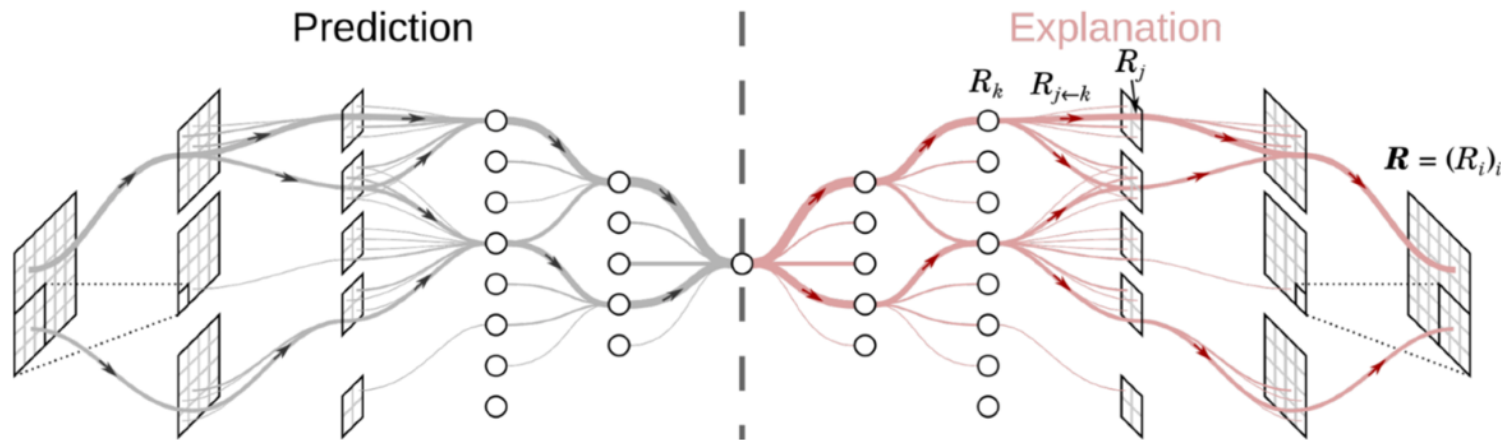


Inspecting AI Like Engineers: From Explanation to Validation with SemanticLens

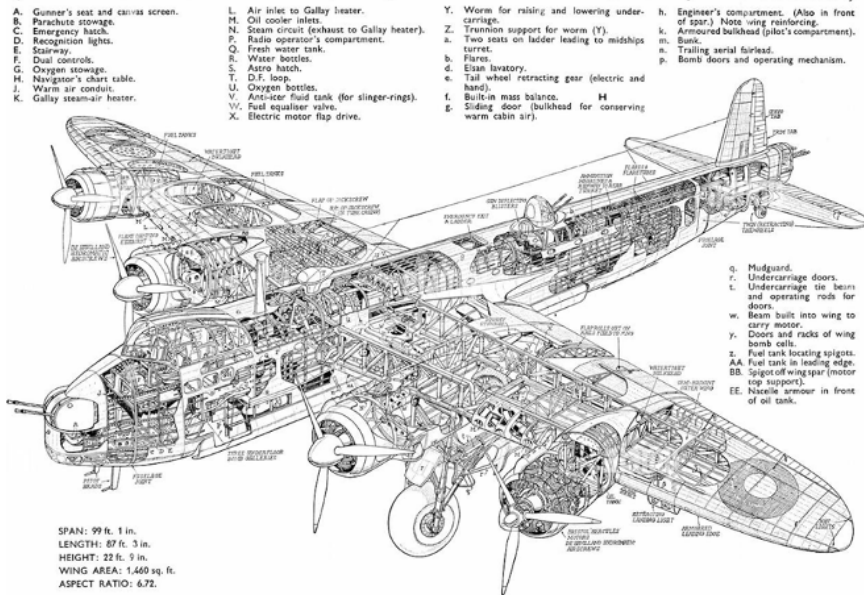
Wojciech Samek
TU Berlin & Fraunhofer HHI



