

MAKE Decisions Medical Information Science for Decision Support



Assoc. Prof. Dr. Andreas HOLZINGER (Medical University Graz)



<https://hci-kdd.org/mini-course-make-decisions-practice>

Day 2 > Part 5 > 20.09.2018

Methods of explainable AI concise overview (details in LV 706.315 TU Graz)

Day 1 - Hot Ideas

01 Information Sciences meets Life Sciences

02 Data, Information and Knowledge

03 Decision Making and Decision Support

04 DSS: from Expert Systems to explainable AI

Day 2 - Cool Practice

05 Methods of Explainable-AI

Groupwork: Planning of a 500 bed Hospital - Bringing AI into the workflows

Plenary: Presenting of the developed concepts

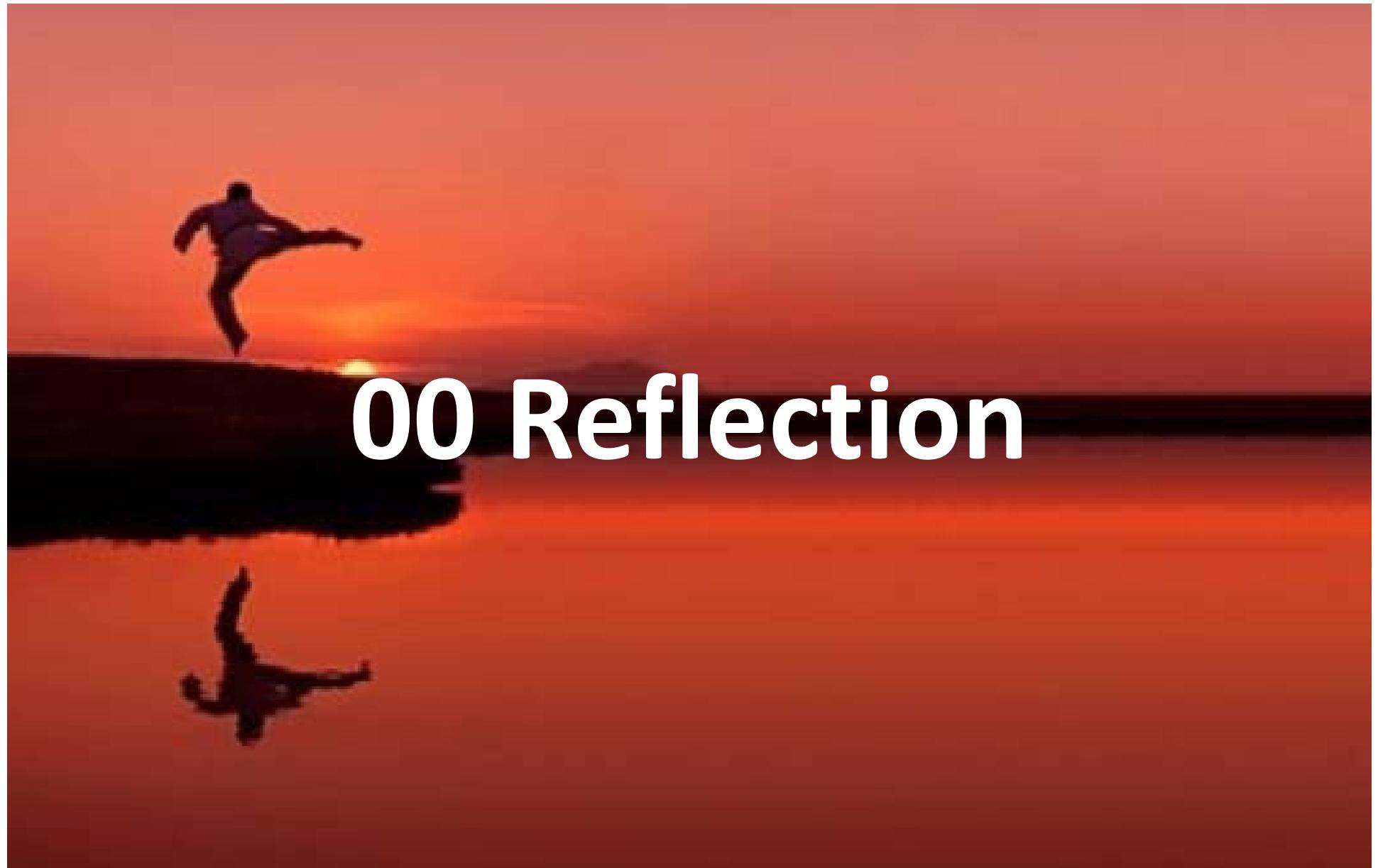
- Explainable Artificial intelligence
- Re-trace > Understand > Explain
- Transparency > Trust > Acceptance
- Fairness > Transparency > Accountability
- Methods of Explainable AI
- Transparent Machine Learning
- Intelligible > Reliability > Trustworthy
- Interpretability > Explicable > Accountable
- Interpretable Deep Learning

- **Ante-hoc Explainability (AHE)** := such models are interpretable by design, e.g. glass-box approaches; typical examples include linear regression, decision trees/lists, random forests, Naive Bayes and fuzzy inference systems; or GAMs, Stochastic AOGs, and deep symbolic networks; they have a long tradition and can be designed from expert knowledge or from data and are useful as framework for the interaction between human knowledge and hidden knowledge in the data.
- **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;
- **Explainability** := motivated by the opaqueness of so called “black-box” approaches it is the ability to provide an explanation on why a machine decision has been reached (e.g. why is it a cat what the deep network recognized). Finding an appropriate explanation is difficult, because this needs understanding the context and providing a description of causality and consequences of a given fact. (German: Erklärbarkeit; siehe auch: Verstehbarkeit, Nachvollziehbarkeit, Zurückverfolgbarkeit, Transparenz)
- **Explanation** := set of statements to describe a given set of facts to clarify causality, context and consequences thereof and is a core topic of knowledge discovery involving “why” questionss (“Why is this a cat?”). (German: Erklärung, Begründung)
- **Explanatory power** := is the ability of a set hypothesis to effectively explain the subject matter it pertains to (opposite: explanatory impotence).
- **European General Data Protection Regulation (EU GDPR)** := Regulation EU 2016/679 – see the EUR-Lex 32016R0679 , will make black-box approaches difficult to use, because they often are not able to explain why a decision has been made (see explainable AI).

- **Interactive Machine Learning (iML)** := machine learning algorithms which can interact with – partly human – agents and can optimize its learning behaviour through this interaction. 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.
- **Inverse Probability** := an older term for the probability distribution of an unobserved variable, and was described by De Morgan 1837, in reference to Laplace's (1774) method of probability.
- **Post-hoc Explainability (PHE)** := such models are designed for interpreting black-box models and provide local explanations for a specific decision and re-enact on request, typical examples include LIME, BETA, LRP, or Local Gradient Explanation Vectors, prediction decomposition or simply feature selection.
- **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)

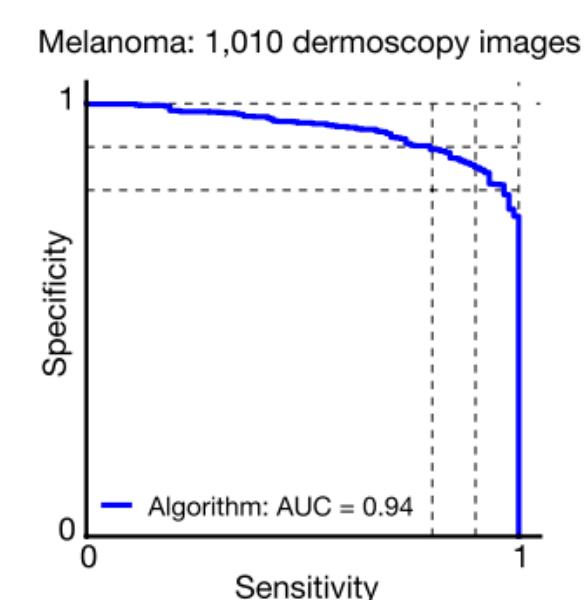
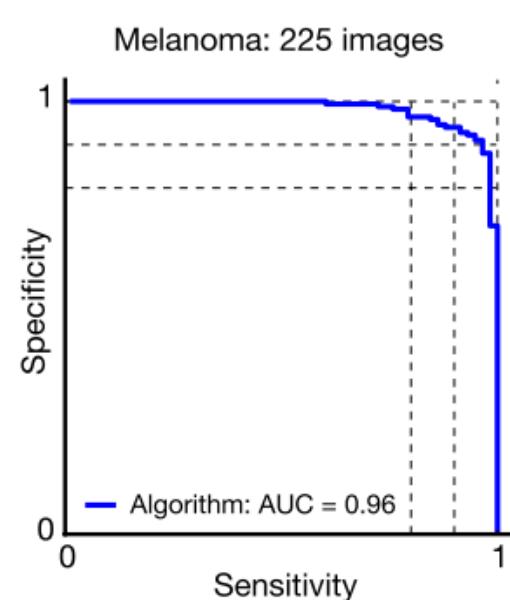
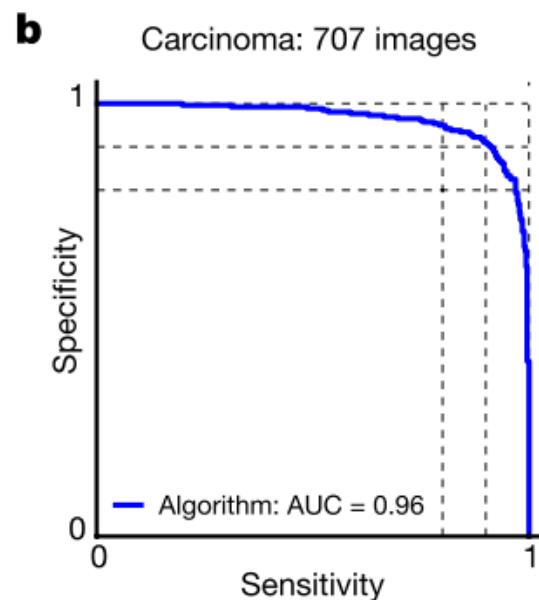
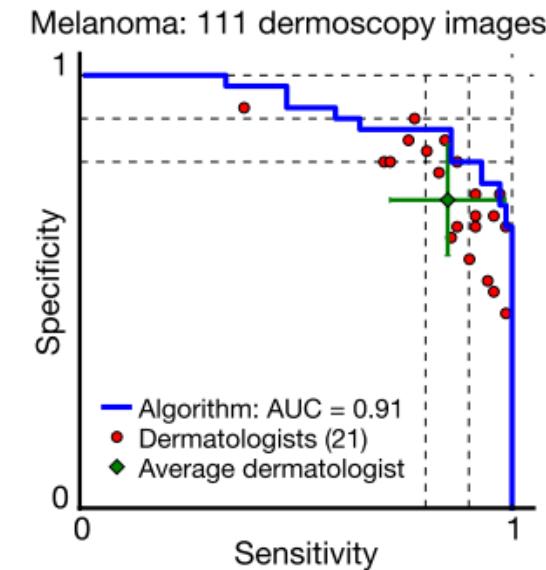
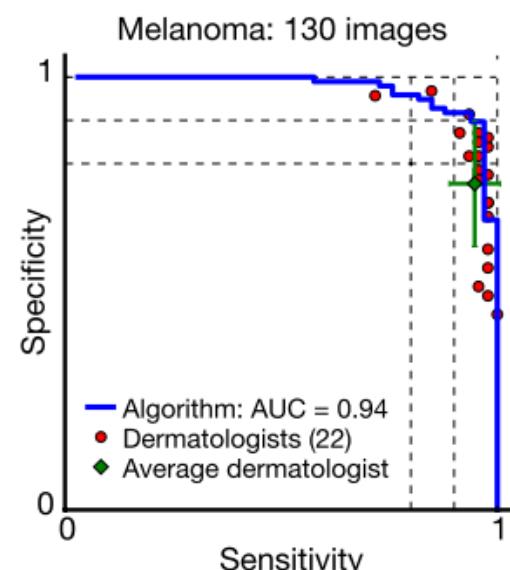
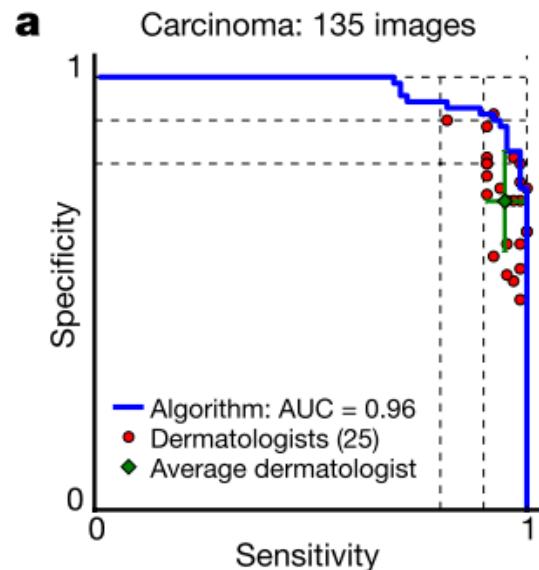
- ... appreciate why this thematic is of increasing importance for the international community
- ... understand some basic concepts of what explainability, transparency, understandable is
- ... are aware of some methods of ex-AI
- ... comprehend the complexity of context and the grand goal of context adaptive systems
- ... see how important causal reasoning is
- ... get a feeling on how important this field is for the medical domain for building trust, acceptance and reliability

- 00 Reflection – follow-up from last lecture
- 01 Towards explainable AI
- 02 Causal Reasoning
- 03 Tradeoff: Explainability vs. Accuracy
- 04 (Some) Current State of the Art Methods
- 05 Stochastic And-Or-Graphs



00 Reflection

How do you explain this ...



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-computer interface and interaction; and (g) system integration into existing clinical legacy and proprietary environments, e.g. the enterprise hospital information system; to mention only a few.

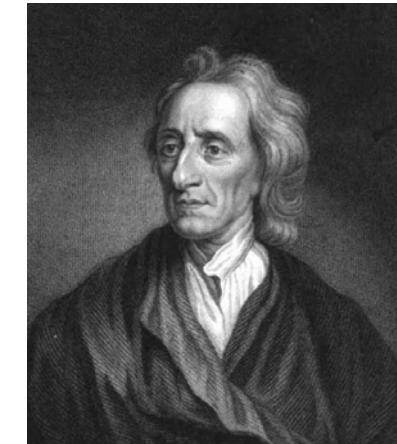
01 Towards Explainable AI

- 1) Wrong decisions can be very costly ...
- 2) Insights from new perspectives ...
- 3) Raising trust, acceptance, responsibility, ...
- 4) Compliance to legislation ...



- Understanding ... psychological concept correlated with the ability to make inferences.
- Mechanistic understanding
 - fMRI – which area is stimulated if one moves a finger
- Functional understanding $f : \mathbb{R}^d \rightarrow \mathbb{R}$
 - what is the relation between input and output
- Causal understanding
 - what is cause and effect (-> Judea Pearl)

John Locke 1841. An essay concerning human understanding, Teaside, Thomas TEGG.



“It is not a human move ... I've never seen such a move ...”



人機世紀之戰

Google

AlphaGo

vs.

李世石

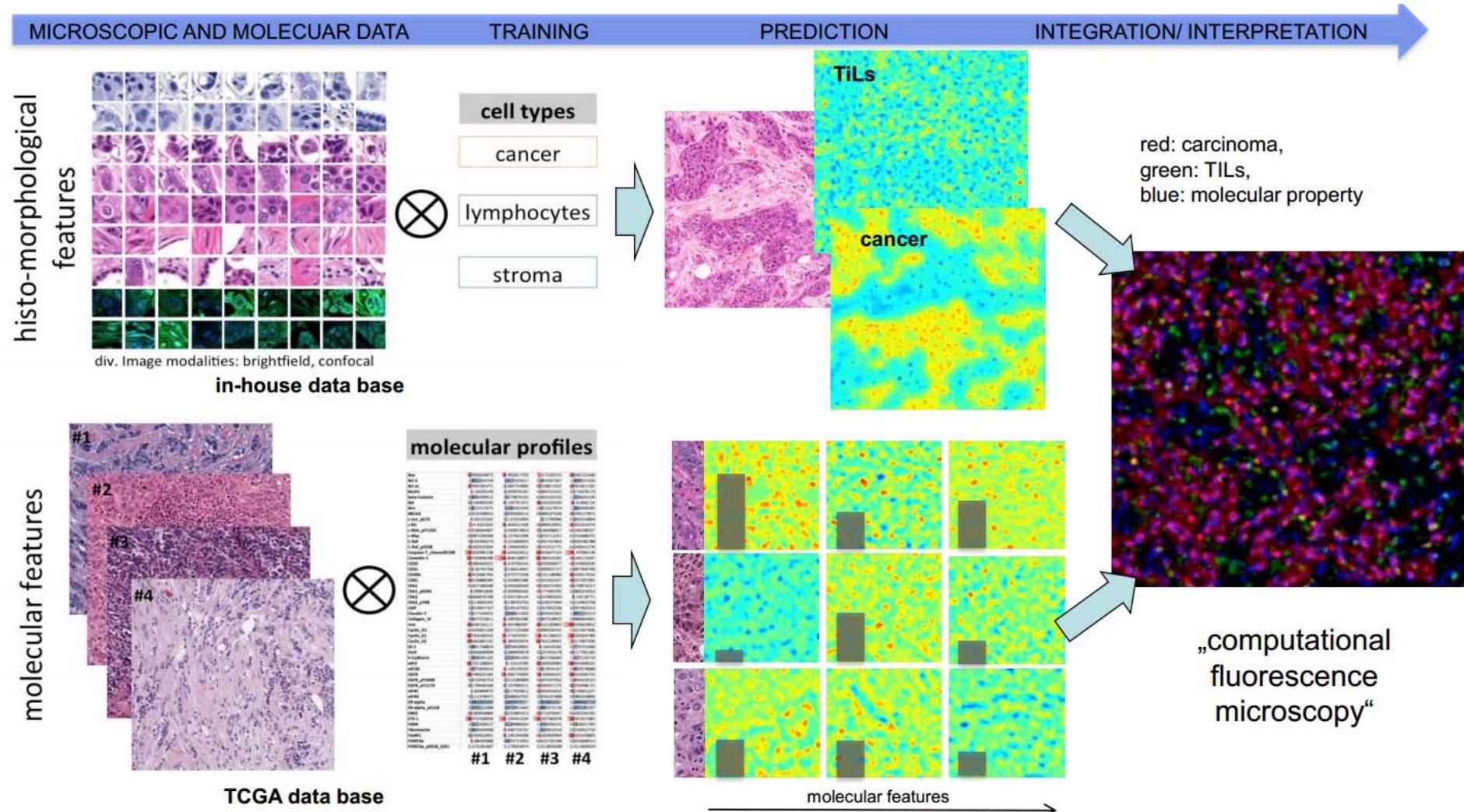
3/10 (四) 12:00 Round 2

LIVE

必POTV

全程中文直播

Integration and Interpretation



Alexander Binder, Michael Bockmayr, Miriam Hägele, Stephan Wienert, Daniel Heim, Katharina Hellweg, Albrecht Stenzinger, Laura Parlow, Jan Budczies & Benjamin Goeppert 2018. Towards computational fluorescence microscopy: Machine learning-based integrated prediction of morphological and molecular tumor profiles. arXiv:1805.11178vl.

