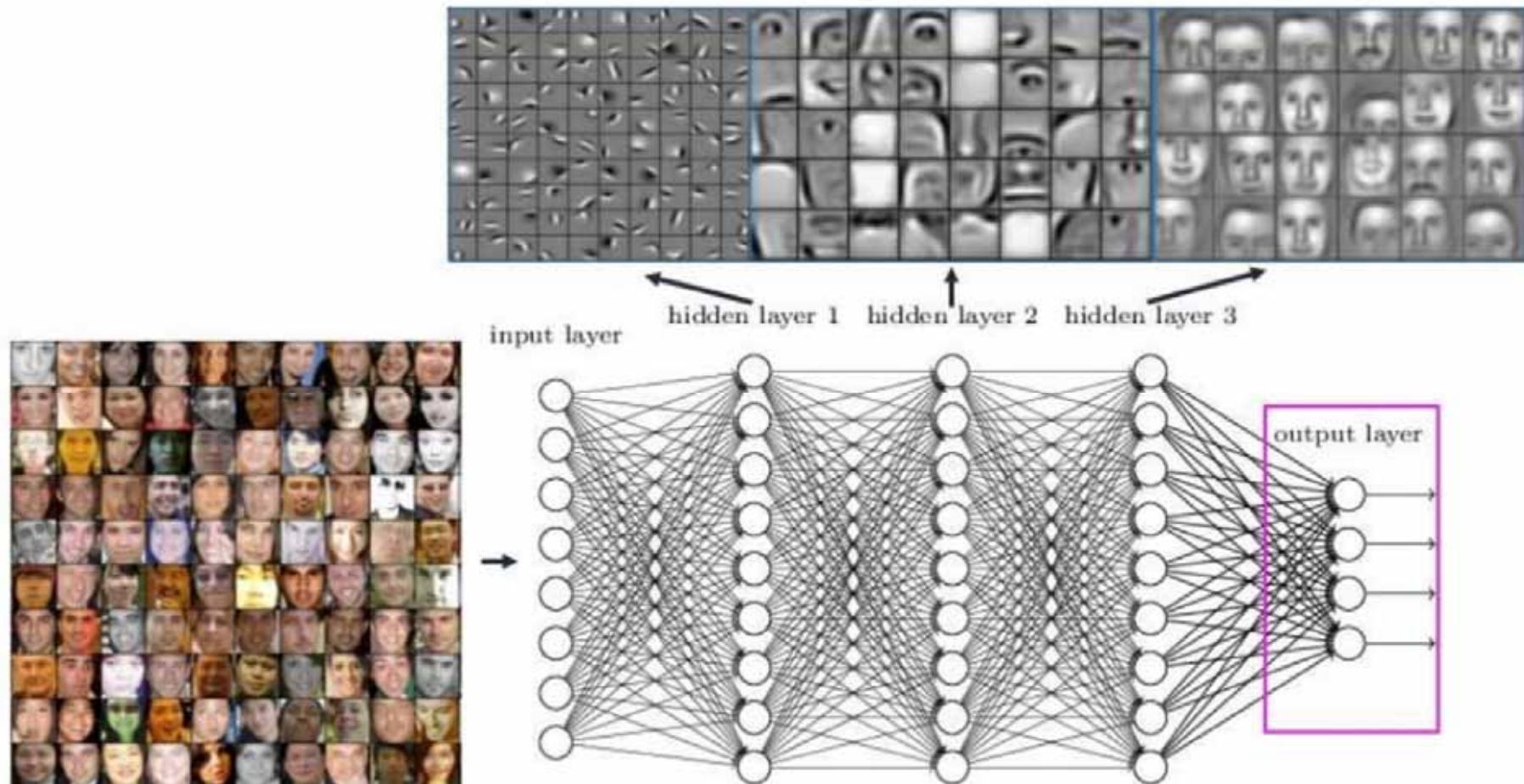
Matthew D. Zeiler & Rob Fergus 2013. Visualizing and Understanding Convolutional Networks. arXiv:1311.2901