Inspecting AI Like Engineers

Human-Engineered System: Built of known components

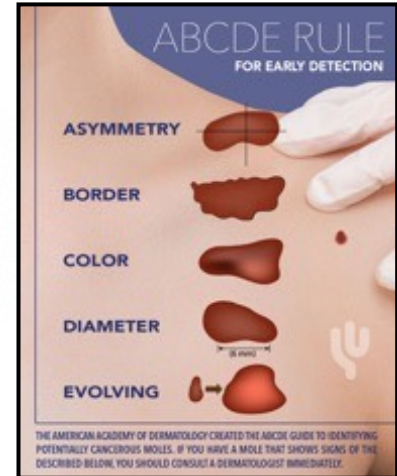
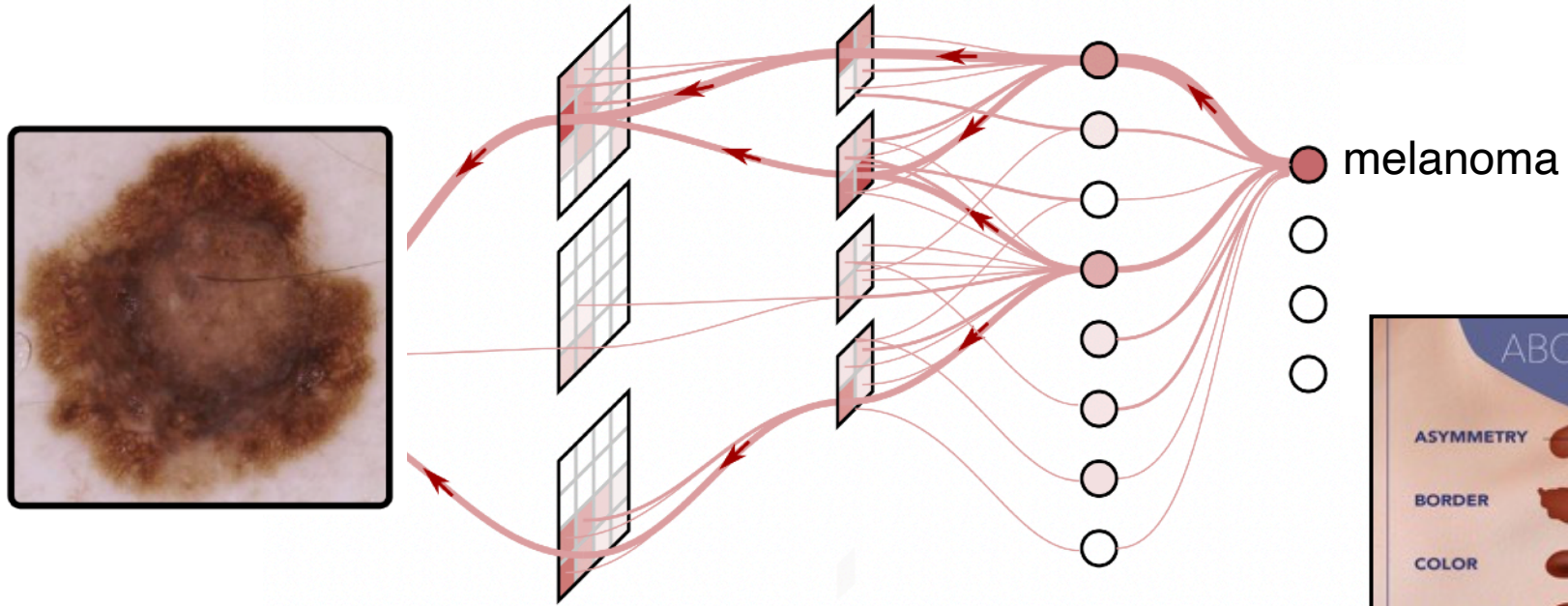
(THE SHORT STIRLING Four 1,600 h.p. Bristol Hercules motors, H.D. Hydromatic airscrews)



Modern AI: Function of components unknown



Does my AI model follow the ABCDE rule ?