Context understanding is difficult



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)

02 Causal Reasoning

Example: Discovery of causal relationships from data ...

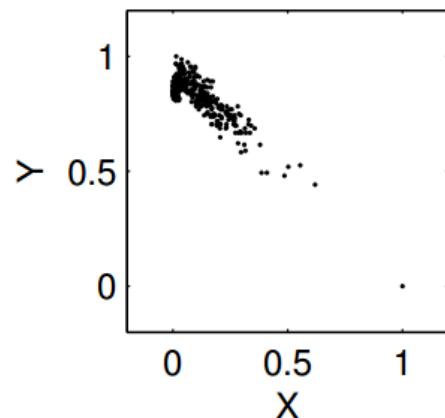
Hans Holbein d.J., 1533,
The Ambassadors,
London: National Gallery

Lopez-Paz, D., Muandet, K., Schölkopf, B. & Tolstikhin, I. 2015. Towards a learning theory of cause-effect inference. Proceedings of the 32nd International Conference on Machine Learning, JMLR, Lille, France.

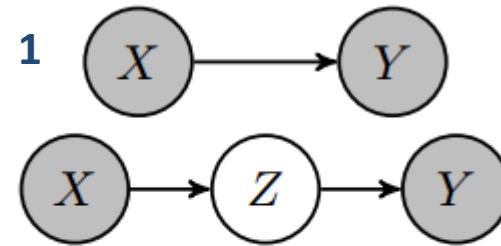


<https://www.youtube.com/watch?v=9KiVNIUMmCc>

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

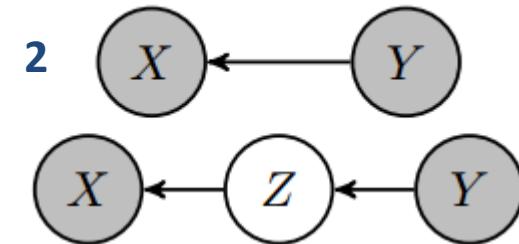


Joris M. Mooij, Jonas Peters, Dominik Janzing, Jakob Zscheischler & Bernhard Schölkopf
 2016. Distinguishing cause from effect using observational data: methods and benchmarks. The Journal of Machine Learning Research, 17, (1), 1103-1204.



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

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



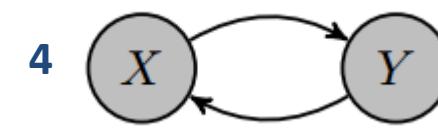
$$\mathbb{P}_Y = \mathbb{P}_{Y | \text{do}(x)} \neq \mathbb{P}_{Y | x}$$

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



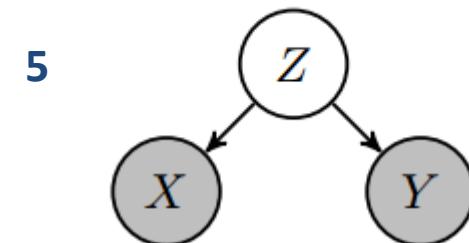
$$\mathbb{P}_Y = \mathbb{P}_{Y | \text{do}(x)} = \mathbb{P}_{Y | x}$$

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



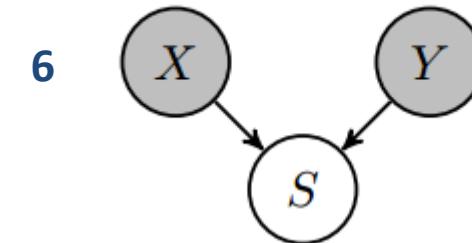
$$\mathbb{P}_Y \neq \mathbb{P}_{Y | \text{do}(x)} \neq \mathbb{P}_{Y | x}$$

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



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

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



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

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

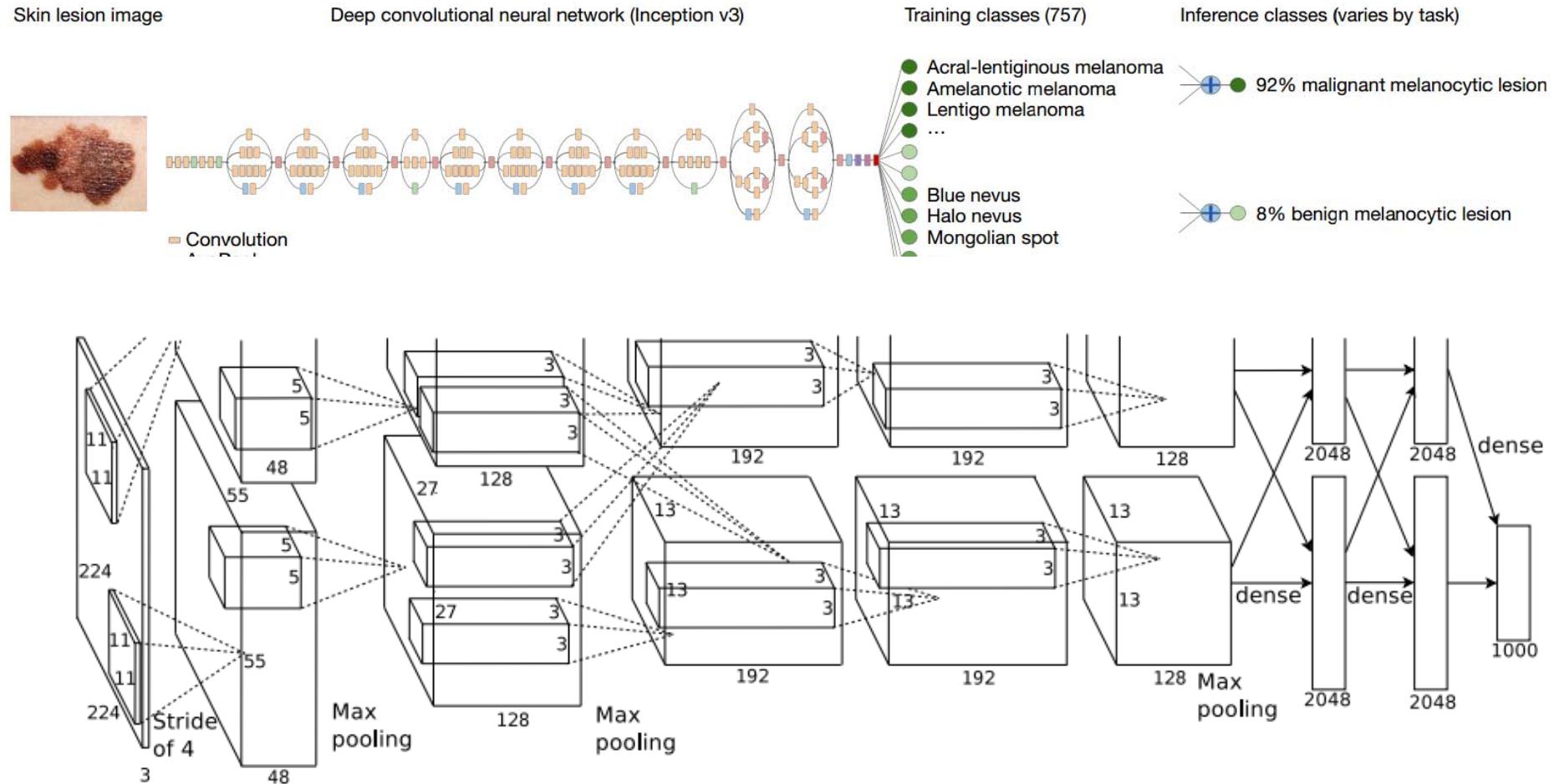
- “How do humans generalize from few examples?”
 - Learning relevant representations
 - Disentangling the explanatory factors
 - Finding the shared underlying explanatory factors, in particular between $P(x)$ and $P(Y|X)$, with a causal link between $Y \rightarrow X$

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.

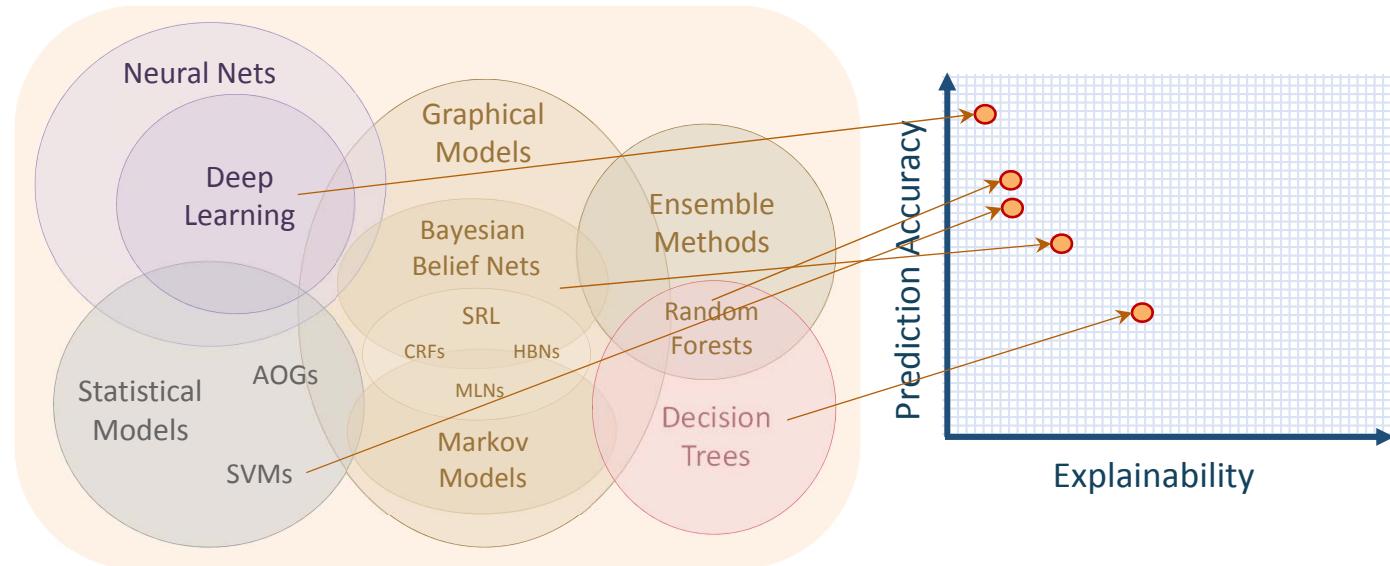
Tenenbaum, J. B., Kemp, C., Griffiths, T. L. & Goodman, N. D. 2011. How to grow a mind: Statistics, structure, and abstraction. *Science*, 331, (6022), 1279-1285, doi:10.1126/science.1192788.

03 Tradeoff: Explainability vs. Accuracy

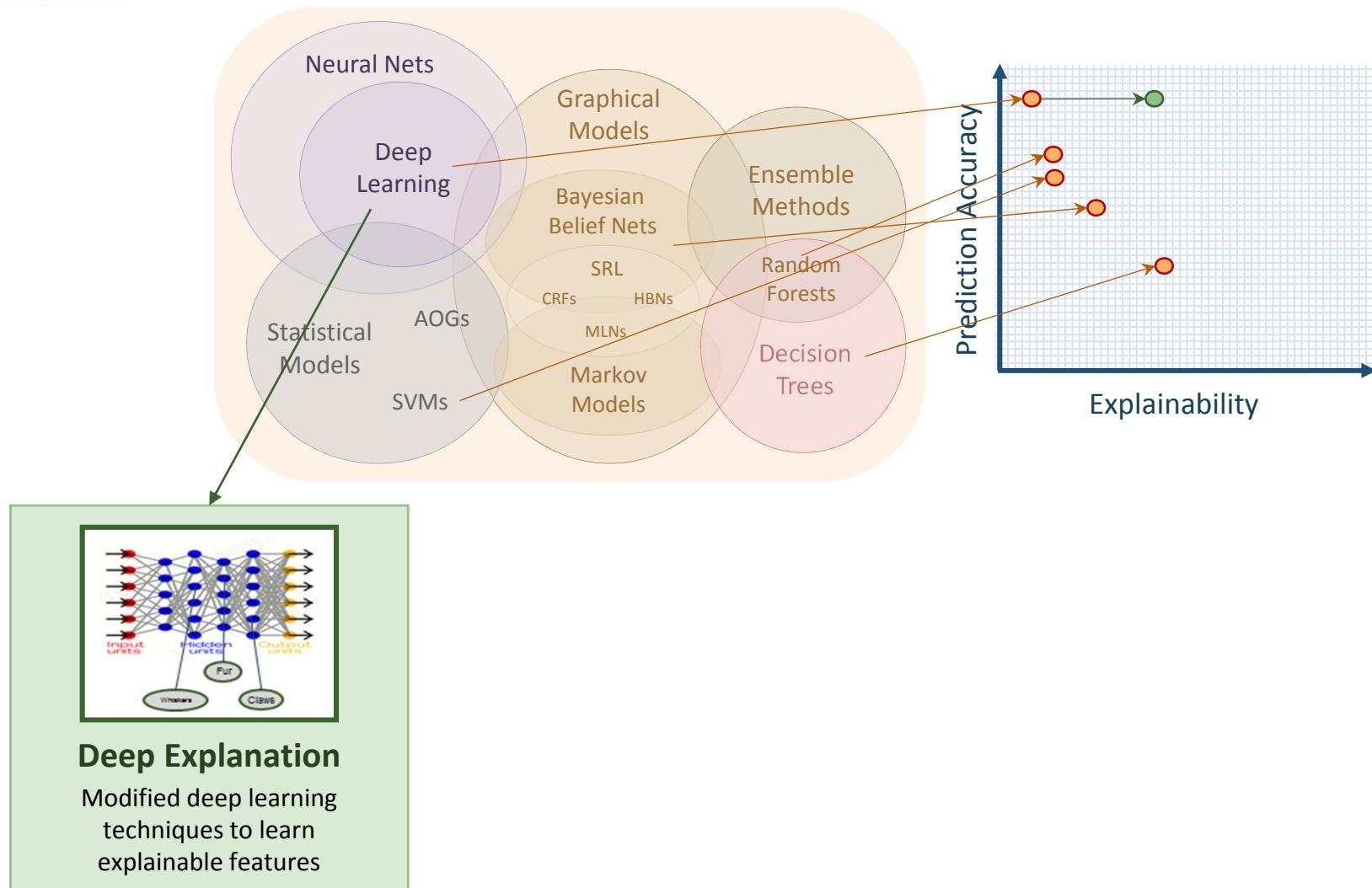
Deep Convolutional Neural Network Pipeline



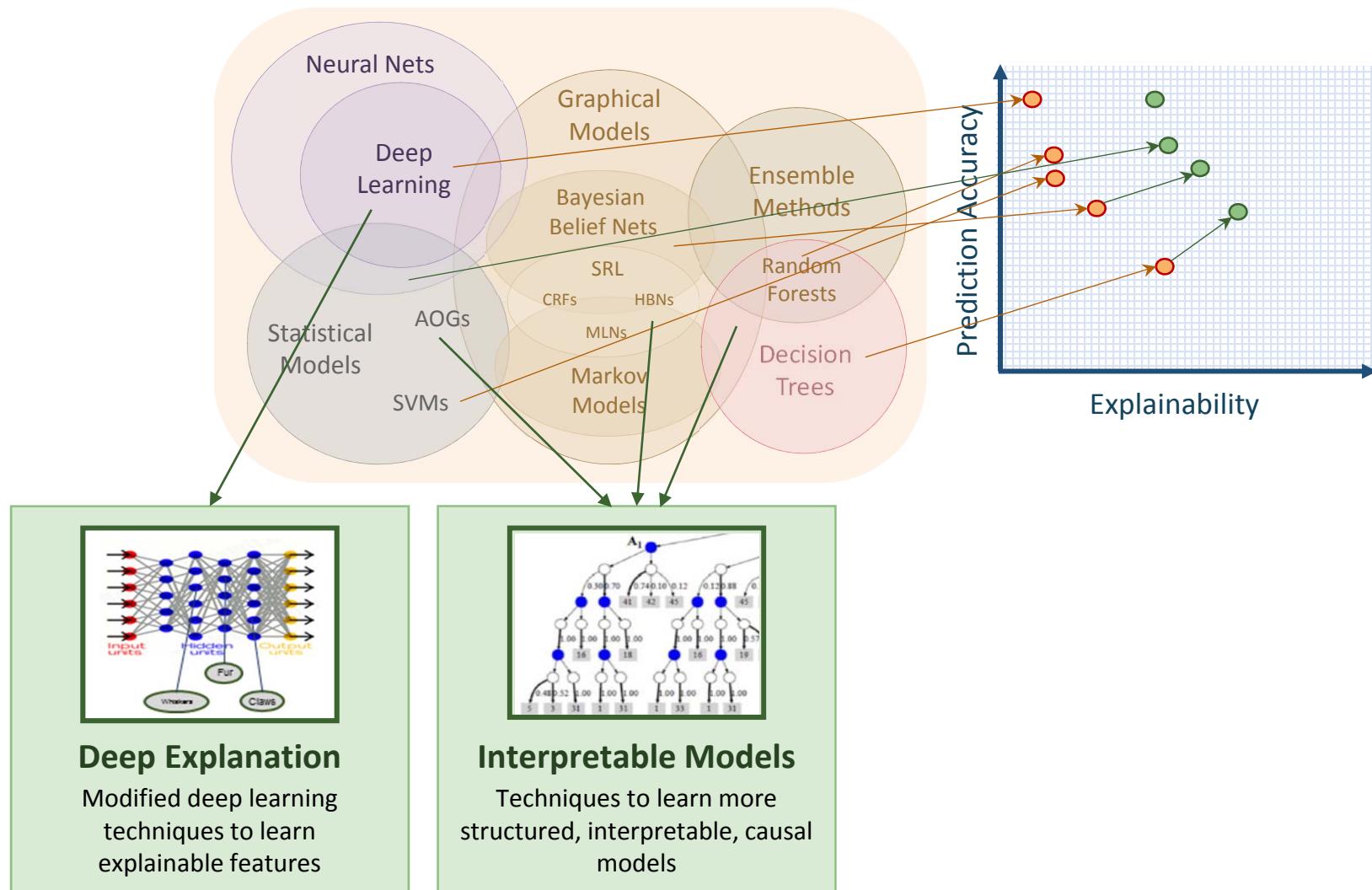
Krizhevsky, A., Sutskever, I. & Hinton, G. E. Imagenet classification with deep convolutional neural networks. In: Pereira, F., Burges, C. J. C., Bottou, L. & Weinberger, K. Q., eds. Advances in neural information processing systems (NIPS 2012), 2012 Lake Tahoe. 1097-1105.