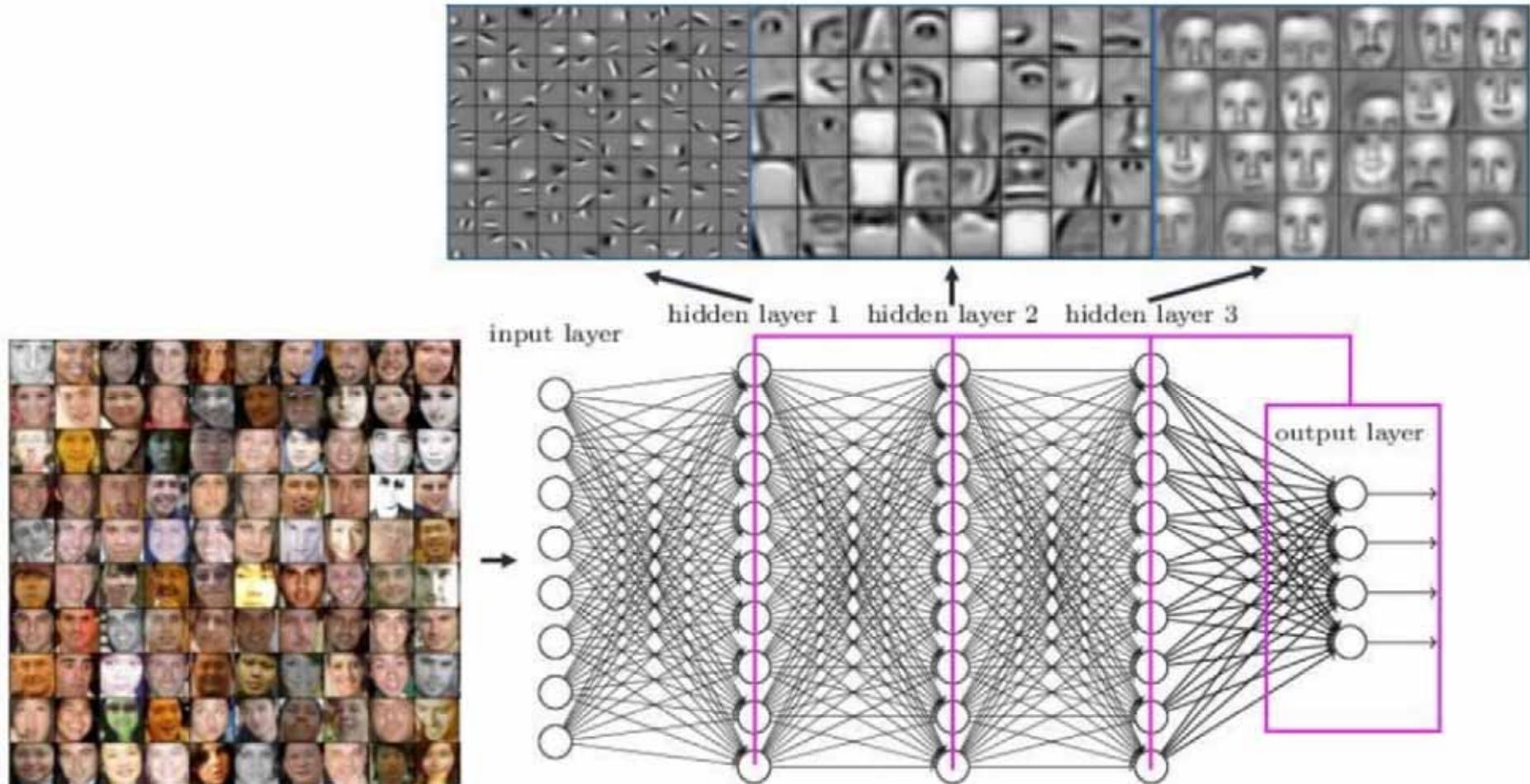












Testing with Concept Activation Vectors (TCAV)

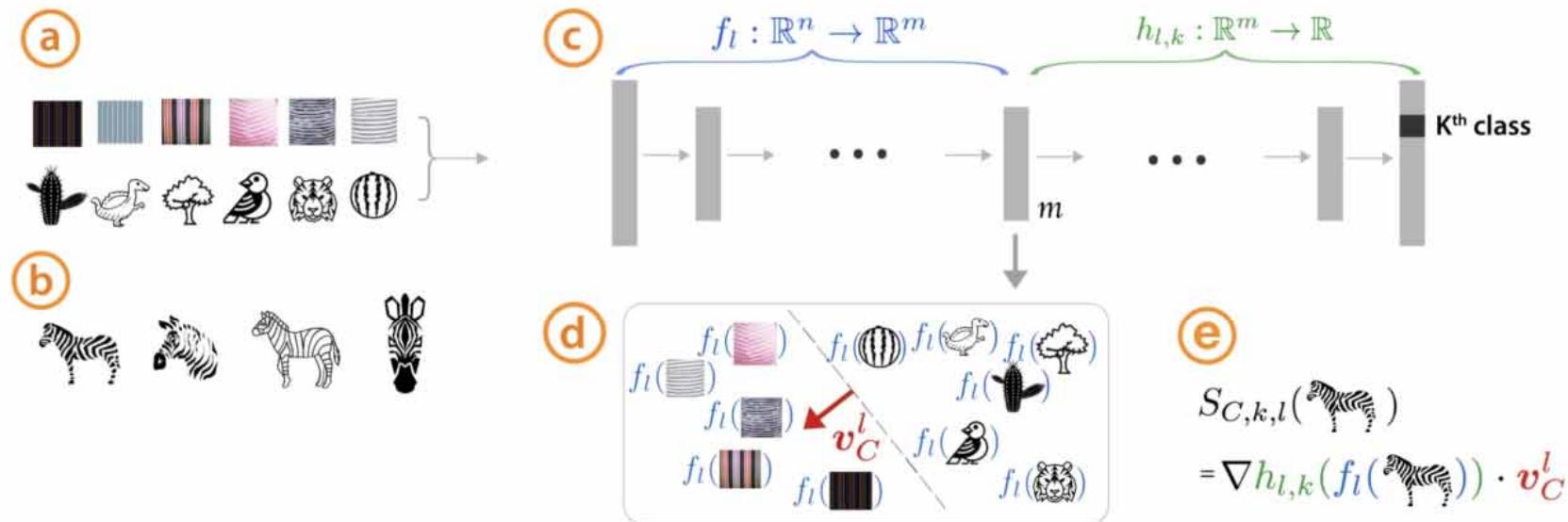
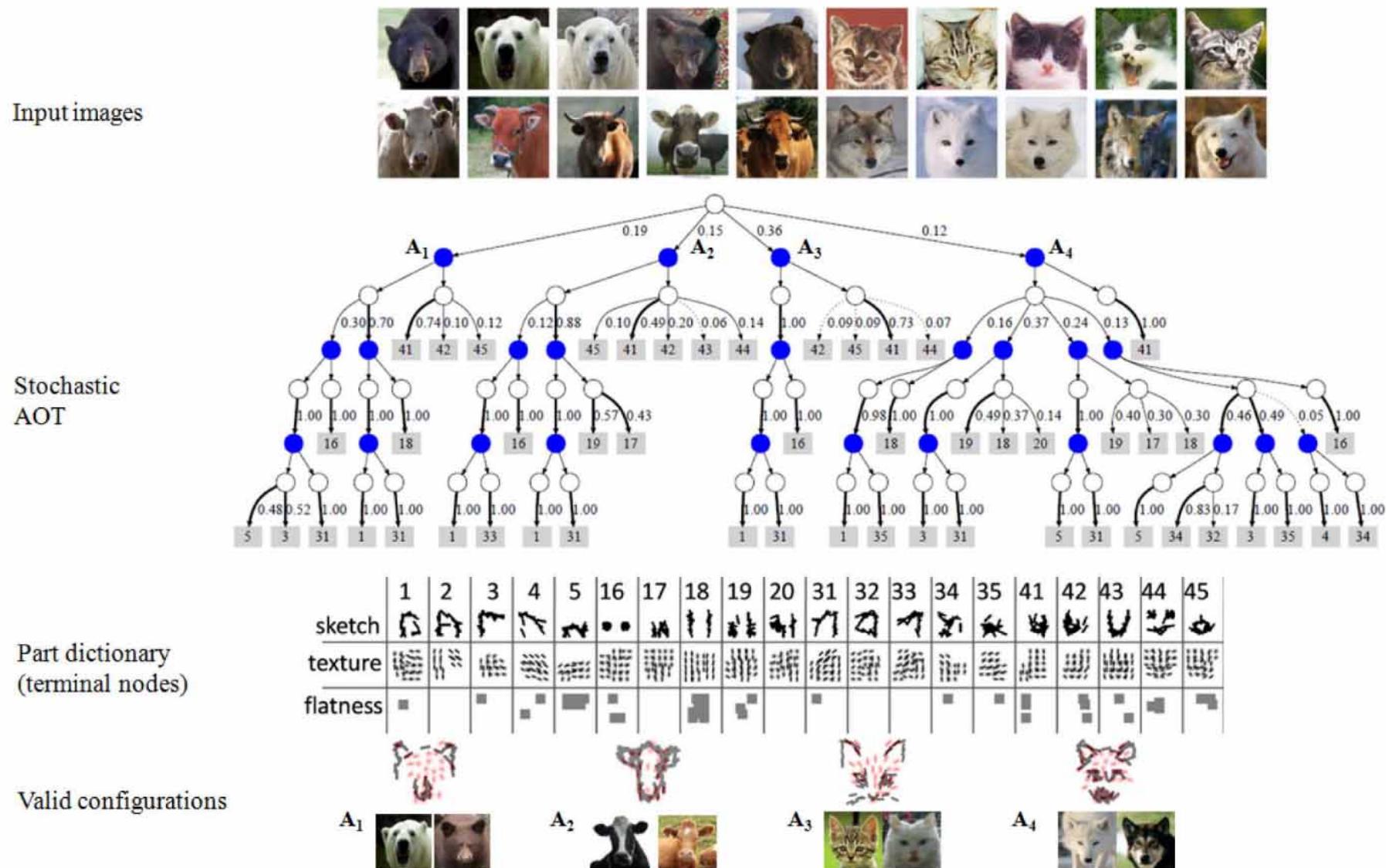