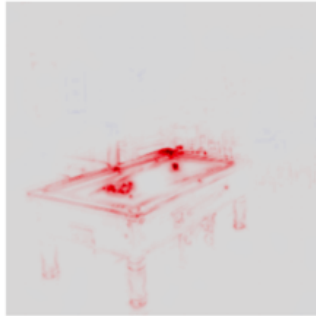
To trust or not to trust
AI; that is the question

We need to to understand
the “Black Box” at
component-level

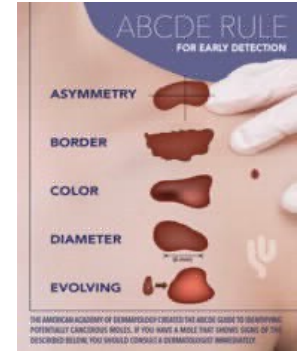
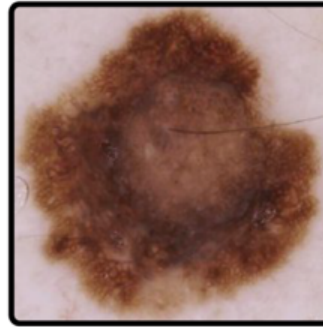


Explainable AI Research

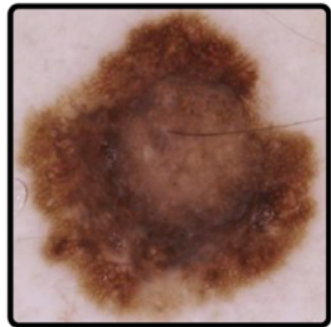
Relevance-Based: Where does AI look at ?



Concept-Based: Which concepts / rules does AI use?



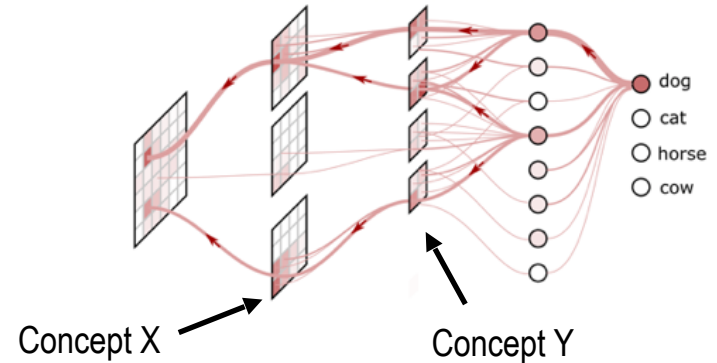
Example-based: Are there similar cases ?