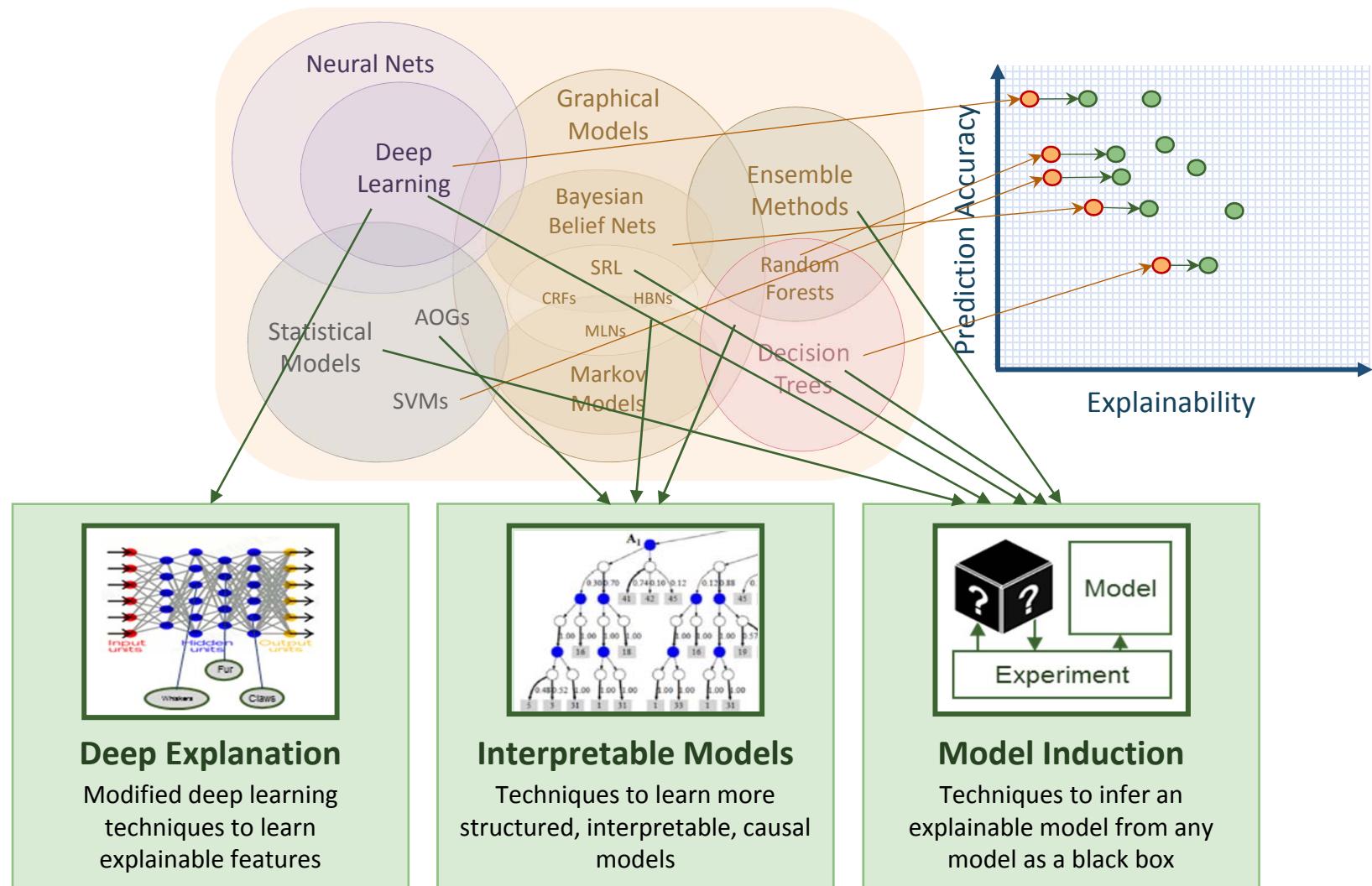
David Gunning 2016. Explainable artificial intelligence (XAI): Technical Report Defense Advanced Research Projects Agency DARPA-BAA-16-53, Arlington, USA, DARPA.



David Gunning 2016. Explainable artificial intelligence (XAI): Technical Report Defense Advanced Research Projects Agency DARPA-BAA-16-53, Arlington, USA, DARPA.



David Gunning 2016. Explainable artificial intelligence (XAI): Technical Report Defense Advanced Research Projects Agency DARPA-BAA-16-53, Arlington, USA, DARPA.



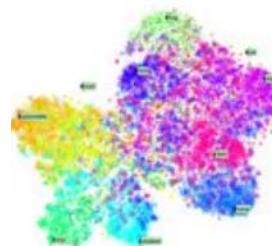
David Gunning 2016. Explainable artificial intelligence (XAI): Technical Report Defense Advanced Research Projects Agency DARPA-BAA-16-53, Arlington, USA, DARPA.

Post-hoc: Select a model and develop a technique to make it transparent



$$f(x) = \text{DeepNet}(x)$$

Different dimensions of “interpretability”



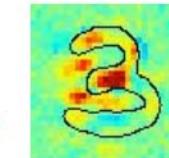
data

“Which dimensions of the data are most relevant for the task.”

Ante-hoc: Select a model that is already transparent and optimize it

$$f(x) = \sum_{i=1}^d g_i(x_i)$$

contribution of i th variable



prediction

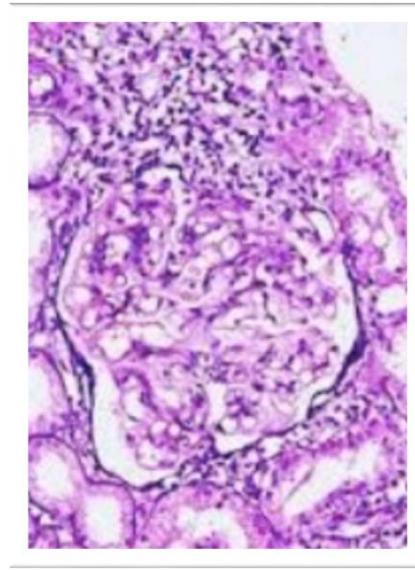
“Explain why a certain pattern x has been classified in a certain way $f(x)$.”

model

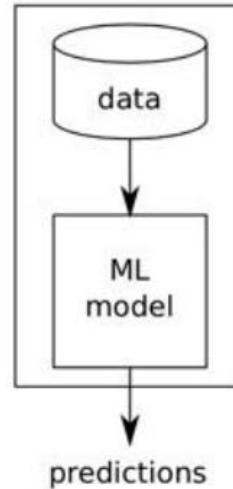
“What would a pattern belonging to a certain category typically look like according to the model.”



Montavon, G., Samek, W. & Müller, K.-R. 2017. Methods for interpreting and understanding deep neural networks. arXiv:1706.07979.

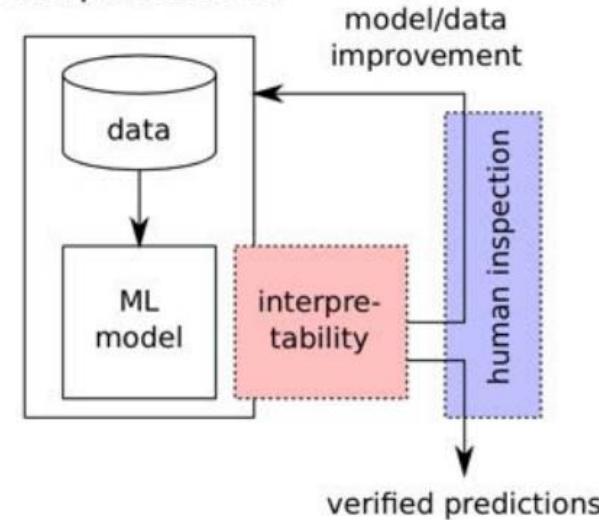


Standard ML



Generalization error

Interpretable ML



Generalization error + human experience

Image credit to: Samek, Montavon & Müller Tutorial at ICASSEP 2017

- Wrong decisions can be costly and dangerous!
- Verify that classifier works as expected
- Improve classifier continuously
- Human learning inspired by machine learning

- Interpretability as a novel kind for supporting teaching, learning and knowledge discovery,
- Particularly in abstract fields (informatics)
- Compliance to European Law “the right of explanation”
- Check for bias in machine learning results
- Fostering trust, acceptance, making clear the reliability ... “can you trust your results”

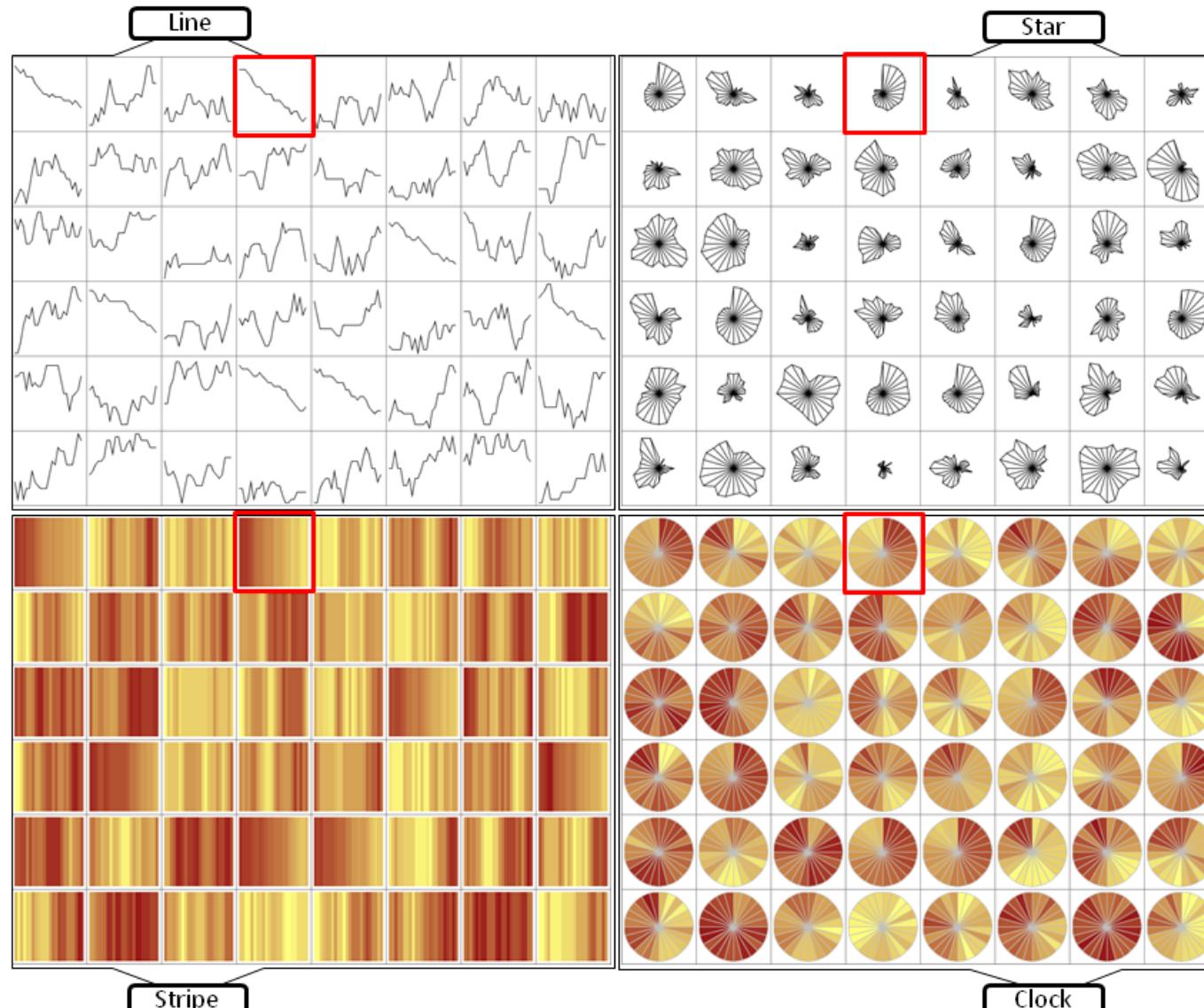
Katharina Holzinger, Klaus Mak, Peter Kieseberg & Andreas Holzinger 2018. Can we trust Machine Learning Results? Artificial Intelligence in Safety-Critical decision Support. ERCIM News, 112, (1), 42-43.

What is interpretable for humans?

What is interpretable for humans?

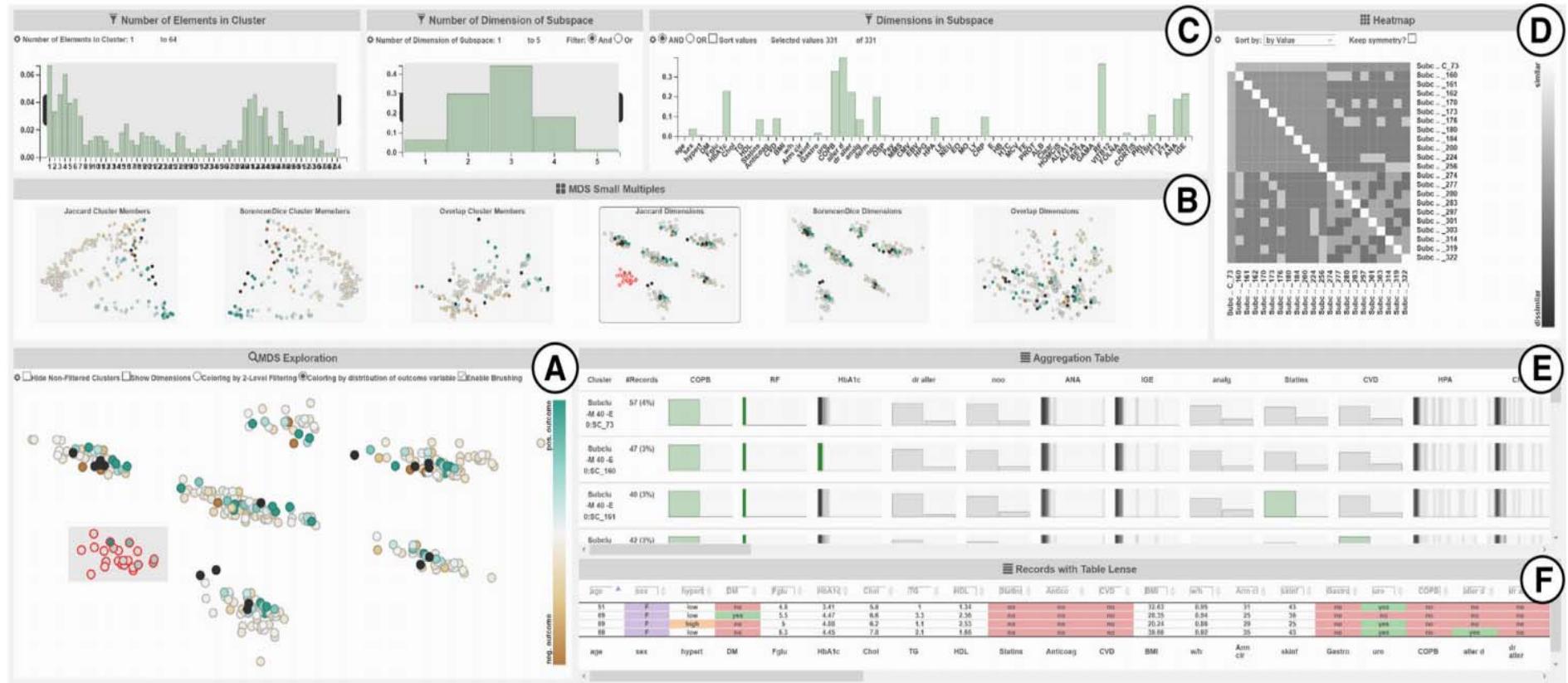
1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	0.953	0.894	0.620	0.699	0.629	0.546	0.540	1.000	0.526	1.000	0.522	0.483	0.471	1.000	0.522	0.576	0.658						
1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	0.722	0.638	1.000	0.785	0.743	0.792	0.801	0.875	0.712	1.000	0.444	0.947	0.431	1.000	0.793	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	0.576	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	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	0.505	0.744	1.000	0.485	
1.000	1.000	1.000	1.000	1.000	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	0.370	0.376	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	0.350	0.350	0.350	0.350	0.354	0.318	0.462	0.491	0.426	0.510	0.578	0.538	
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	0.370	0.370	0.370	0.370	0.370	0.370	0.370	0.370	0.370	0.370	0.370	0.370	0.370	0.370	0.370
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	0.620	0.620	0.620	0.620	0.620	0.620	0.620	0.620	0.620	0.620	0.620	0.620	0.620	0.620	0.620	0.620	
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	0.733	0.731	0.729	0.729	0.729	0.729	0.729	0.729	0.729	0.729	0.729	0.729	0.729	0.729	0.729	0.729	0.729	0.729	
1.000	1.000	0.590	0.509	0.486	0.445	0.411	0.372	0.569	0.675	0.732	0.741	0.741	0.741	0.741	0.741	0.741	0.741	0.741	0.741	0.741	0.741	0.741	0.741	0.741	0.741	0.741	0.741	0.741	0.741		
1.000	0.924	0.554	0.517	0.450	0.416	0.449	0.373	0.585	0.7	0.727	0.736	0.747	0.747	0.747	0.747	0.747	0.747	0.747	0.747	0.747	0.747	0.747	0.747	0.747	0.747	0.747	0.747	0.747	0.747		
1.000	1.000	0.557	0.517	0.457	0.396	0.390	0.390	0.635	0.658	0.707	0.719	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	0.380	0.380	0.380		
1.000	1.000	0.556	0.4	0.42	0.36	0.52	0.623	0.635	0.670	0.711	0.748	0.771	0.775	0.772	0.724	0.598	0.440	0.434	0.378	0.354	0.414	0.307	0.282	0.278	0.402	0.402	0.402	0.402	0.402	0.402	0.402
0.763	1.000	0.611	0.515	0.490	0.408	0.383	0.484	0.590	0.646	0.687	0.718	0.724	0.748	0.717	0.659	0.530	0.500	0.560	0.594	0.483	0.499	0.472	0.273	0.234	0.279	0.306	0.306	0.306	0.306	0.306	
1.000	1.000	0.750	0.600	0.544	0.328	0.490	0.550	0.623	0.593	0.515	0.521	0.616	0.616	0.620	0.598	0.502	0.431	0.338	0.279	0.295	0.330	0.446	0.446	0.446	0.446	0.446	0.446	0.446	0.446	0.446	
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	0.200	0.200	0.200	0.200
0.929	0.672	0.503	0.654	0.388	0.335	0.306	0.475	0.416	0.475	0.46	0.4	0.74	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	0.263	0.263	0.263	0.263		
1.000	0.758	0.639	0.726	0.931	0.330	0.299	0.398	0.54	0.5	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	0.252	0.252	0.252		
1.000	0.790	0.907	0.701	0.897	0.382	0.296	0.358	0.63	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	0.254	0.254	0.254	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	0.237	0.237	0.237	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	0.243	0.243	0.243	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	0.256	0.256	0.256	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	0.237	0.237	0.237	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	0.273	0.273	0.273	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	0.228	0.228	0.228	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	0.252	0.252	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	0.248	0.248	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	0.269	0.269	0.269	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	0.623	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	0.529	0.529	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	0.546	0.546	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	0.546	0.546	0.546		

What is understandable, interpretable, intelligible?