Figure 1. Testing with Concept Activation Vectors: Given a user-defined set of examples for a concept (e.g., ‘striped’), and random examples ④, labeled training-data examples for the studied class (zebras) ⑤, and a trained network ③, TCAV can quantify the model’s sensitivity to the concept for that class. CAVs are learned by training a linear classifier to distinguish between the activations produced by a concept’s examples and examples in any layer ④. The CAV is the vector orthogonal to the classification boundary (v_C^l , red arrow). For the class of interest (zebras), TCAV uses the directional derivative $S_{C,k,l}(x)$ to quantify conceptual sensitivity ⑥.

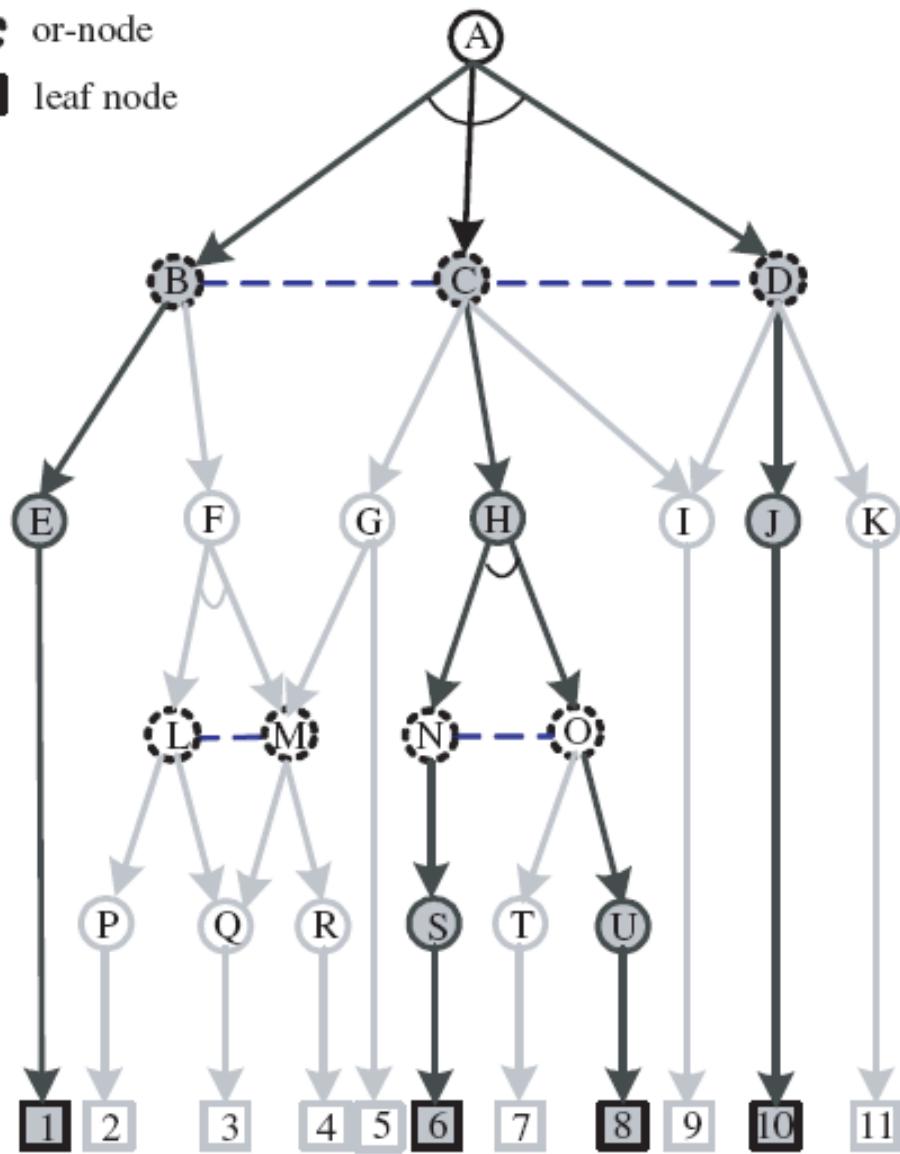
<https://github.com/tensorflow/tcav>

Been Kim, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler & Fernanda Viegas. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (TCAV). International Conference on Machine Learning (ICML), 2018. 2673-2682.



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.