similar →

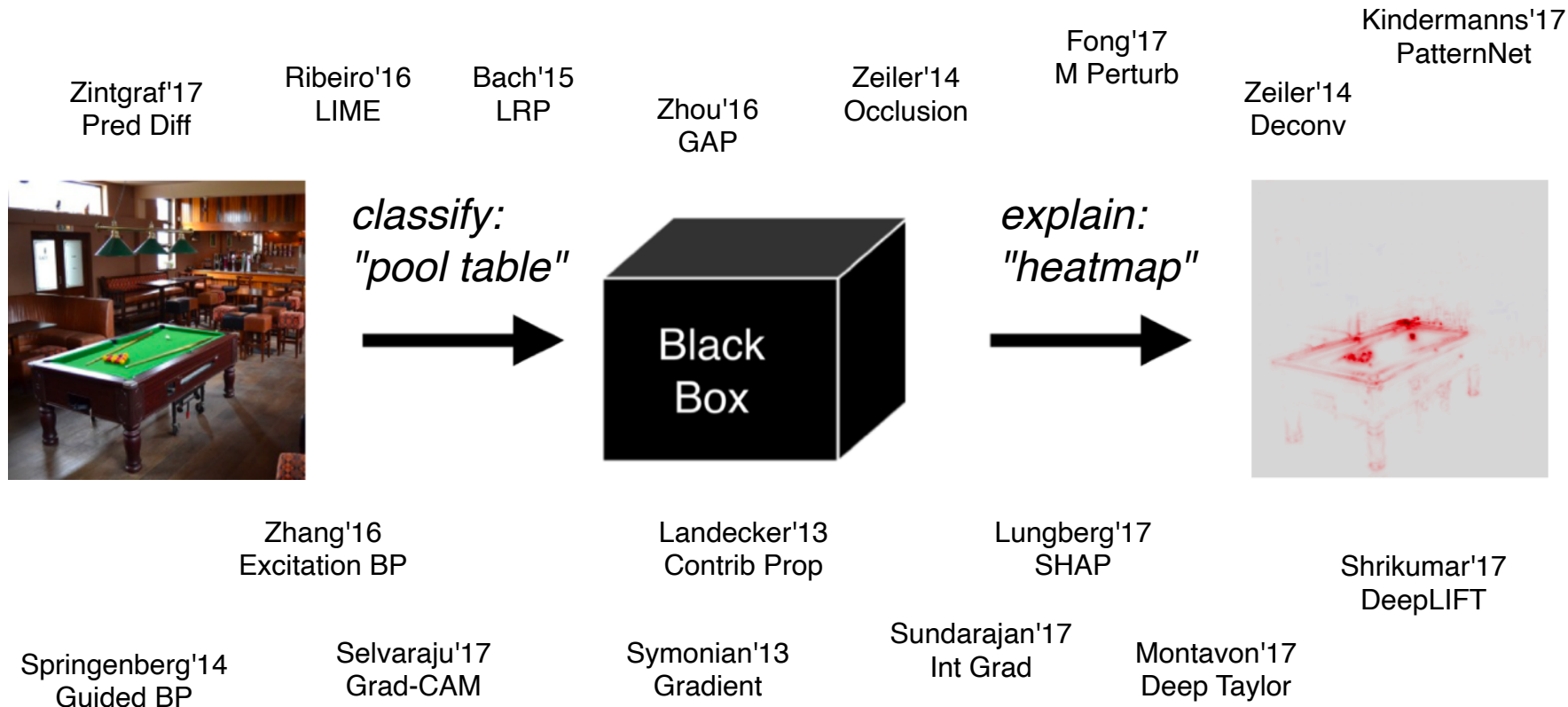


Model-based: What does model represent internally ?

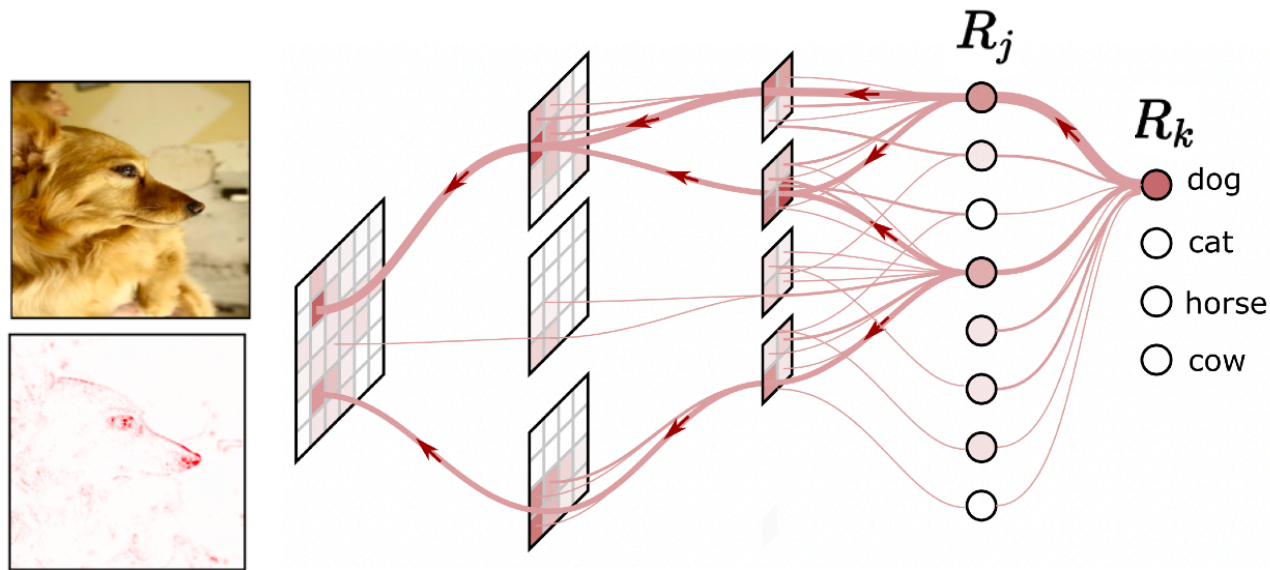


First Wave of XAI: "Understand Prediction"