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

Explainable AI is a huge challenge for visualization



Michael Hund, Dominic Boehm, Werner Sturm, Michael Sedlmair, Tobias Schreck, Torsten Ullrich, Daniel A. Keim, Ljiljana Majnaric & Andreas Holzinger 2016. Visual analytics for concept exploration in subspaces of patient groups: Making sense of complex datasets with the Doctor-in-the-loop. *Brain Informatics*, 3, (4), 233-247, doi:10.1007/s40708-016-0043-5.

04 Current State of the Art

Example for an Explanation Interface

The screenshot shows the 'Message Predictor' application window. At the top, there are several controls: 'Move message to folder...', 'Only show predictions that just changed' (OFF), 'Search Stanley', and 'Clear'. Below this is a 'Folders' section with a list of messages in the 'Unknown' folder (1,180 messages). A message from Harold Zazula is selected, showing it was predicted to be about 'Hockey' with 99% confidence. The message content discusses an octopus on the ice after a game. To the right, a detailed explanation is provided, highlighting words like 'baseball', 'hockey', 'stanley', and 'tiger' as important, with a ratio of 2.3 times more likely for Hockey than Baseball. Below this, a bar chart shows the importance of various words, with 'baseball' and 'hockey' being the most prominent. A legend indicates that blue bars represent 'Hockey' and green bars represent 'Baseball'.

A Unknown (1,180 messages)
corresponding to 8/8 correct predictions

B

C

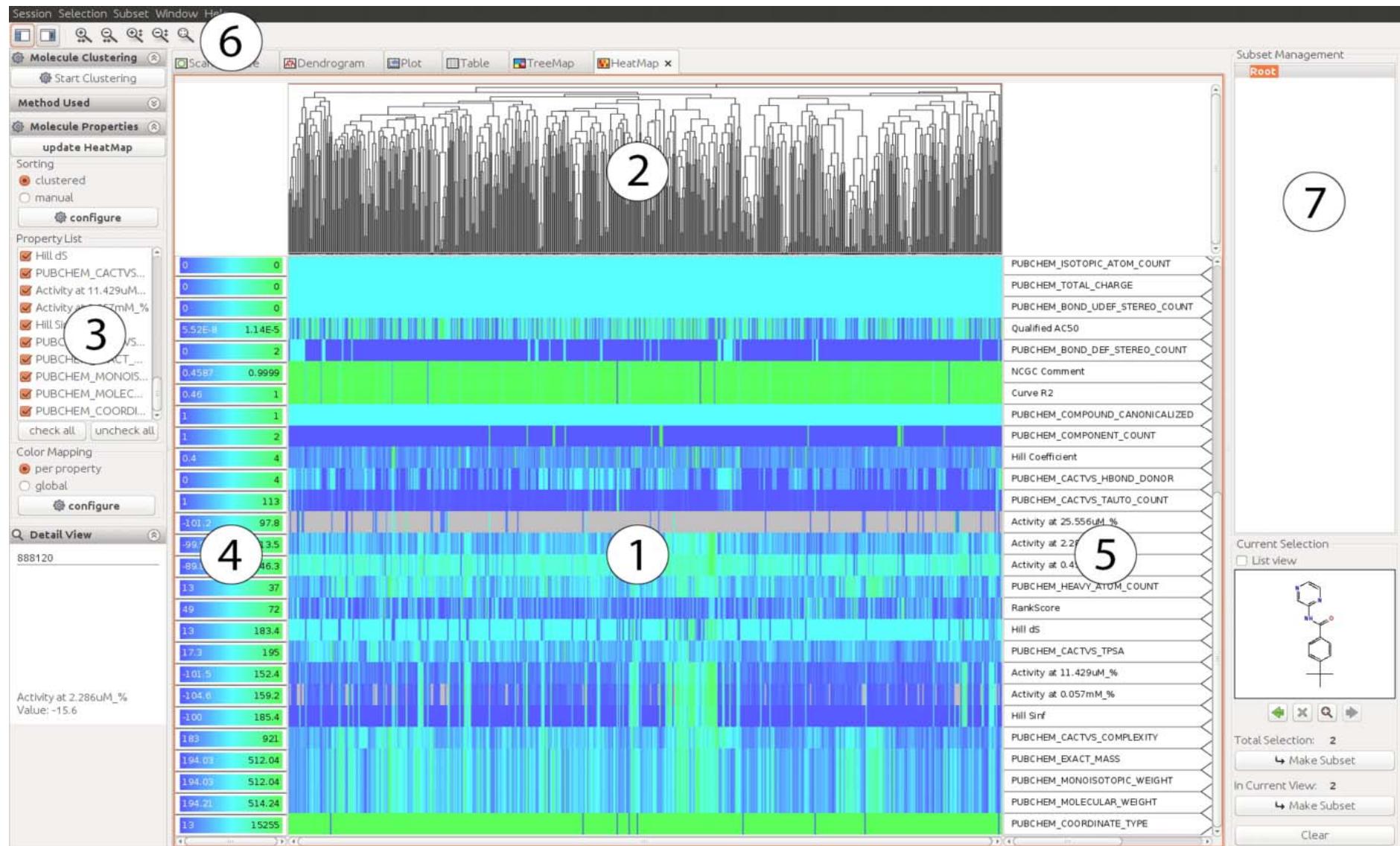
D

E

F

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.

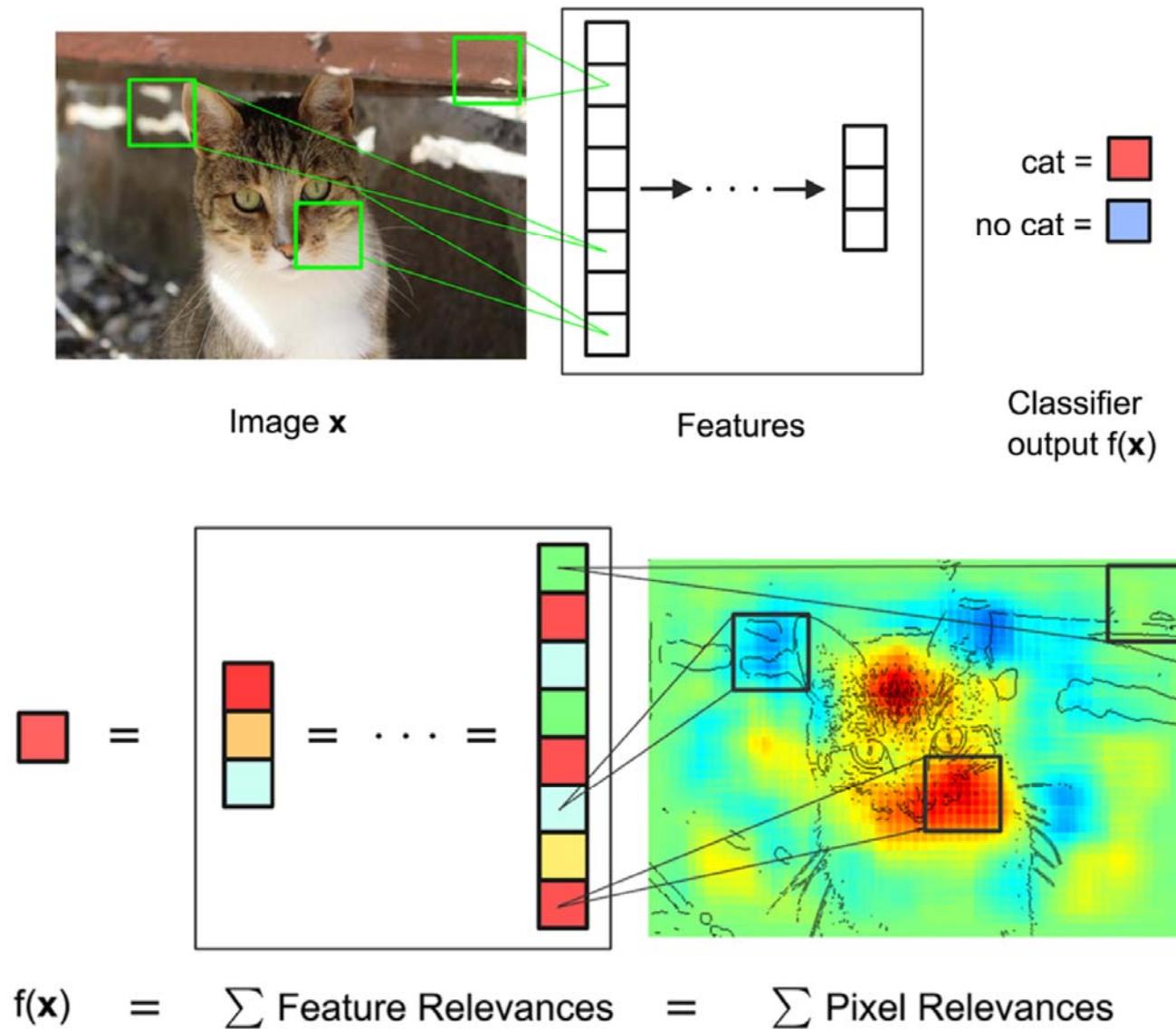
Example for an Explanation Interface - open work 😊



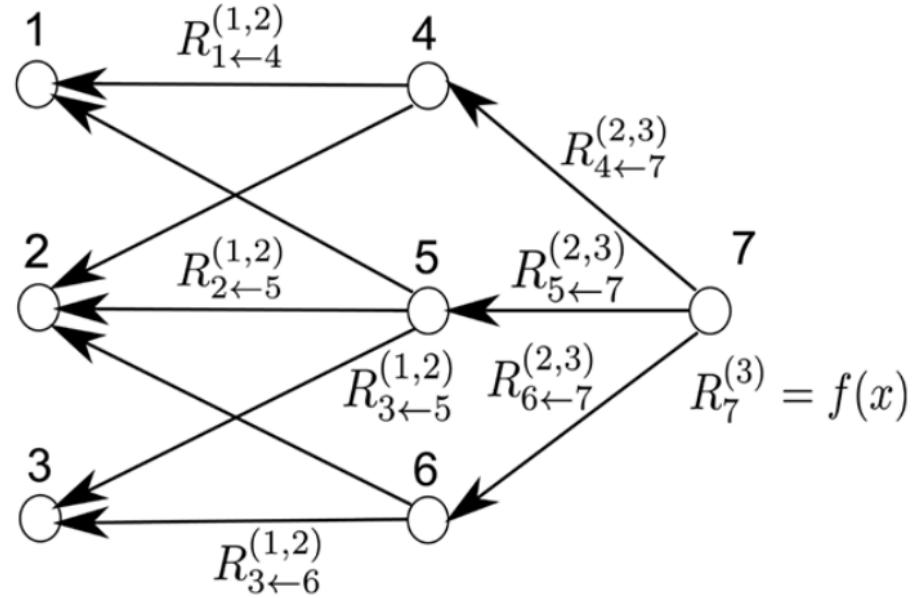
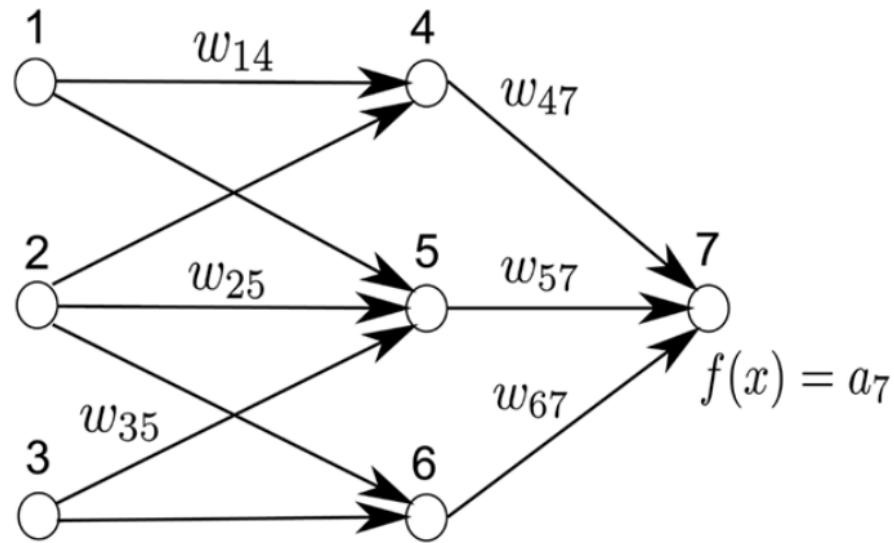
Werner Sturm, Till Schaefer, Tobias Schreck, Andreas Holzinger & Torsten Ullrich. Extending the Scaffold Hunter Visualization Toolkit with Interactive Heatmaps In: Borgo, Rita & Turkay, Cagatay, eds. EG UK Computer Graphics & Visual Computing CGVC 2015, 2015 University College London (UCL). Euro Graphics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

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.



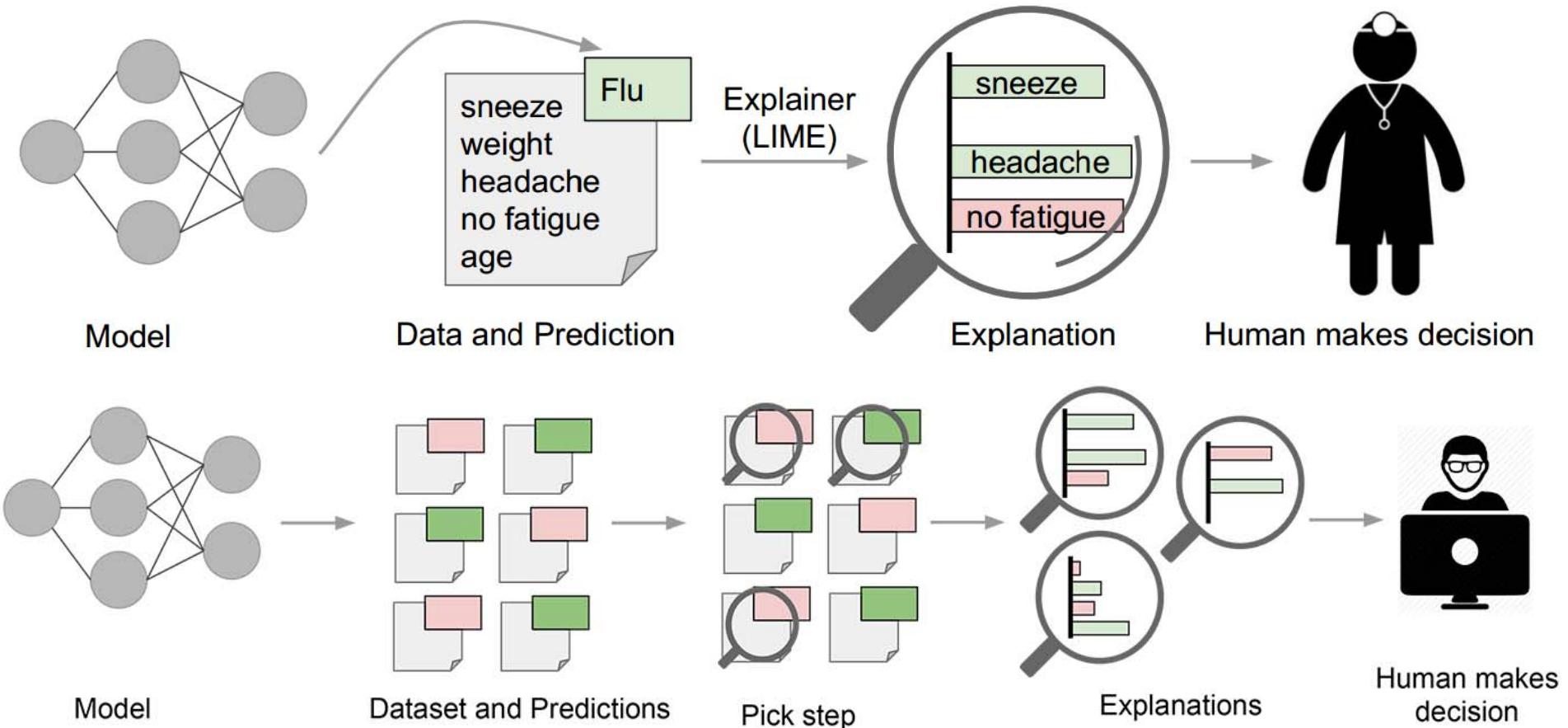
A NN-classifier during prediction time



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

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.

LIME – Local Interpretable Model Agnostic Explanations



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.