- and-node
- or-node
- leaf node



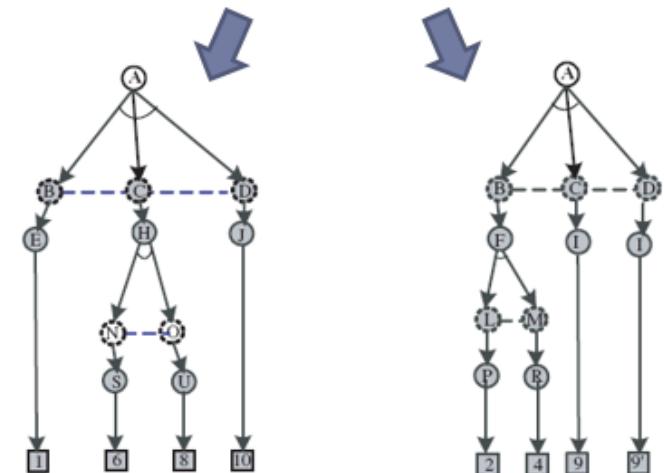
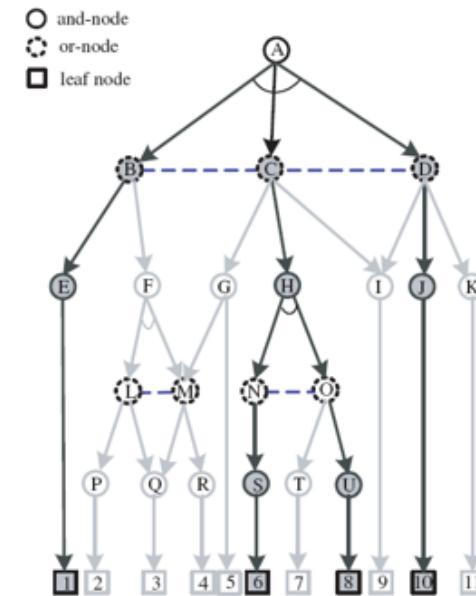
- Algorithm for this framework
 - Top-down/bottom-up computation
- Generalization of small sample
 - Use Monte Carlos simulation to synthesis more configurations
- Fill semantic gap

Images credit to Zhaoyin Jia (2009)

- ▶ Terminal (leaf) node: $T(pg)$
- ▶ And-Or node: $V^{or}(pg), V^{and}(pg)$
- ▶ Set of links: $E(pg)$
- ▶ Switch variable at Or-node: $w(t)$
- ▶ Attributes of primitives: $\alpha(t)$

$$p(pg; \Theta, R, \Delta) = \frac{1}{Z(\Theta)} \exp(-\xi(pg))$$

$$\begin{aligned} \xi(pg) = & \sum_{v \in V^{or}(pg)} \lambda_v(w(v)) + \sum_{v \in V^{and}(pg) \cup T(pg)} \lambda_t(\alpha(t)) \\ & + \sum_{(i,j) \in E(pg)} \lambda_{ij}(v_i, v_j, \gamma_{ij}, \rho_{ij}) \end{aligned}$$

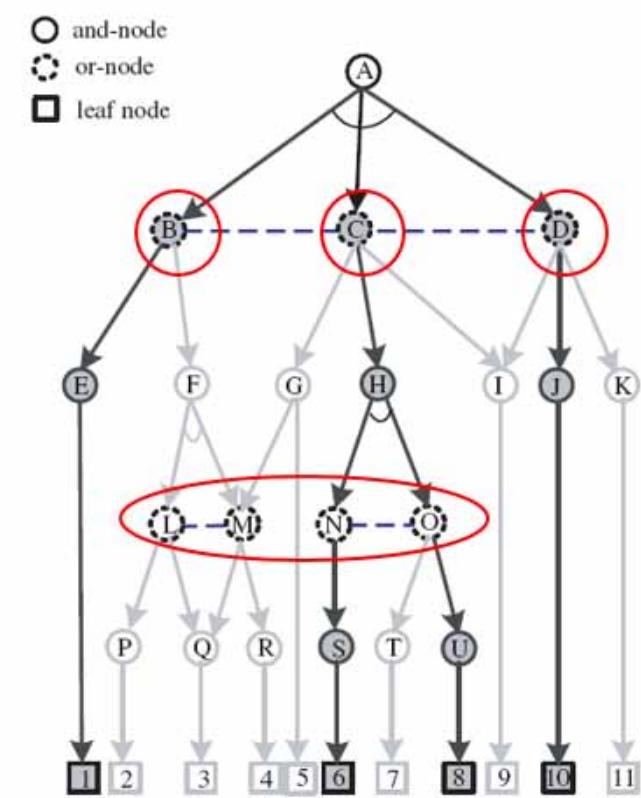


- ▶ Terminal (leaf) node: $T(pg)$
- ▶ And-Or node: $V^{or}(pg), V^{and}(pg)$
- ▶ Set of links: $E(pg)$
- ▶ Switch variable at Or-node: $w(t)$
- ▶ Attributes of primitives: $\alpha(t)$

$$p(pg; \Theta, R, \Delta) = \frac{1}{Z(\Theta)} \exp(-\xi(pg))$$

$$\begin{aligned} \xi(pg) &= \sum_{v \in V^{or}(pg)} \lambda_v(w(v)) + \sum_{v \in V^{and}(pg) \cup T(pg)} \lambda_t(\alpha(t)) \\ &+ \sum_{(i,j) \in E(pg)} \lambda_{ij}(v_i, v_j, \gamma_{ij}, \rho_{ij}) \end{aligned}$$