First Wave of Explainable AI



Layer-wise Relevance Propagation (LRP)



(1) decompose

$$R_{j \leftarrow k} = \frac{z_{jk}}{z_k} R_k$$



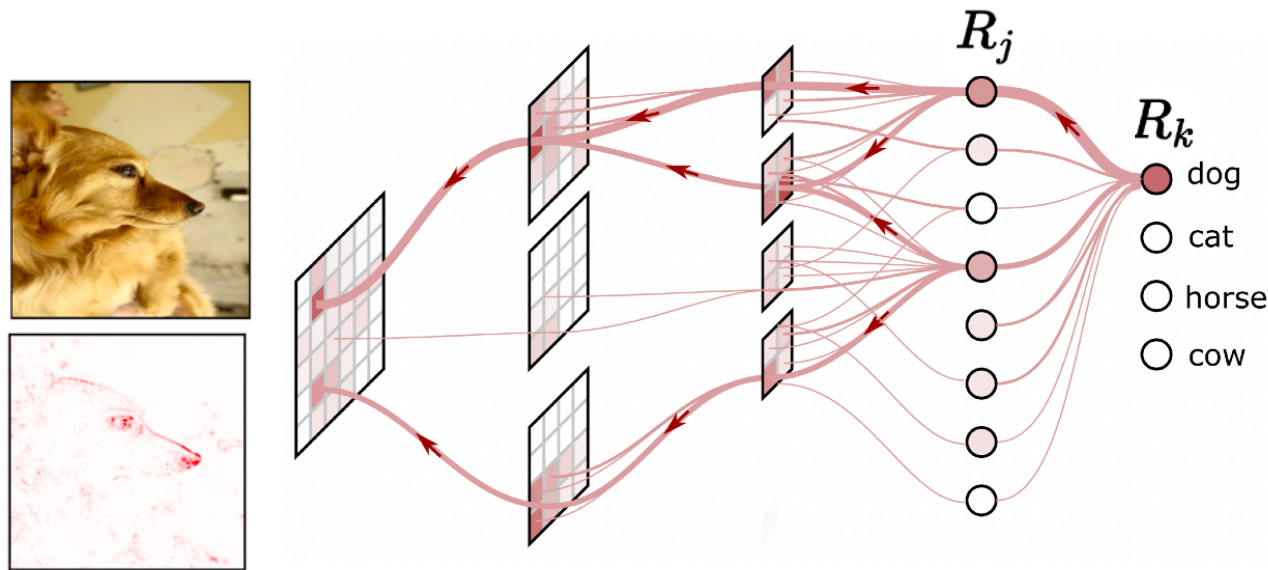
(2) aggregate

$$R_j = \sum R_{j \leftarrow k}$$

z_{jk} measures how much j has contributed to activation of k

(Bach et al. 2015)

Layer-wise Relevance Propagation (LRP)



Layer-wise relevance conservation

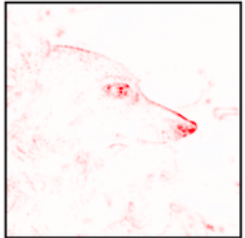
$$\sum_i R_i = \dots = \sum_i R_i^{(l)} = \sum_j R_j^{(l+1)} = \dots = f(x)$$

Advantages

- efficient & faithful
- relevance values for all elements of NN
- applicable to non-differentiable layers (no gradient shattering)

Which redistribution rule is the right one (i.e. how to best measure z_{jk})?

Layer-wise Relevance Propagation (LRP)