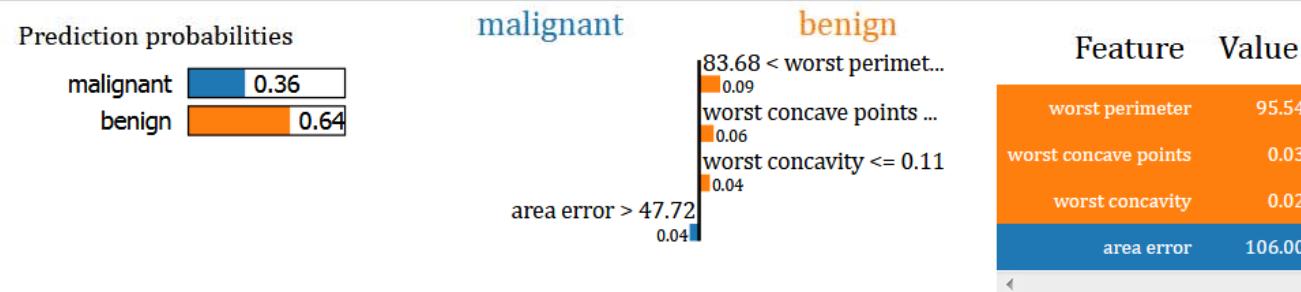
Example LIME – Model Agnostic Explanation

```
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])

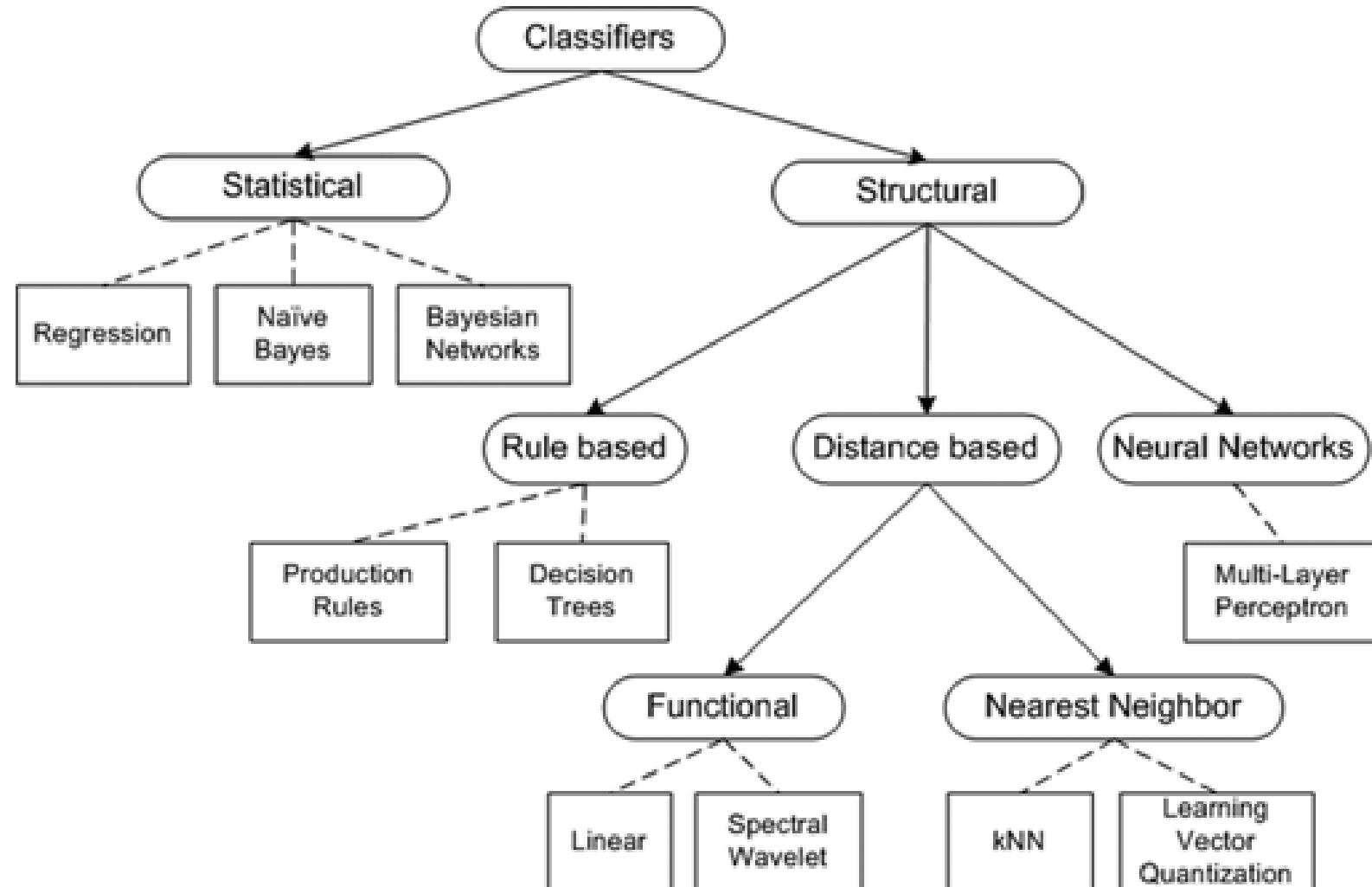
exp = explainer.explain_instance(X_test[i], random_forest.predict_proba, num_features=4)
exp.show_in_notebook(show_table=True, show_all=False)
```



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:

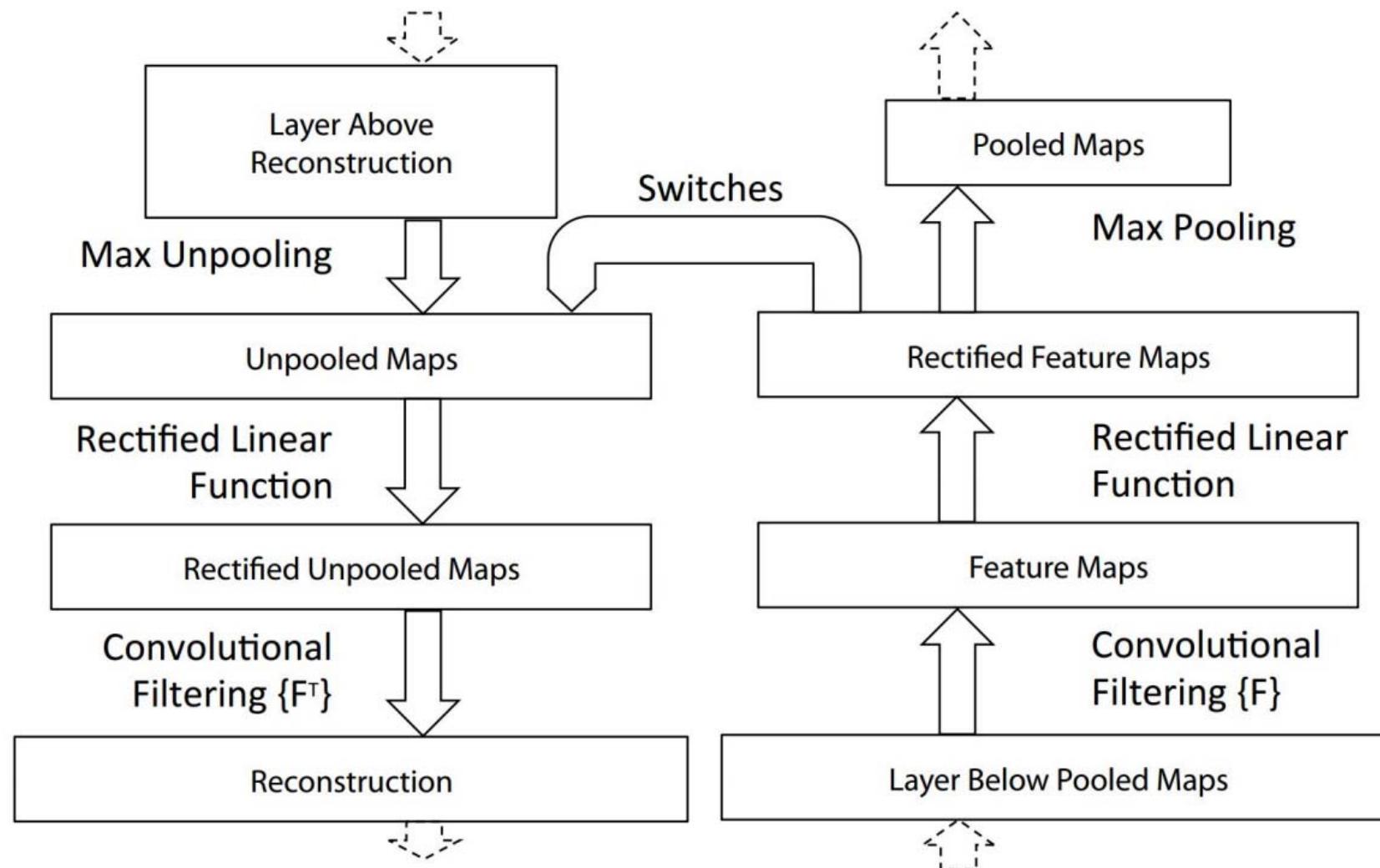
If Past-Depression =Yes and Tiredness =No and Exercise =No and Insomnia =Yes, then Depression

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

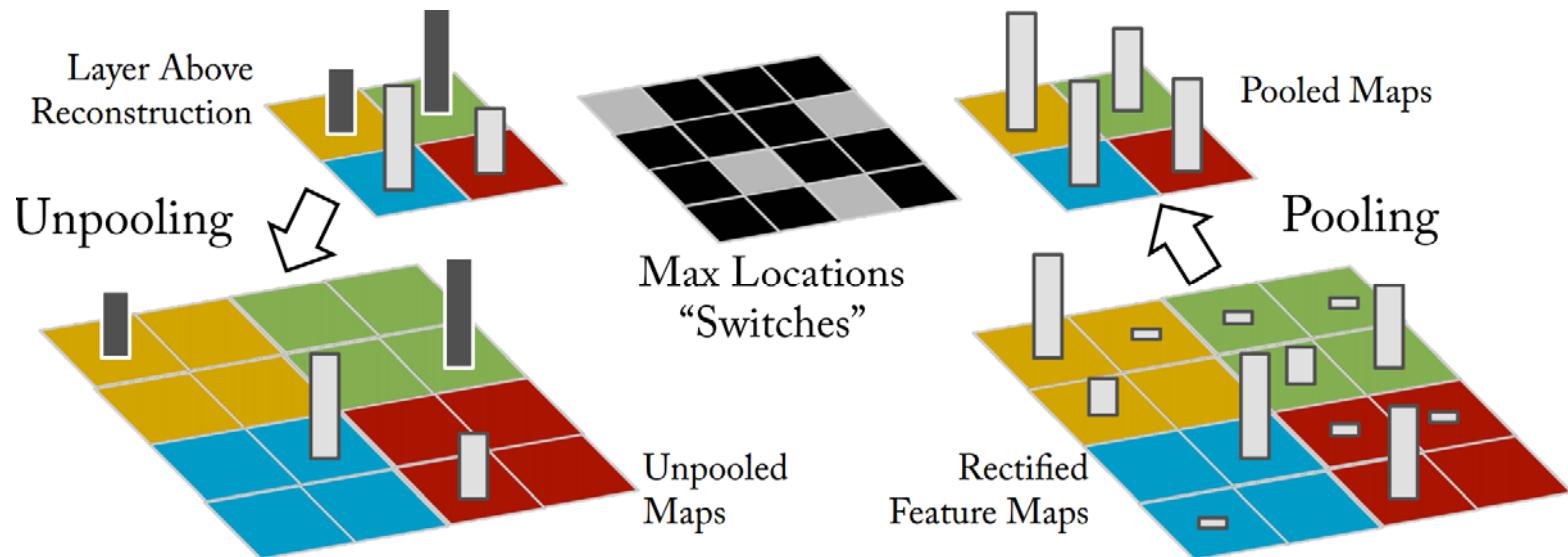
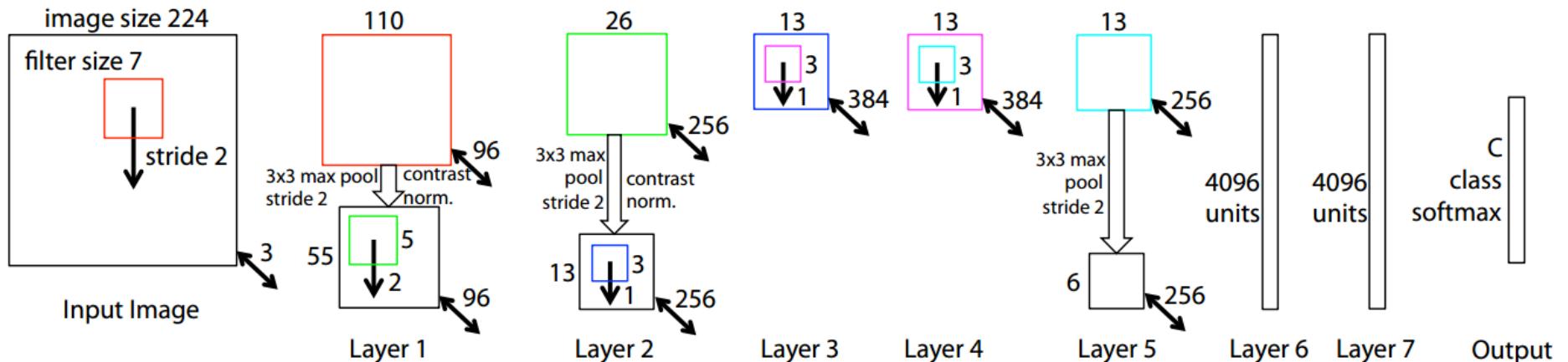
Himabindu Lakkaraju, Ece Kamar, Rich Caruana & Jure Leskovec 2017. Interpretable and Explorable Approximations of Black Box Models. arXiv:1707.01154.

Example: Interpretable Deep Learning Model

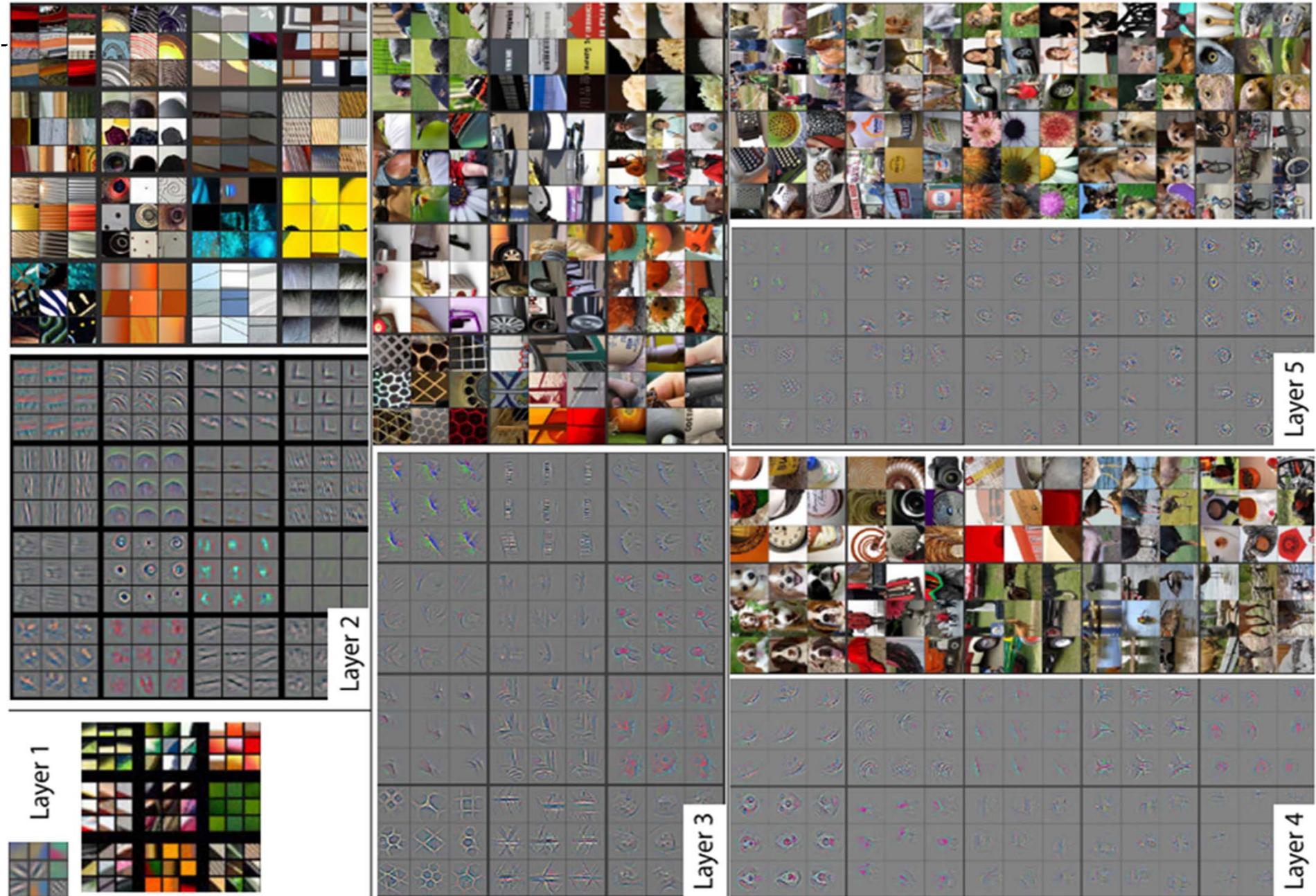


Matthew D. Zeiler & Rob Fergus 2013. Visualizing and Understanding Convolutional Networks. arXiv:1311.2901.

Visualizing a Conv Net with a De-Conv Net

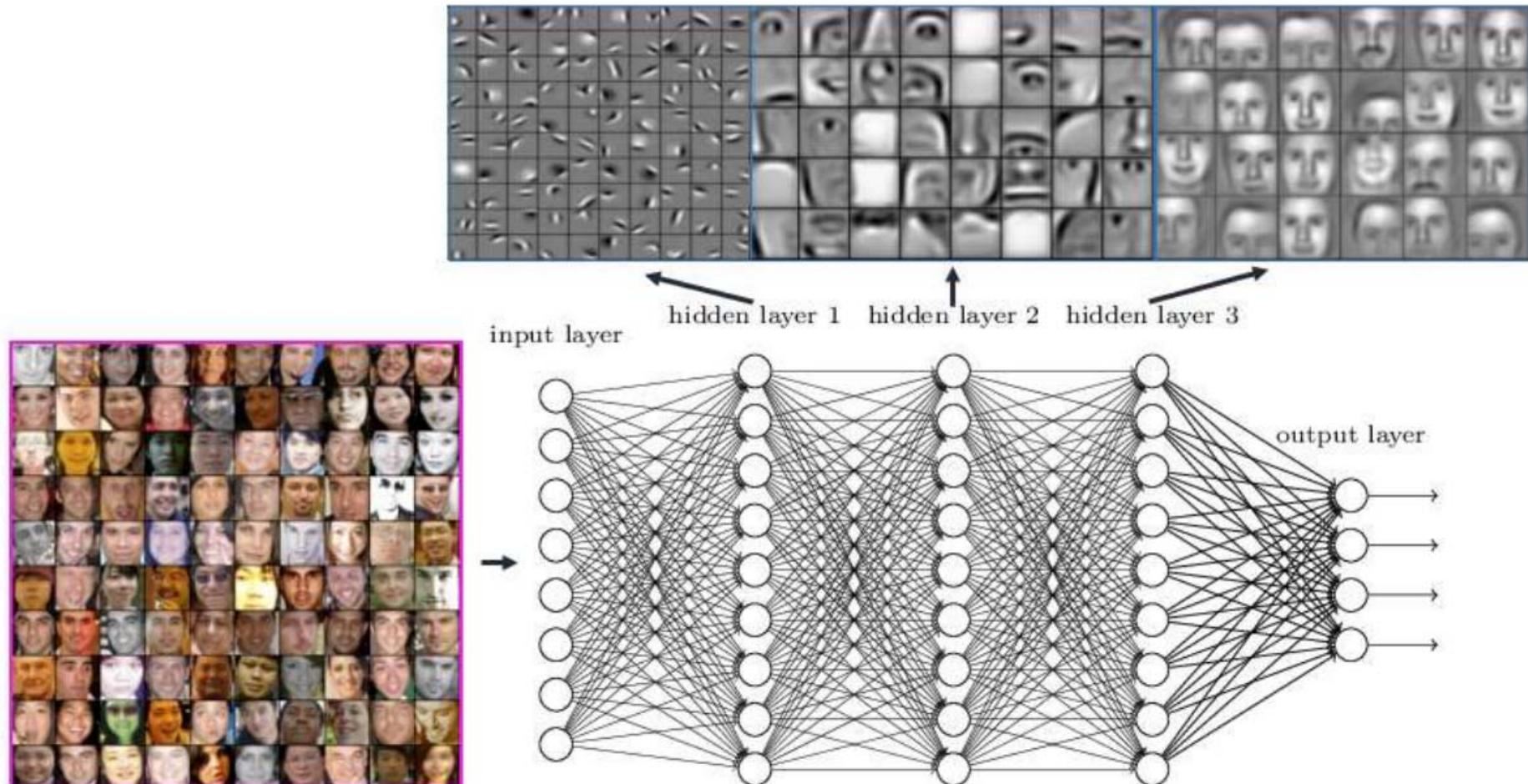


Matthew D. Zeiler & Rob Fergus 2014. Visualizing and understanding convolutional networks. In: D., Fleet, T., Pajdla, B., Schiele & T., Tuytelaars (eds.) ECCV, Lecture Notes in Computer Science LNCS 8689. Cham: Springer, pp. 818-833, doi:10.1007/978-3-319-10590-1_53.