SCFG: weigh the frequency at the children of or-nodes

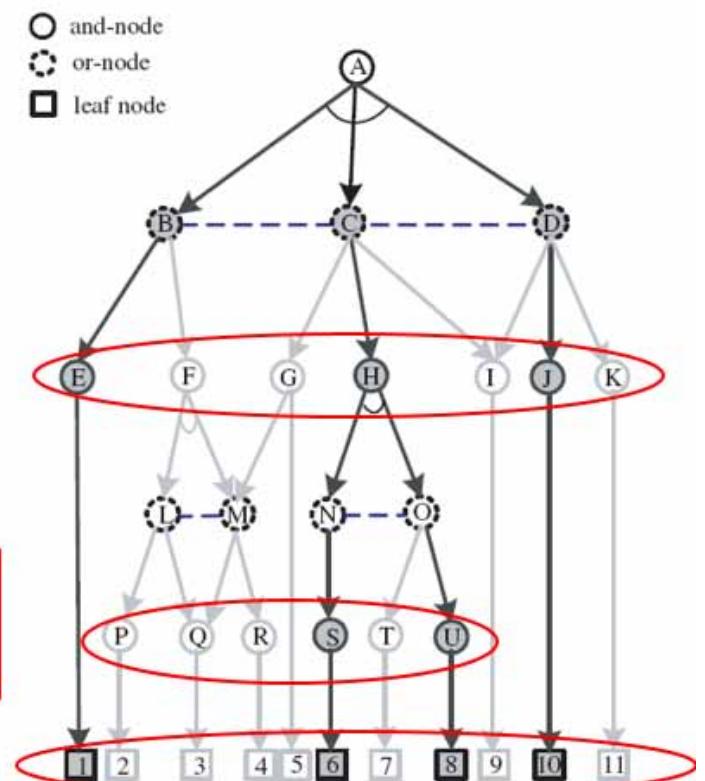


- ▶ Terminal (leaf) node: $T(pg)$
- ▶ And-Or node: $V^{or}(pg), V^{and}(pg)$
- ▶ Set of links: $E(pg)$
- ▶ Switch variable at Or-node: $w(t)$
- ▶ Attributes of primitives: $\alpha(t)$

$$p(pg; \Theta, R, \Delta) = \frac{1}{Z(\Theta)} \exp(-\xi(pg))$$

$$\begin{aligned} \xi(pg) = & \sum_{v \in V^{or}(pg)} \lambda_v(w(v)) + \boxed{\sum_{v \in V^{and}(pg) \cup T(pg)} \lambda_t(\alpha(t))} \\ & + \sum_{(i,j) \in E(pg)} \lambda_{ij}(v_i, v_j, \gamma_{ij}, \rho_{ij}) \end{aligned}$$

Weigh the local compatibility of primitives (geometric and appearance)

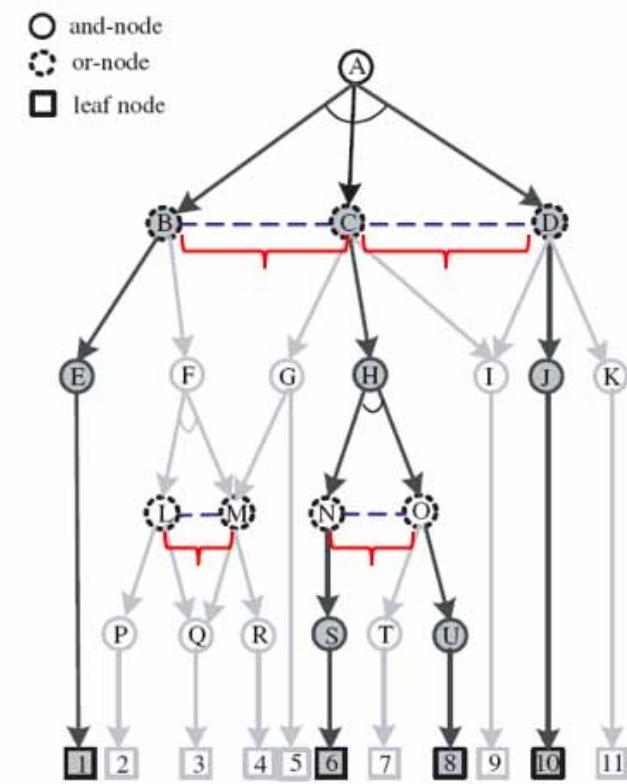


- ▶ Terminal (leaf) node: $T(pg)$
- ▶ And-Or node: $V^{or}(pg), V^{and}(pg)$
- ▶ Set of links: $E(pg)$
- ▶ Switch variable at Or-node: $w(t)$
- ▶ Attributes of primitives: $\alpha(t)$

$$p(pg; \Theta, R, \Delta) = \frac{1}{Z(\Theta)} \exp(-\xi(pg))$$

$$\begin{aligned} \xi(pg) = & \sum_{v \in V^{or}(pg)} \lambda_v(w(v)) + \sum_{v \in V^{and}(pg) \cup T(pg)} \lambda_t(\alpha(t)) \\ & + \sum_{(i,j) \in E(pg)} \lambda_{ij}(v_i, v_j, \gamma_{ij}, \rho_{ij}) \end{aligned}$$

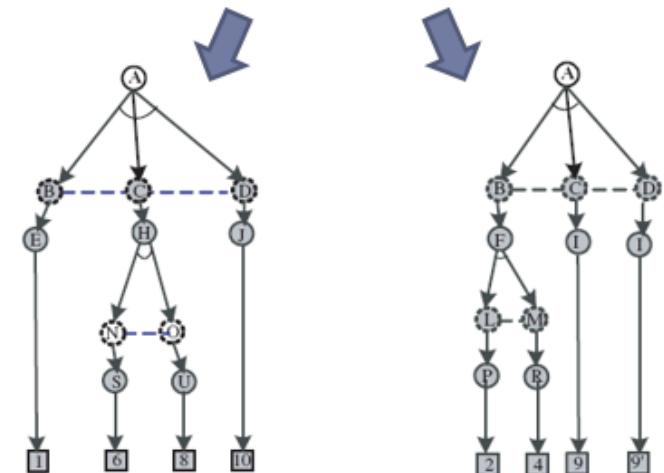
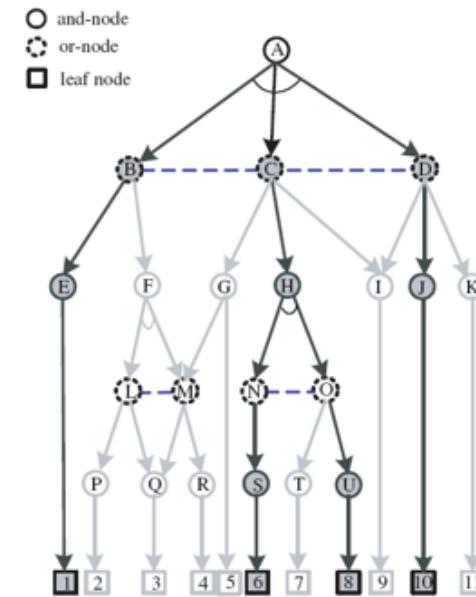
Spatial and appearance between primitives (parts or objects)