Layer-wise relevance
 $\sum_i R_i = \dots$

Convolutional NN

Name	Formula	Usage	DTD
LRP-0 [7]	$R_j = \sum_k \frac{a_j w_{jk}}{\sum_{0,j} a_j w_{jk}} R_k$	upper layers	✓
LRP- ϵ [7]	$R_j = \sum_k \frac{a_j w_{jk}}{\epsilon + \sum_{0,j} a_j w_{jk}} R_k$	middle layers	✓
LRP- γ	$R_j = \sum_k \frac{a_j (w_{jk} + \gamma w_{jk}^+)}{\sum_{0,j} a_j (w_{jk} + \gamma w_{jk}^+)} R_k$	lower layers	✓
LRP- $\alpha\beta$ [7]	$R_j = \sum_k \left(\alpha \frac{(a_j w_{jk})^+}{\sum_{0,j} (a_j w_{jk})^+} - \beta \frac{(a_j w_{jk})^-}{\sum_{0,j} (a_j w_{jk})^-} \right) R_k$	lower layers	\times^*
flat [30]	$R_j = \sum_k \frac{1}{\sum_j 1} R_k$	lower layers	\times
w^2 -rule [36]	$R_i = \sum_j \frac{w_{ij}^2}{\sum_i w_{ij}^2} R_j$	first layer (\mathbb{R}^d)	✓
z^B -rule [36]	$R_i = \sum_j \frac{x_i w_{ij} - l_i w_{ij}^+ - h_i w_{ij}^-}{\sum_i x_i w_{ij} - l_i w_{ij}^+ - h_i w_{ij}^-} R_j$	first layer (pixels)	✓

(* DTD interpretation only for the case $\alpha = 1, \beta = 0$.)

Images

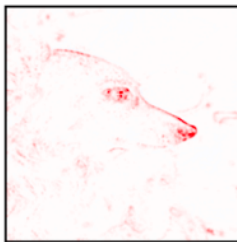
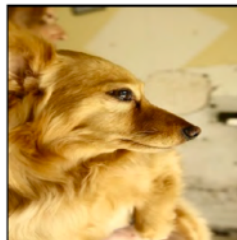
Int & faithful

Relevance values for
 elements of NN

Not applicable to non-
 differentiable layers (no
 backshattering)

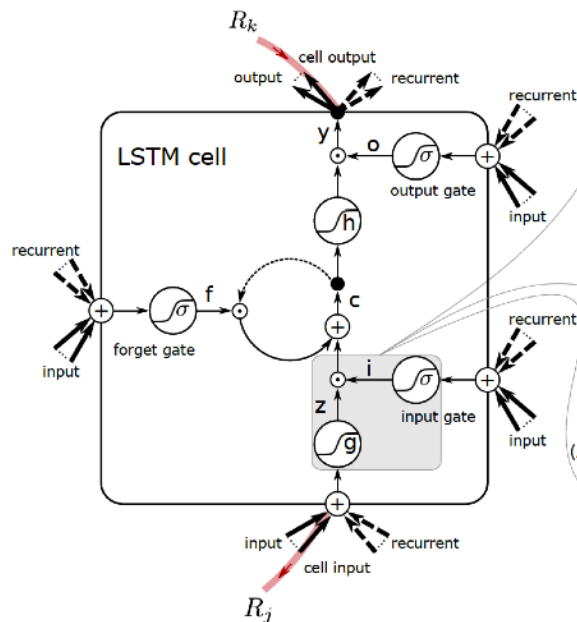
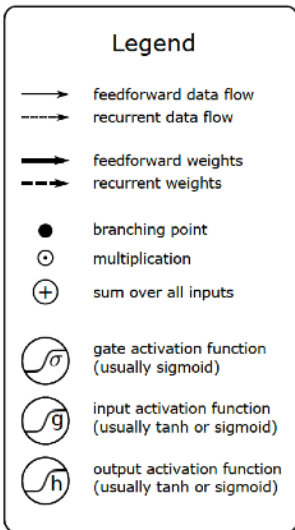
Not a distribution rule
 (not one (i.e. how
 to measure z_{jk})?)

Layer-wise Relevance Propagation (LRP)

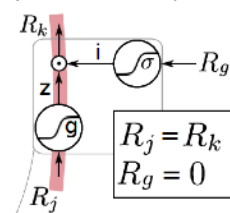


Layer-wise relevance
 $\sum_i R_i = .$

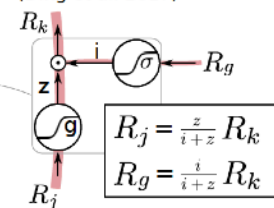
LSTM



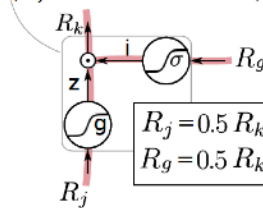
LRP-all (Arras et al. 2017)



LRP-prop (Ding et al. 2017)



LRP-half (Arjona-Medina et al. 2018)



Advantages

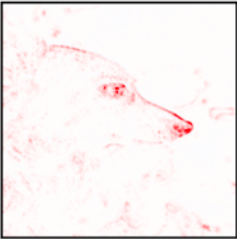
Efficient & faithful

Relevance values for elements of NN

Applicable to non-differentiable layers (no gradient shattering)

Which redistribution rule is the right one (i.e. how best measure zjk)?