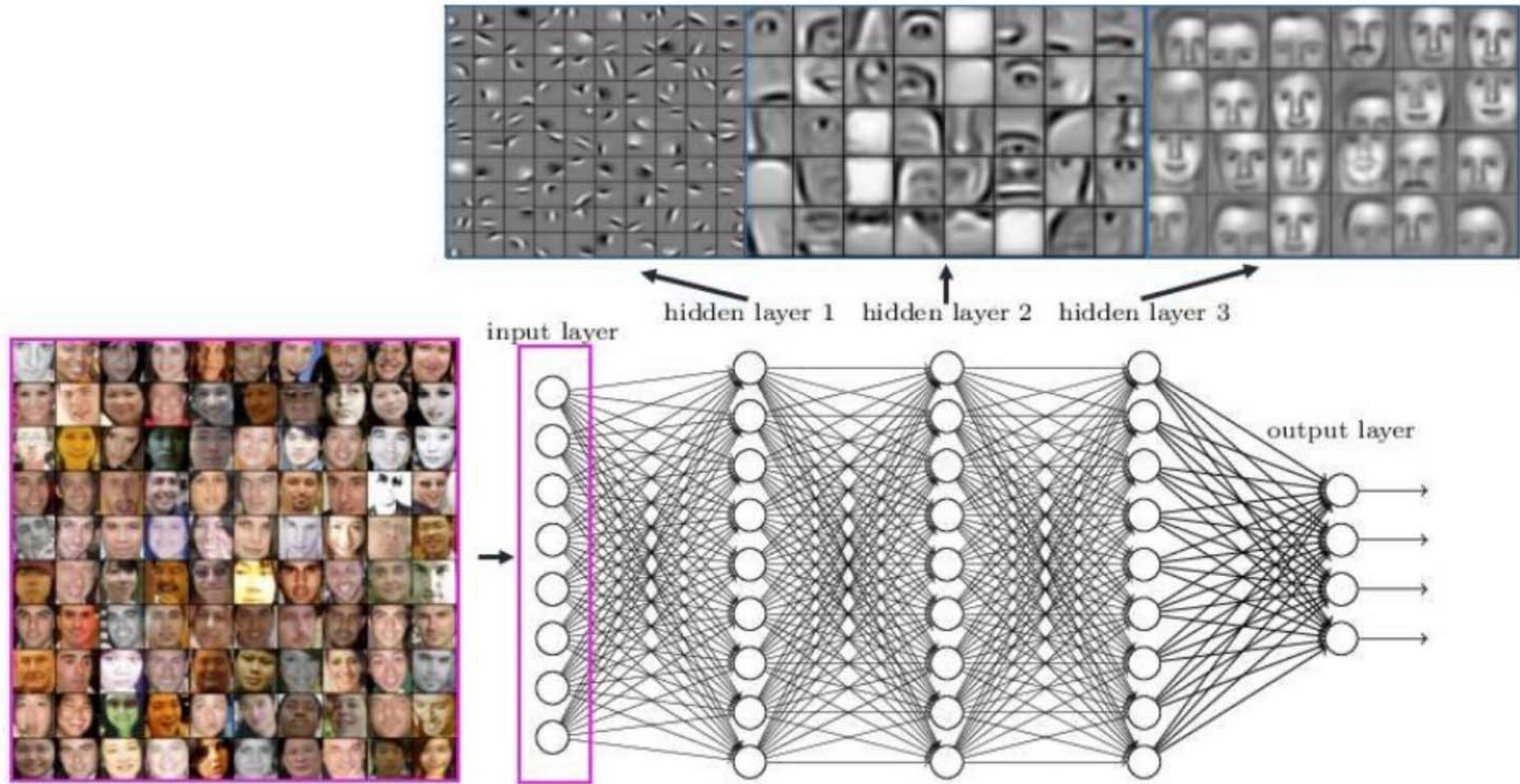
Matthew D. Zeiler & Rob Fergus 2013. Visualizing and Understanding Convolutional Networks. arXiv:1311.2901.
Health Informatics

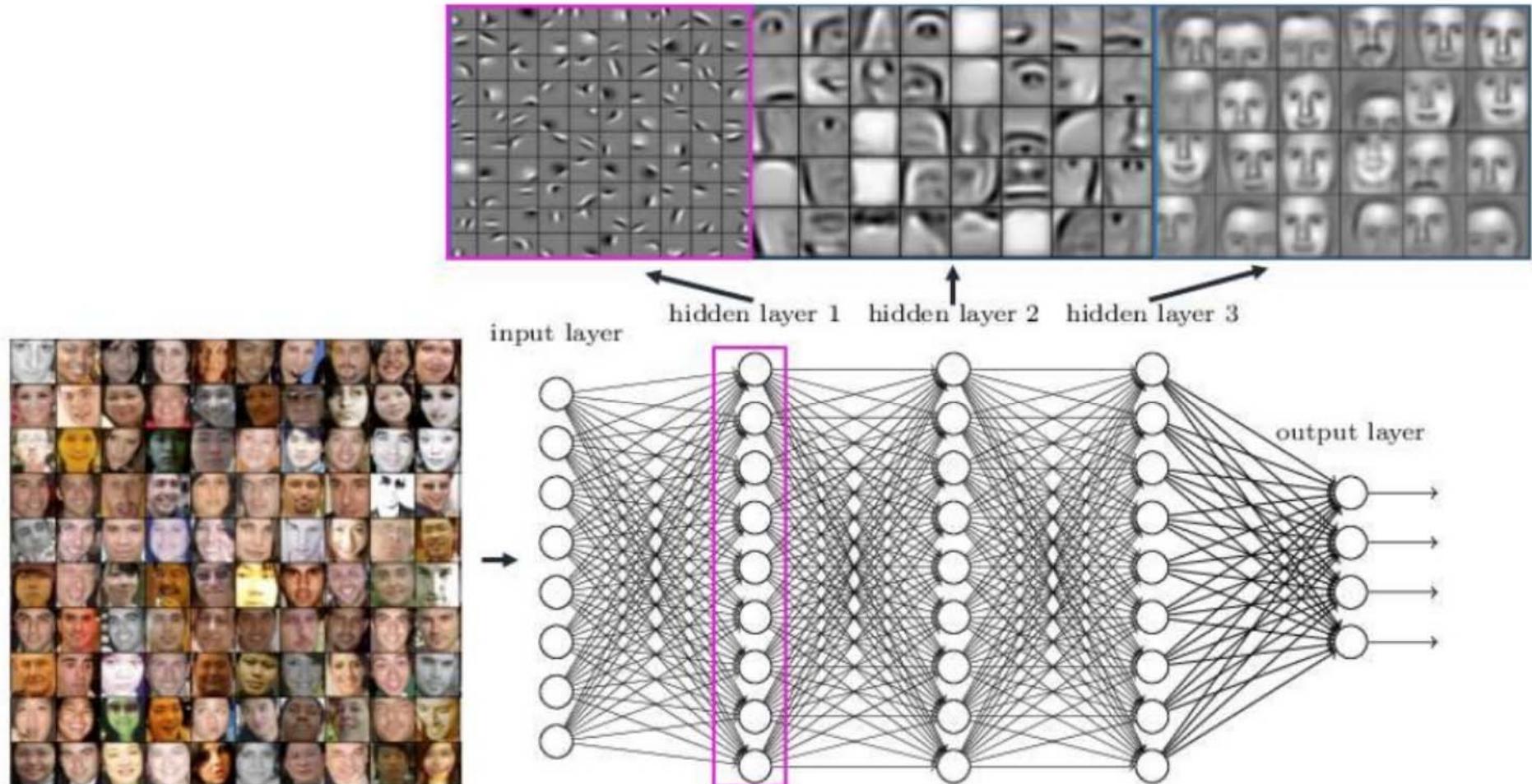
The world is compositional (Yann LeCun)

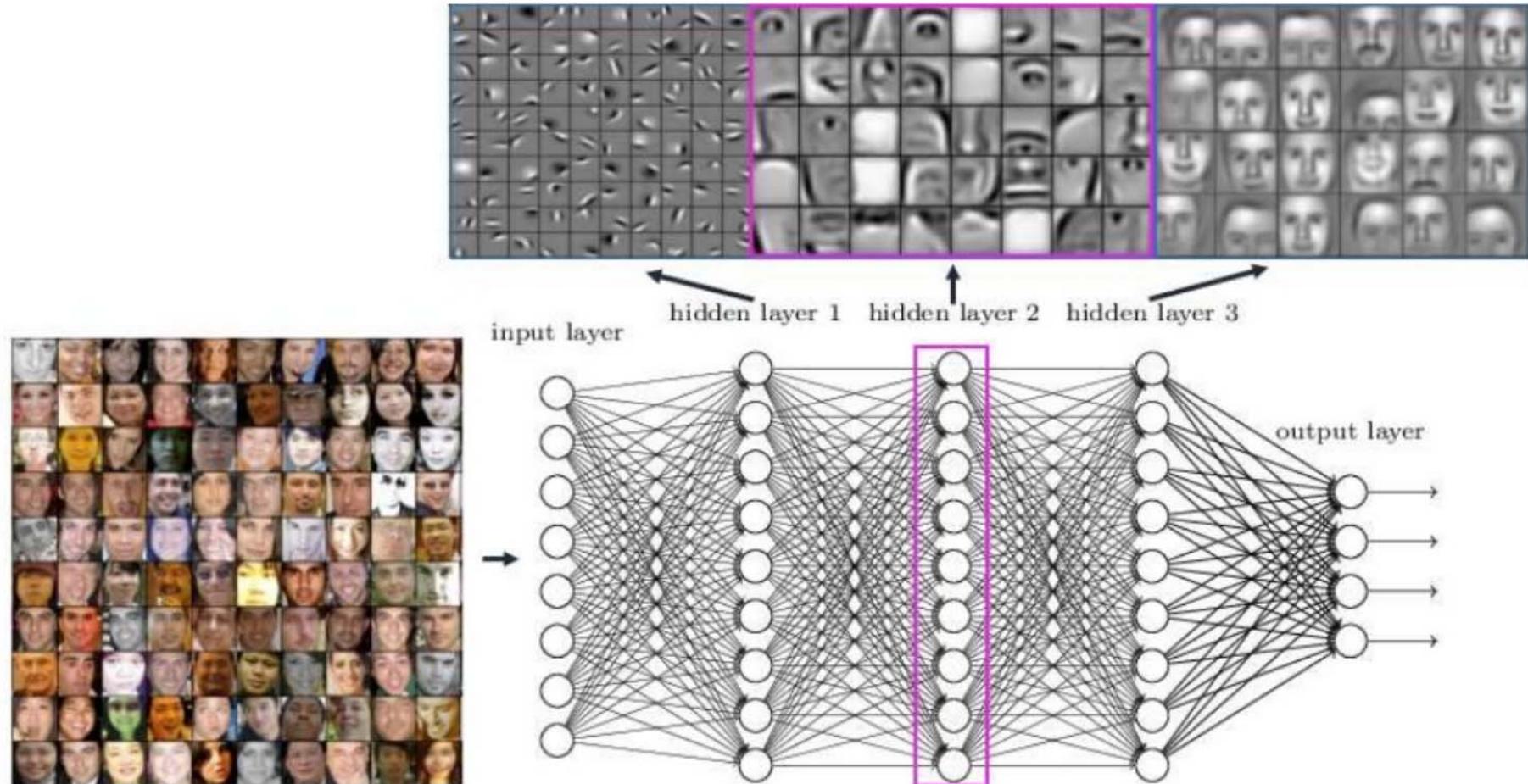


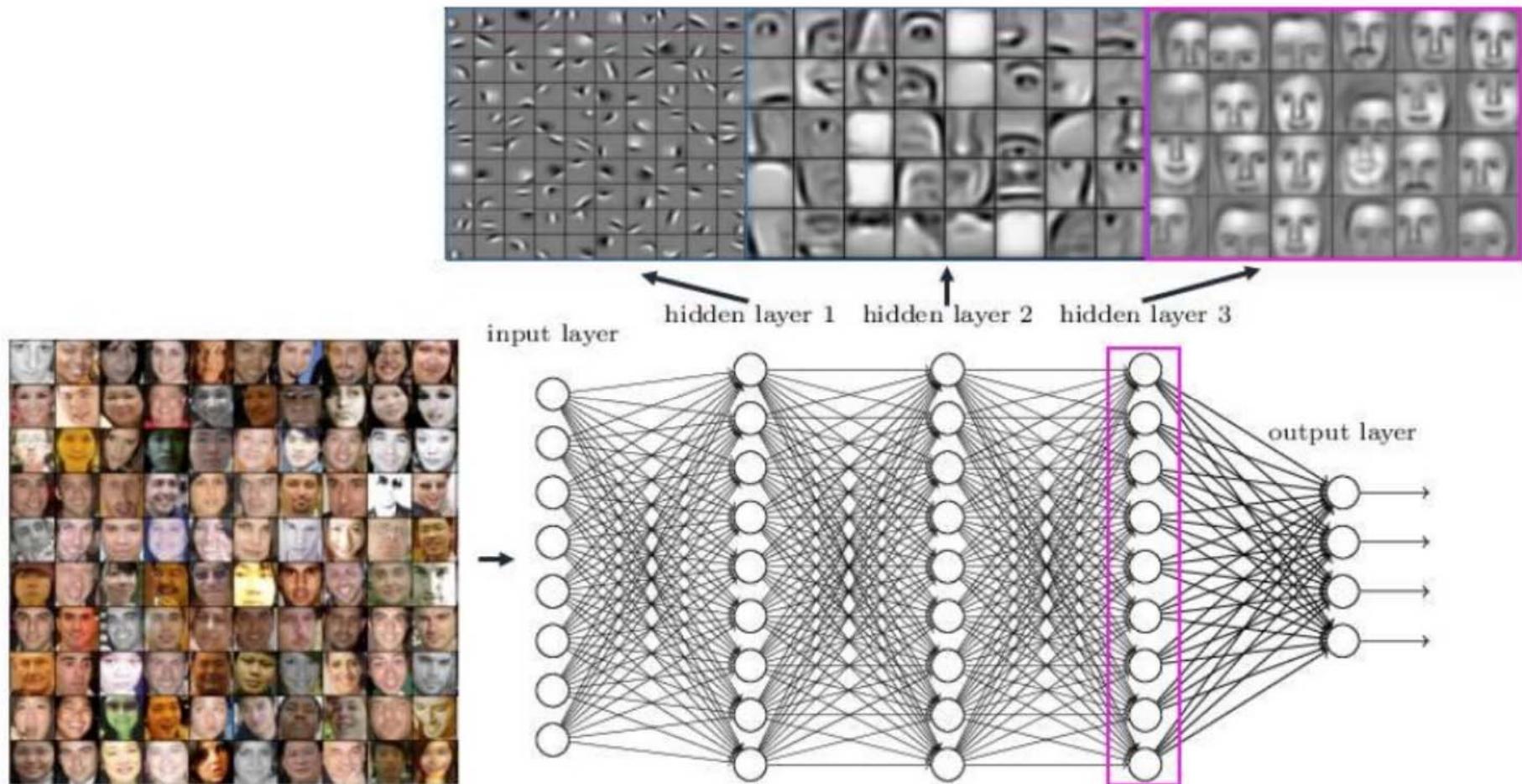
Matthew D. Zeiler & Rob Fergus 2013. Visualizing and Understanding Convolutional Networks. arXiv:1311.2901