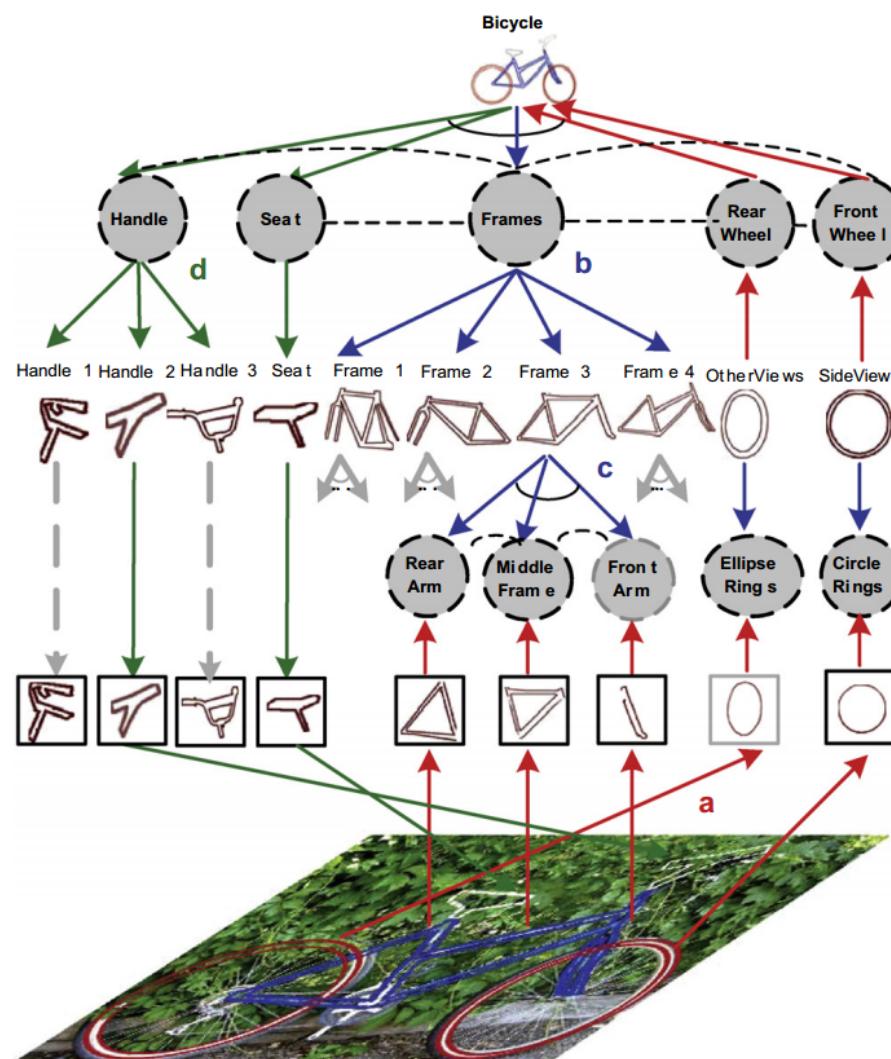


- ▶ Terminal (leaf) node: $T(pg)$
- ▶ And-Or node: $V^{or}(pg), V^{and}(pg)$
- ▶ Set of links: $E(pg)$
- ▶ Switch variable at Or-node: $w(t)$
- ▶ Attributes of primitives: $\alpha(t)$

$$p(pg; \Theta, R, \Delta) = \frac{1}{Z(\Theta)} \exp(-\xi(pg))$$

$$\begin{aligned} \xi(pg) = & \sum_{v \in V^{or}(pg)} \lambda_v(w(v)) + \sum_{v \in V^{and}(pg) \cup T(pg)} \lambda_t(\alpha(t)) \\ & + \sum_{(i,j) \in E(pg)} \lambda_{ij}(v_i, v_j, \gamma_{ij}, \rho_{ij}) \end{aligned}$$





Input: an input image I , and a set of constructed And-Or graphs of compositional object categories.

Output: a parsing graph pg_s of the scene that consists of the parsing graphs of detected objects.

- Repeat the following steps

- 1 Schedule the next node A to visit from the candidate parts.

- 2 Call Bottom-up(A) to update A 's **open** list.

- i Detect terminal instances of A from the image.

- ii Bind non-terminal instances of A from its children's **open** or **closed** lists

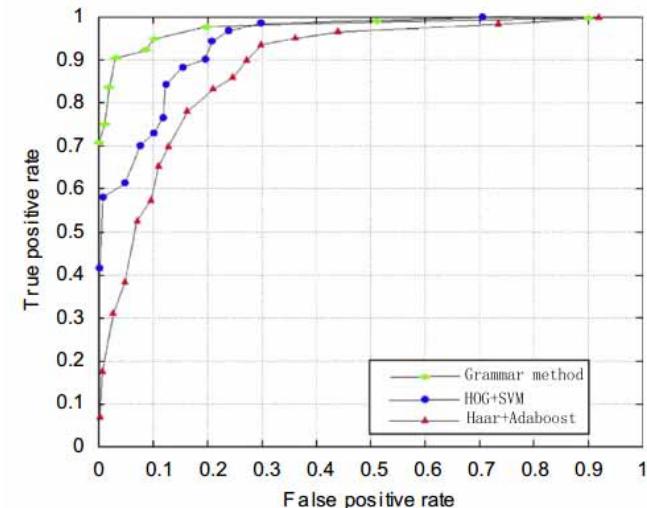
- 3 Call Top-down(A) to update A 's **open** or **closed** lists.

- i Accept hypotheses from A 's **open** list to its **closed** list.

- ii Remove (or disassemble) hypotheses from A 's **closed** list.