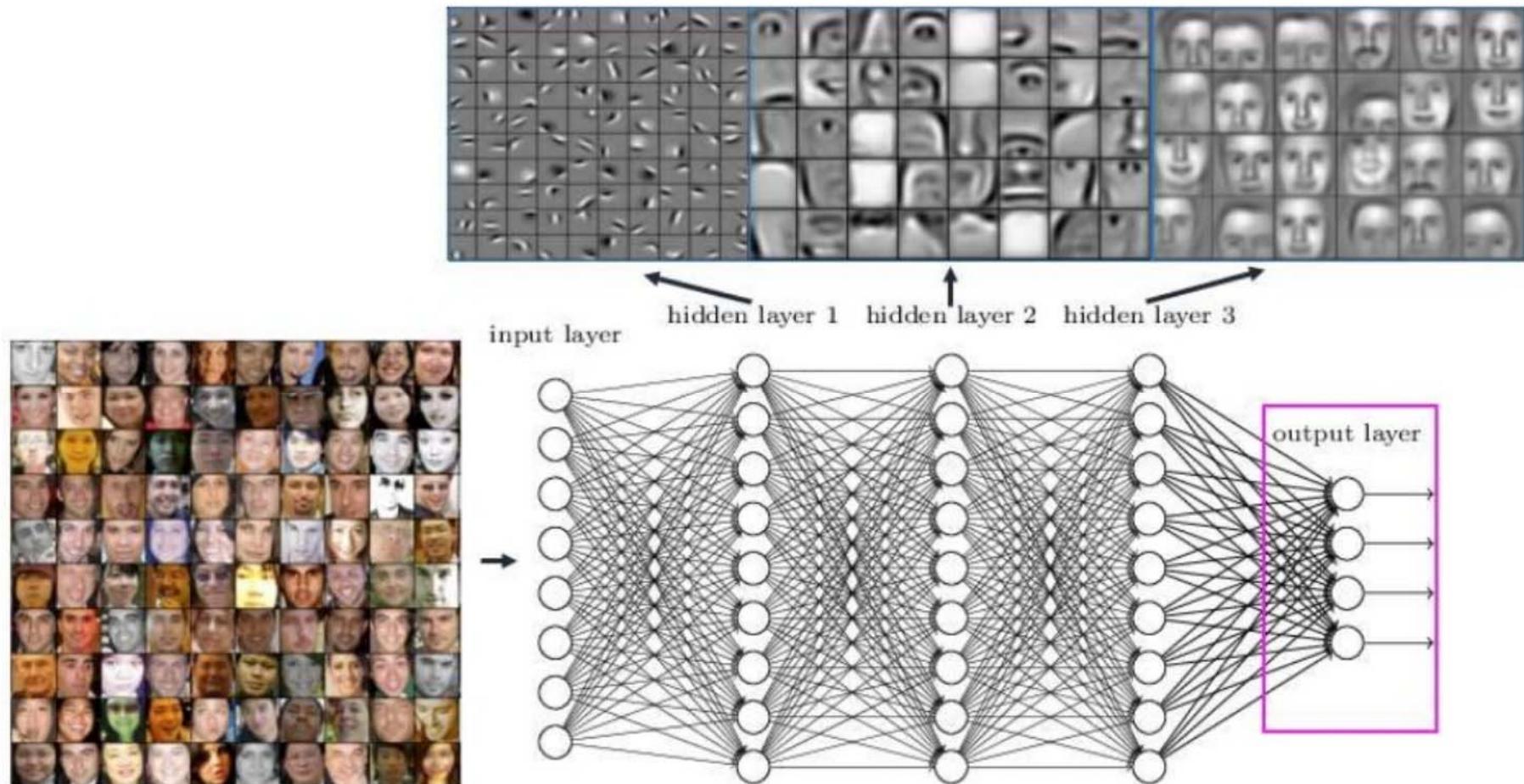
The world is compositional (Yann LeCun)

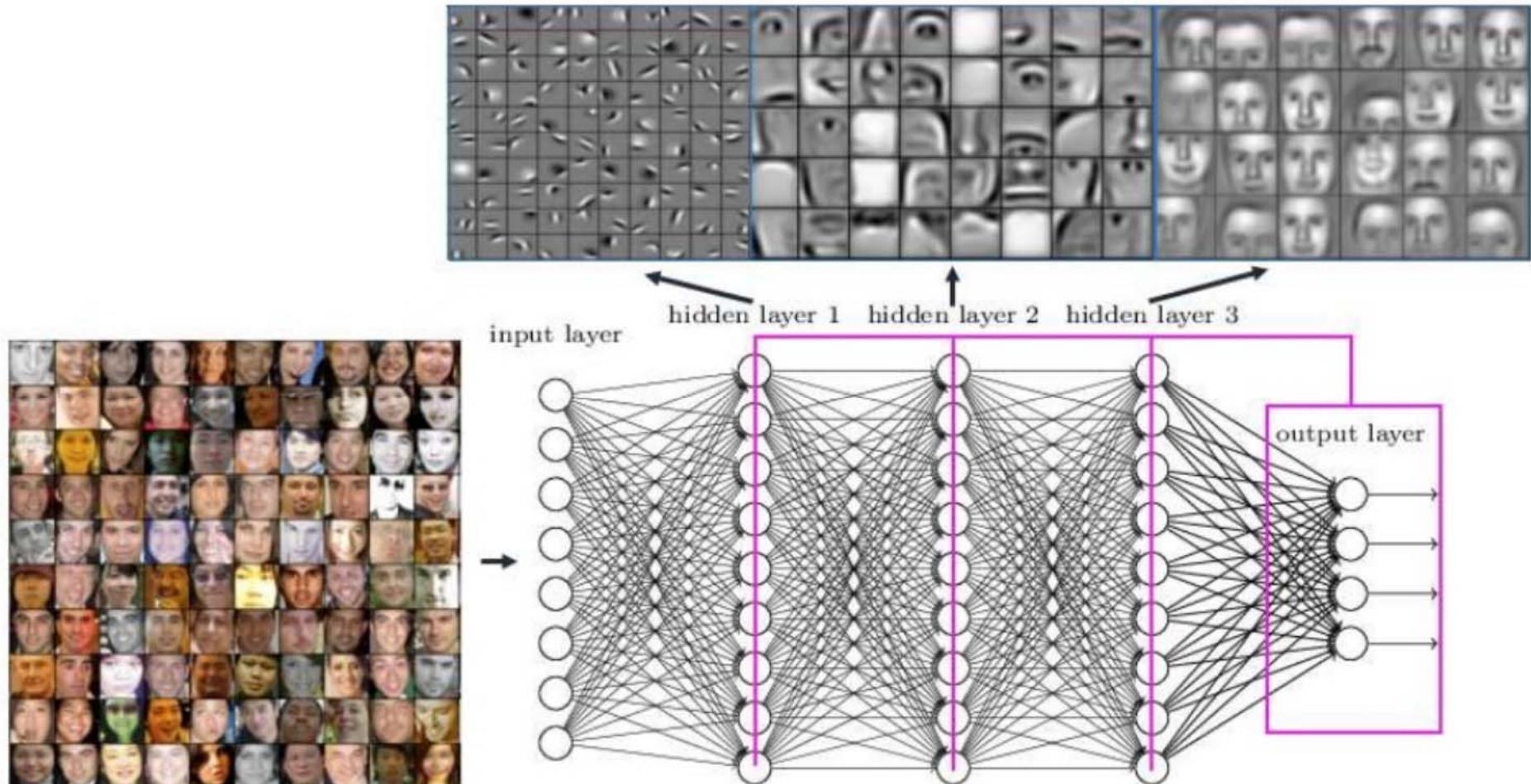






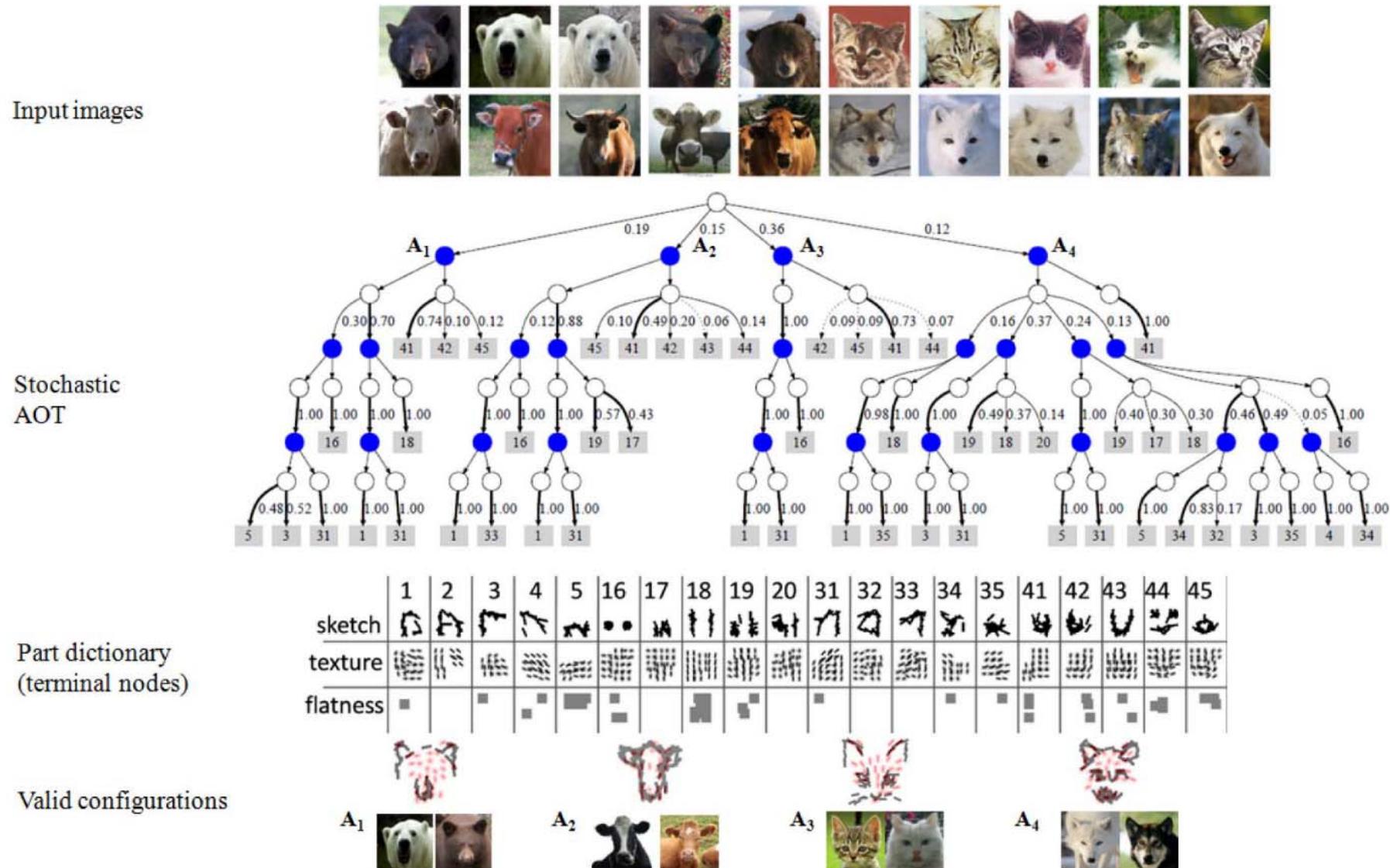






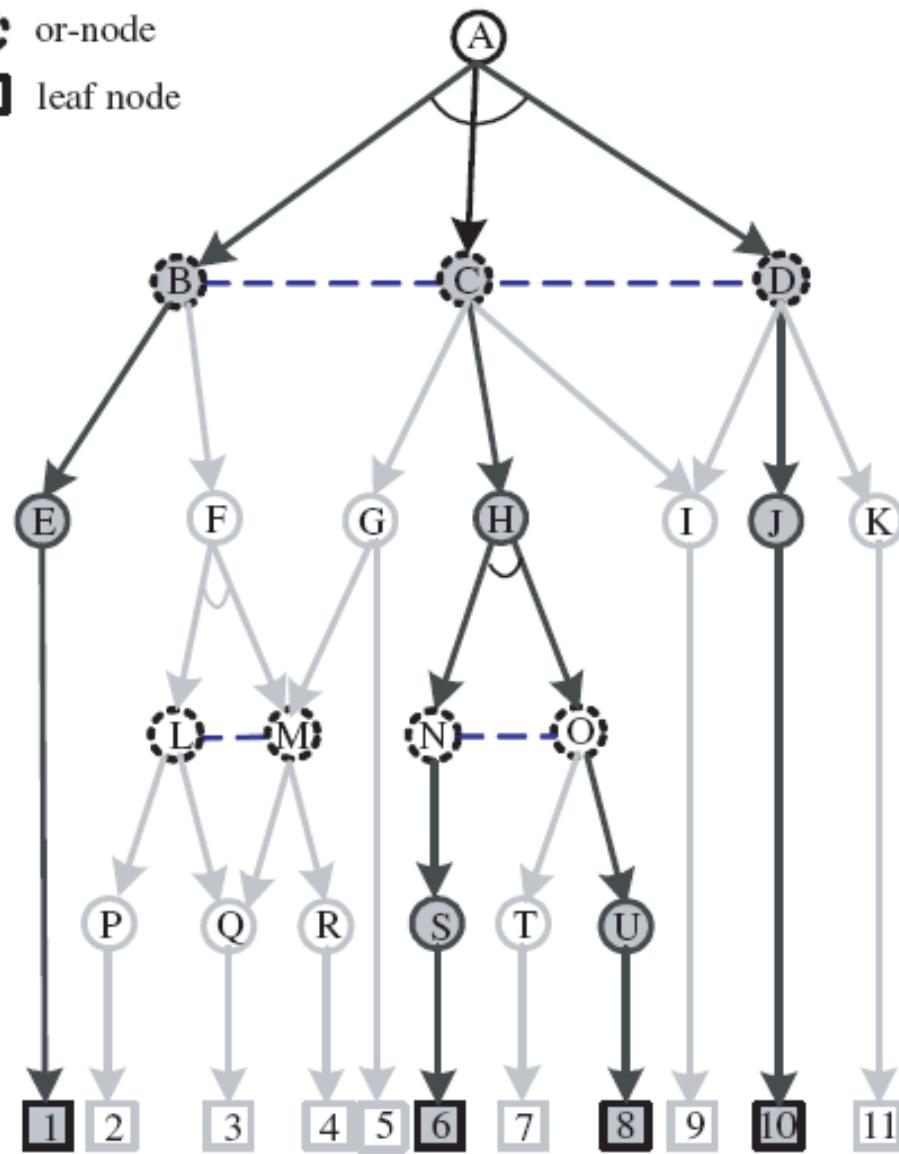
05 Stochastic AOG

Stochastic AND-OR Templates for visual objects



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



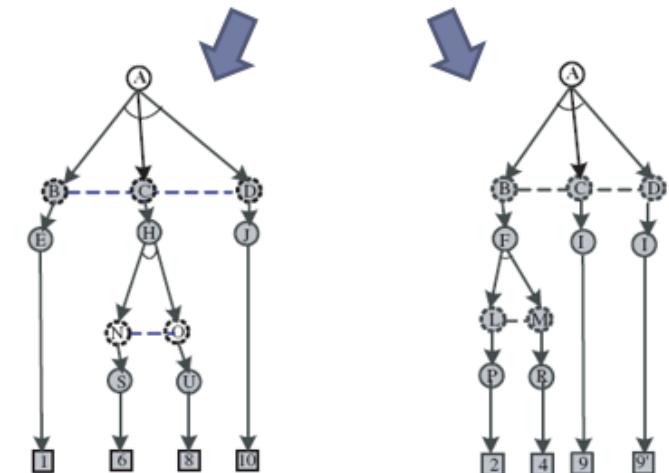
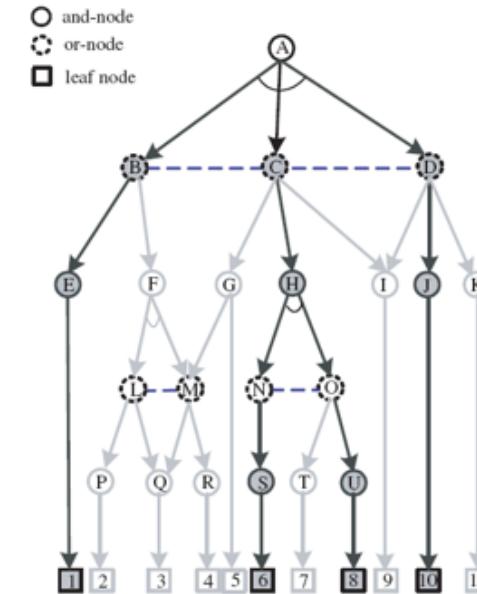
- Algorithm for this framework
 - Top-down/bottom-up computation
- Generalization of small sample
 - Use Monte Carlos simulation to synthesis more configurations
- **Fill semantic gap**

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: $\alpha(t)$

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

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

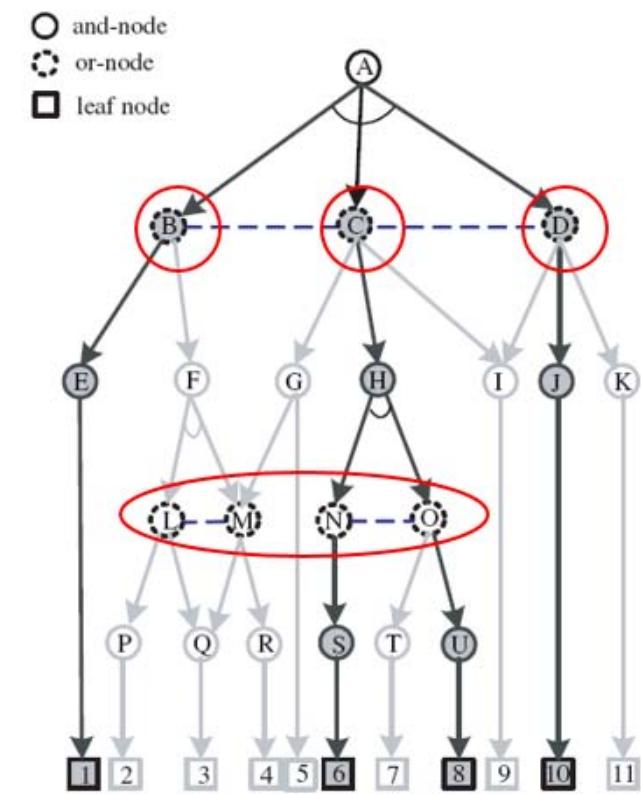


- ▶ Terminal (leaf) node: $T(pg)$
- ▶ And-Or node: $V^{or}(pg), V^{and}(pg)$
- ▶ Set of links: $E(pg)$
- ▶ Switch variable at Or-node: $w(t)$
- ▶ Attributes of primitives: $\alpha(t)$

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

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

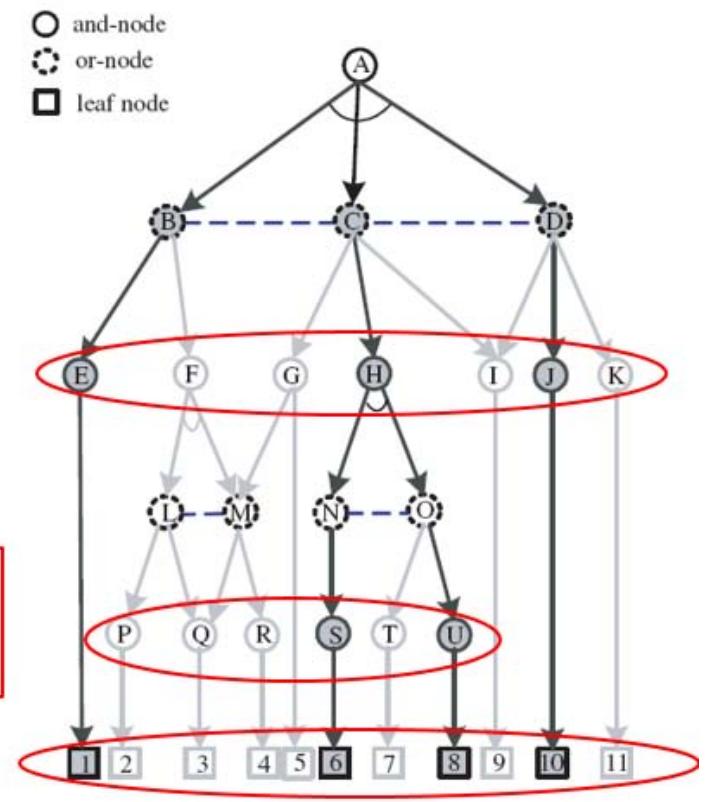
SCFG: weigh the frequency at the children of or-nodes