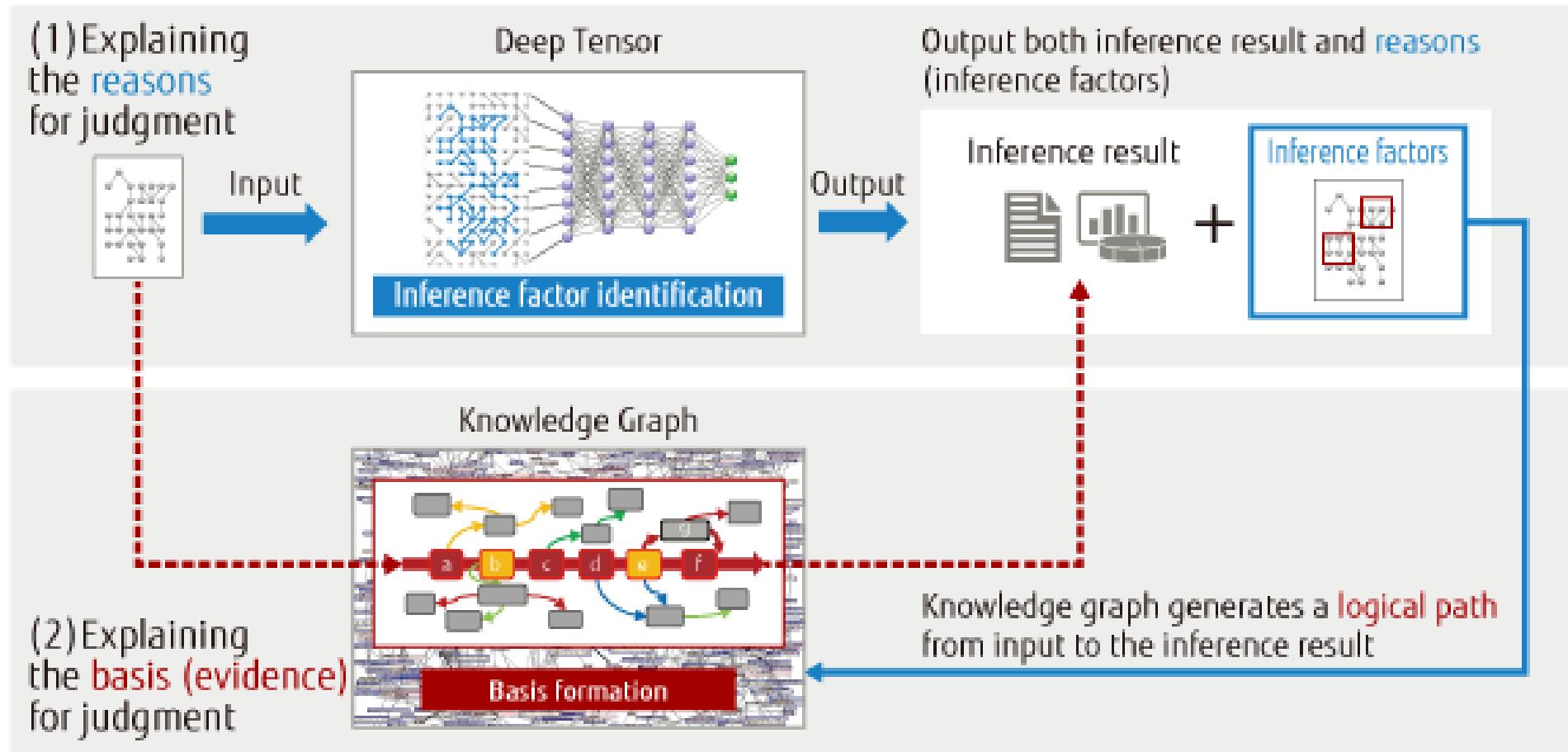
- iii Update the **open** lists for particles that overlap with node A .

- Until the particles in **open** list with weights higher than the empirical threshold are exhausted. Output all parsing graphs whose root nodes are reached.



Liang Lin, Tianfu Wu, Jake Porway & Zijian Xu 2009. A stochastic graph grammar for compositional object representation and recognition. Pattern Recognition, 42, (7), 1297-1307, doi:10.1016/j.patcog.2008.10.033.

Future Work



Explainable AI with Deep Tensor and Knowledge Graph

http://www.fujitsu.com/jp/Images/artificial-intelligence-en_tcm102-3781779.png

- What is a good explanation?
- (obviously if the other did understand it)
- Experiments needed!
- What is explainable/understandable/intelligible?
- When is it enough (Sättigungsgrad – you don't need more explanations – enough is enough)
- But how much is it ...



<https://www.newyorker.com/cartoon/a19697>

Noel C.F. Codella, Michael Hind, Karthikeyan Natesan Ramamurthy, Murray Campbell, Amit Dhurandhar, Kush R. Varshney, Dennis Wei & Aleksandra Mojsilović 2018. Teaching Meaningful Explanations. arXiv:1805.11648.

iv:1805.11648v1 [cs.AI] 29 May 2018

Teaching Meaningful Explanations

Noel C. F. Codella,* Michael Hind,* Karthikeyan Natesan Ramamurthy,*
Murray Campbell, Amit Dhurandhar, Kush R. Varshney, Dennis Wei,
Aleksandra Mojsilović

* These authors contributed equally.

IBM Research

Yorktown Heights, NY 10598

{nccodell,hindm,knatesa,mcam,adhuran,krvarshn,dwei,aleksand}@us.ibm.com

Abstract

The adoption of machine learning in high-stakes applications such as healthcare and law has lagged in part because predictions are not accompanied by explanations comprehensible to the domain user, who often holds ultimate responsibility for decisions and outcomes. In this paper, we propose an approach to generate such explanations in which training data is augmented to include, in addition to features and labels, explanations elicited from domain users. A joint model is then learned to produce both labels and explanations from the input features. This simple idea ensures that explanations are tailored to the complexity expectations and domain knowledge of the consumer. Evaluation spans multiple modeling techniques on a simple game dataset, an image dataset, and a chemical odor dataset, showing that our approach is generalizable across domains and algorithms. Results demonstrate that meaningful explanations can be reliably taught to machine learning algorithms, and in some cases, improve modeling accuracy.