- ▶ 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: $\alpha(t)$

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

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

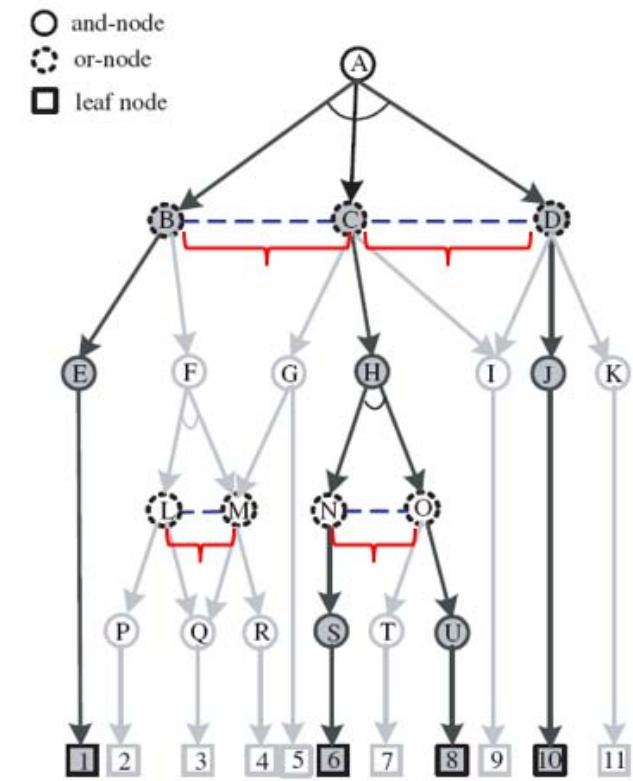


Weigh the local compatibility of primitives (geometric and appearance)

- ▶ 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: $\alpha(t)$

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

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

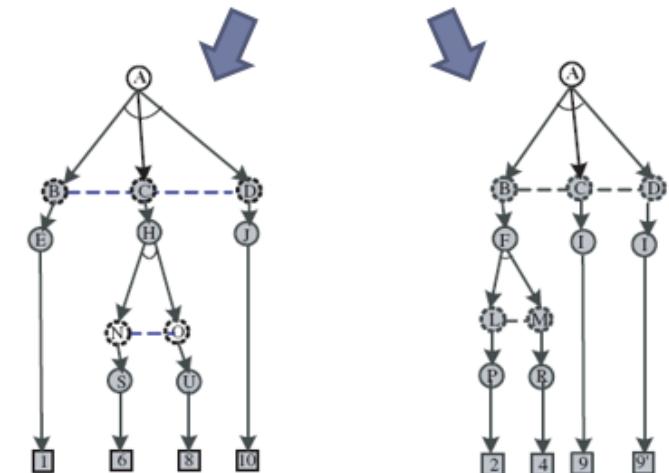
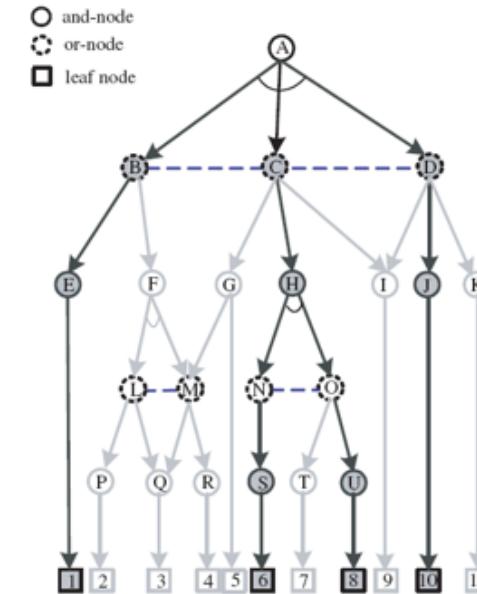


Spatial and appearance between primitives (parts or objects)

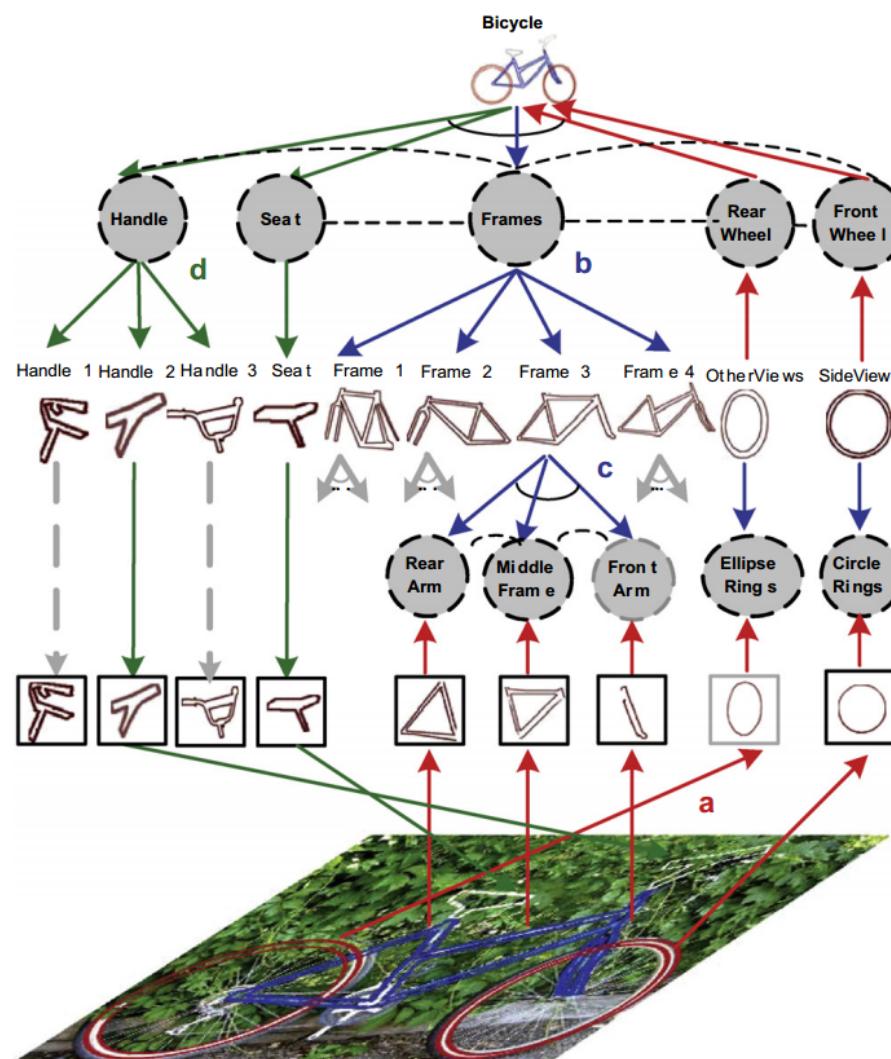
- ▶ Terminal (leaf) node: $T(pg)$
- ▶ And-Or node: $V^{or}(pg), V^{and}(pg)$
- ▶ Set of links: $E(pg)$
- ▶ Switch variable at Or-node: $w(t)$
- ▶ Attributes of primitives: $\alpha(t)$

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

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



Stochastic graph grammar/comp. object representation



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

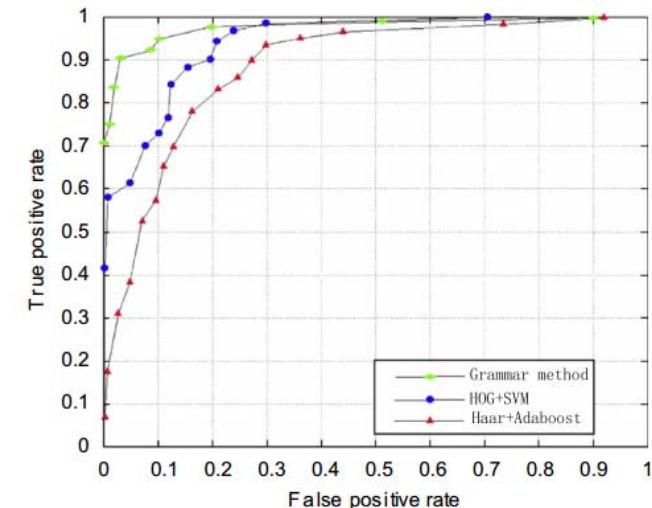
- 3 Call Top-down(A) to update A 's **open** or **closed** lists.

- 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 .

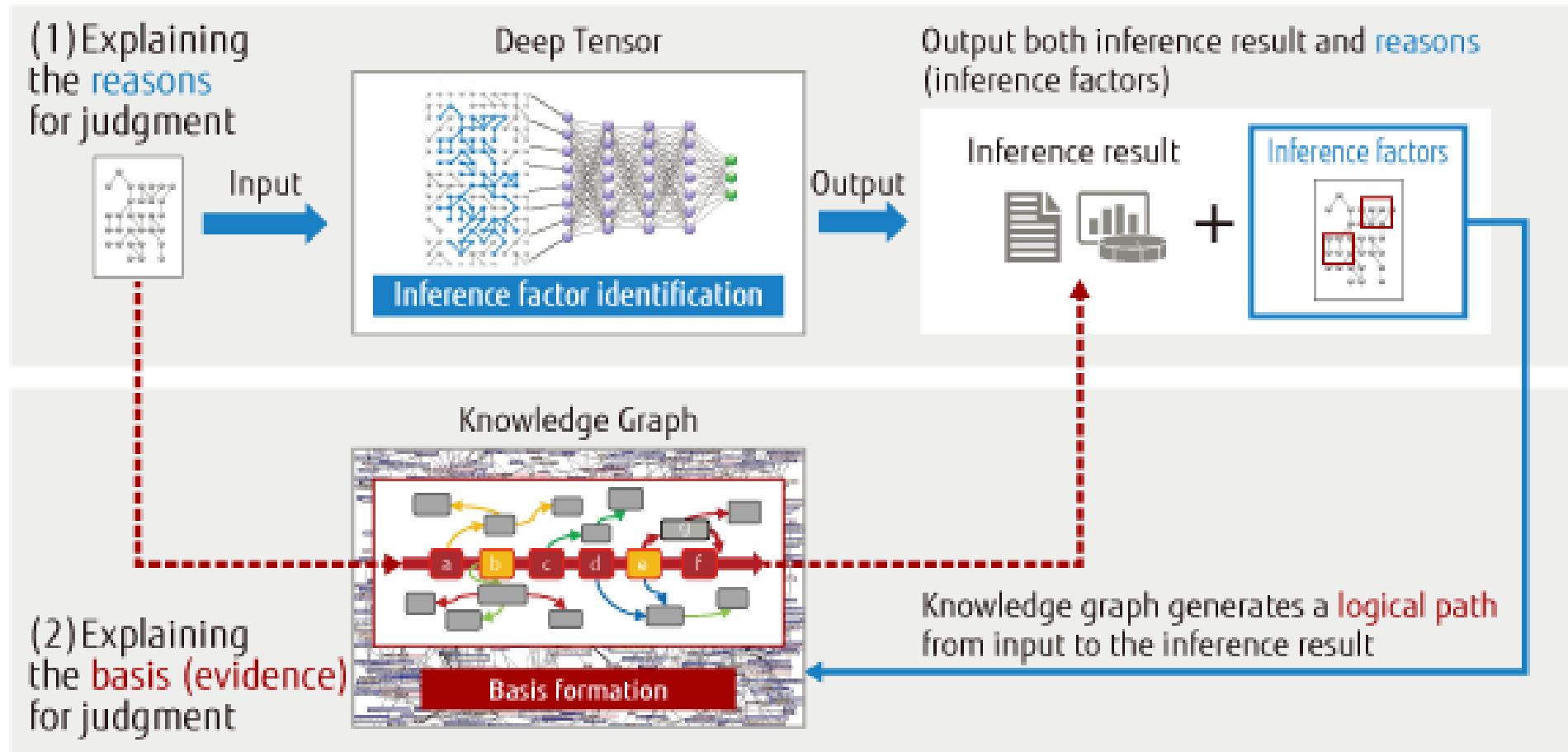
- Until the particles in **open** list with weights higher than the empirical threshold are exhausted. Output all parsing graphs whose root nodes are reached.



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

Combination of Deep Learning with Ontologies



Randy Goebel, Ajay Chander, Katharina Holzinger, Freddy Lecue, Zeynep Akata, Simone Stumpf, Peter Kieseberg & Andreas Holzinger. Explainable AI: the new 42? Springer Lecture Notes in Computer Science LNCS 11015, 2018 Cham. Springer, 295-303, doi:10.1007/978-3-319-99740-7_21.

- When do we need explanations?
 - Of course ONLY in certain situations – most of the time we are happy with automatic approaches!
- 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 enough ... ?



Thank you!

Appendix

Explanations in Artificial Intelligence will be necessary



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

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.

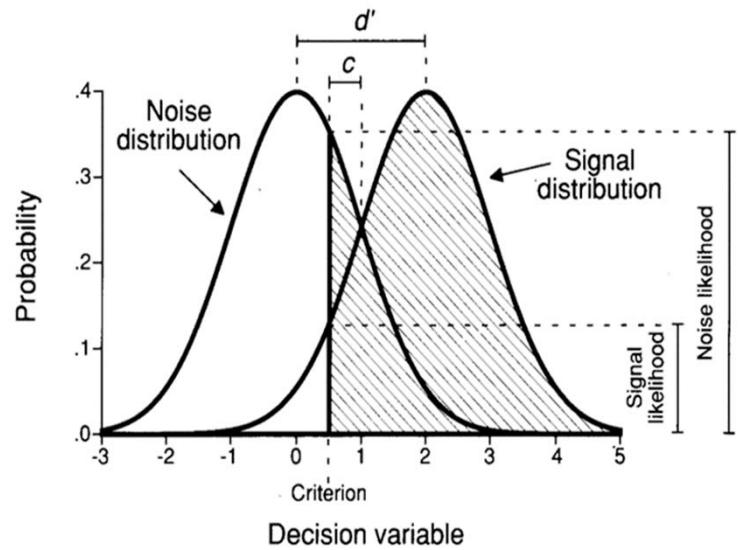
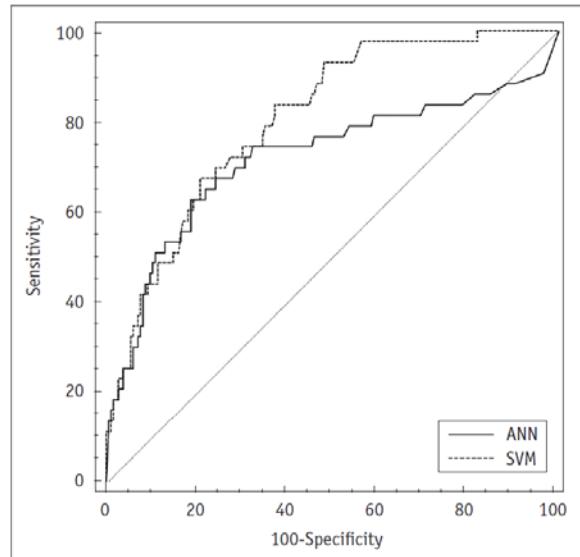
1 Introduction

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

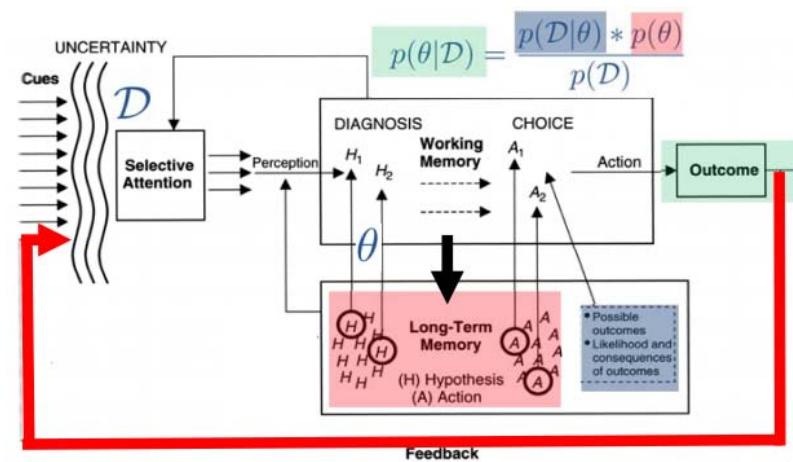
Noel C.F. Codella, Michael Hind, Karthikeyan Natesan Ramamurthy, Murray Campbell, Amit Dhurandhar, Kush R. Varshney, Dennis Wei & Aleksandra Mojsilovic 2018. Teaching Meaningful Explanations. arXiv:1805.11648.

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

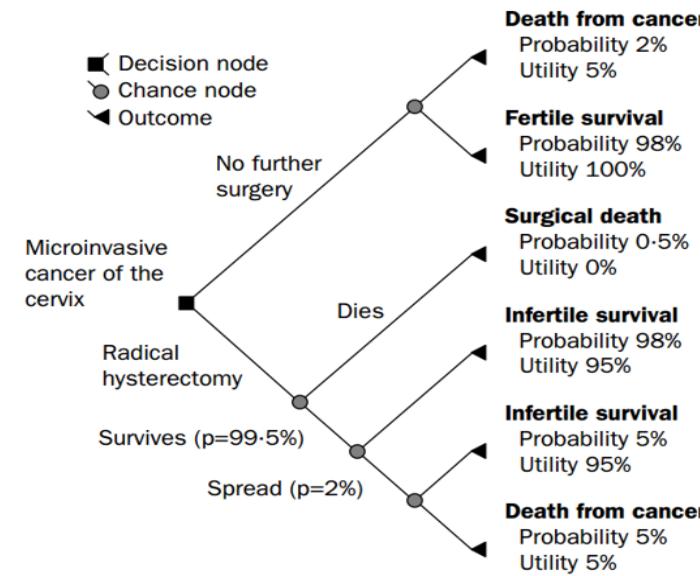
Reflection from last lectures



3



2



4