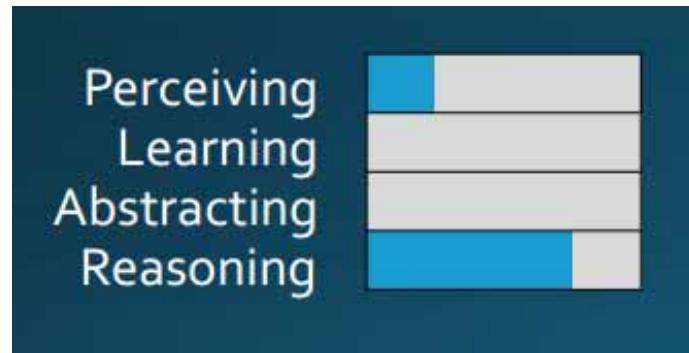
1 Introduction

New regulations call for automated decision making systems to provide “meaningful information” on the logic used to reach conclusions [1–4]. Selbst and Powles interpret the concept of “meaningful information” as information that should be understandable to the audience (potentially individuals



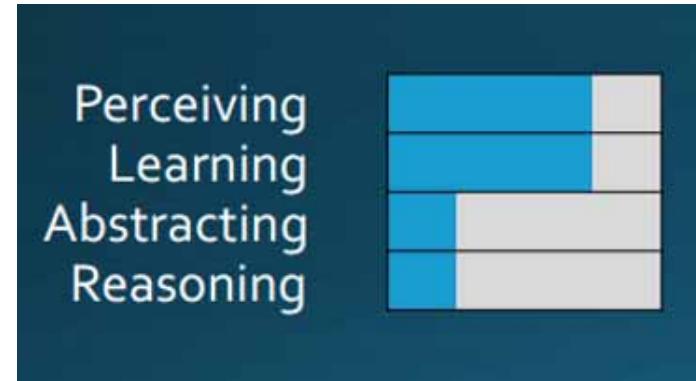
This image is in the Public Domain

- Computational approaches can find in R^n what no human is able to see
- However, still there are many hard problems where a human expert in R^2 can understand the **context** and bring in experience, expertise, knowledge, intuition, ...
- Black box approaches can not explain **WHY** a decision has been made ...



- Engineers create a set of logical rules to represent knowledge (Rule based Expert Systems)
- Advantage: works well in narrowly defined problems of well-defined domains
- Disadvantage: No adaptive learning behaviour and poor handling of $p(x)$

Image credit to John Launchbury



- Engineers create learning models for specific tasks and train them with “big data” (e.g. Deep Learning)
- Advantage: works well for standard classification tasks and has prediction capabilities
- Disadvantage: No contextual capabilities and minimal reasoning abilities

Image credit to John Launchbury



- A contextual model can perceive, learn and understand and abstract and reason
- Advantage: can use transfer learning for adaptation on unknown unknowns
- Disadvantage: Superintelligence ...

Image credit to John Launchbury

- Myth 1a: Superintelligence by 2100 is inevitable!
- Myth 1b: Superintelligence by 2100 is impossible!
- **Fact: We simply don't know it!**
- Myth 2: Robots are our main concern
 - Fact: Cyberthreats are the main concern:
it needs no body – only an Internet connection**
- Myth 3: AI can never control us humans
 - Fact: Intelligence is an enabler for control:
We control tigers by being smarter ...**



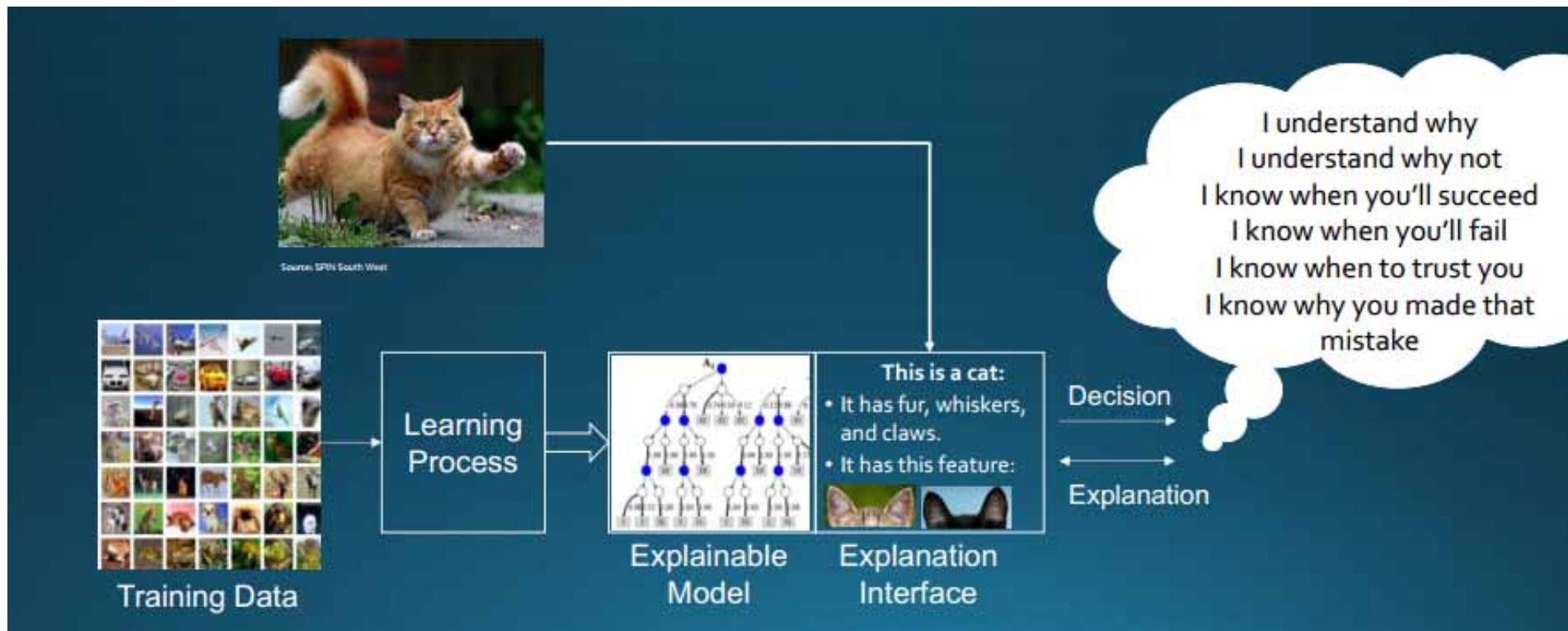


Image credit to John Launchbury



Thank you!