Layer-wise Relevance Propagation (LRP)



Layer-wise relevance
 $\sum_i R_i = \dots$

Transformers

$$\text{LayerNorm}(\mathbf{x}) = \frac{x_j - \mathbb{E}[\mathbf{x}]}{\sqrt{\text{Var}[\mathbf{x}] + \varepsilon}} \gamma_j + \beta_j$$
$$\text{RMSNorm}(\mathbf{x}) = \frac{x_j}{\sqrt{\frac{1}{N} \sum_k x_k^2 + \varepsilon}} \gamma_j$$

$$\mathbf{A} = \text{softmax} \left(\frac{\mathbf{Q} \cdot \mathbf{K}^T}{\sqrt{d_k}} \right)$$
$$\mathbf{O} = \mathbf{A} \cdot \mathbf{V}$$
$$\text{softmax}_j(\mathbf{x}) = \frac{e^{x_j}}{\sum_k e^{x_k}}$$

Proposition 3.3 *Decomposing matrix multiplication with a sequential application of the uniform rule (14) and the ε -rule (8) yields the following relevance propagation rule:*

$$R_{ji}^{l-1}(\mathbf{A}_{ji}) = \sum_p \mathbf{A}_{ji} \mathbf{V}_{ip} \frac{R_{jp}^l}{2 \mathbf{O}_{jp} + \varepsilon} \tag{15}$$

Proposition 3.4 *Decomposing LayerNorm or RMSNorm with a Taylor decomposition (4) with reference point, bias or distributing the bias uniformly) yields the identity relevance propagation rule:*

$$R_i^{l-1} = R_i^l \tag{19}$$

Advantages

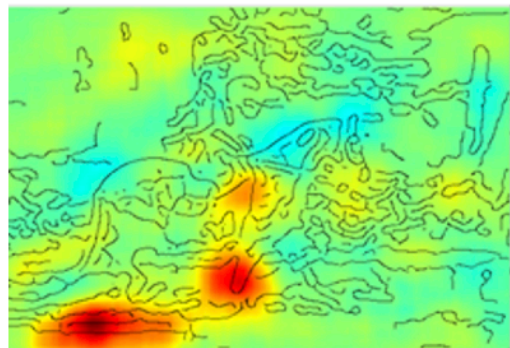
- transparent & faithful
- relevance values for elements of NN
- applicable to non-differentiable layers (no gradient shattering)

redistribution rule
right one (i.e. how to measure zjk)?

What Can We Do ?

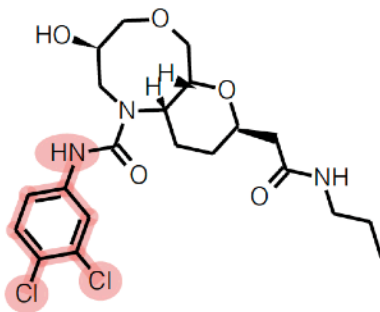
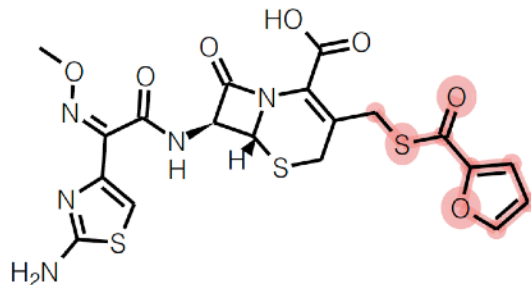
Debug models

(Lapuschkin et al. Nat Comm, 2019)



New insights

(Wong et al. Nature, 2023)



"BLUE XAI"

(Biecek & Samek, ICML, 2024)

Human-values oriented

- Responsible models
- Legal issues
- Trust in predictions
- Ethical issues

Trust in LLMs

(Achtibat et al., ICML, 2024)

Question: In what country is Normandy located?

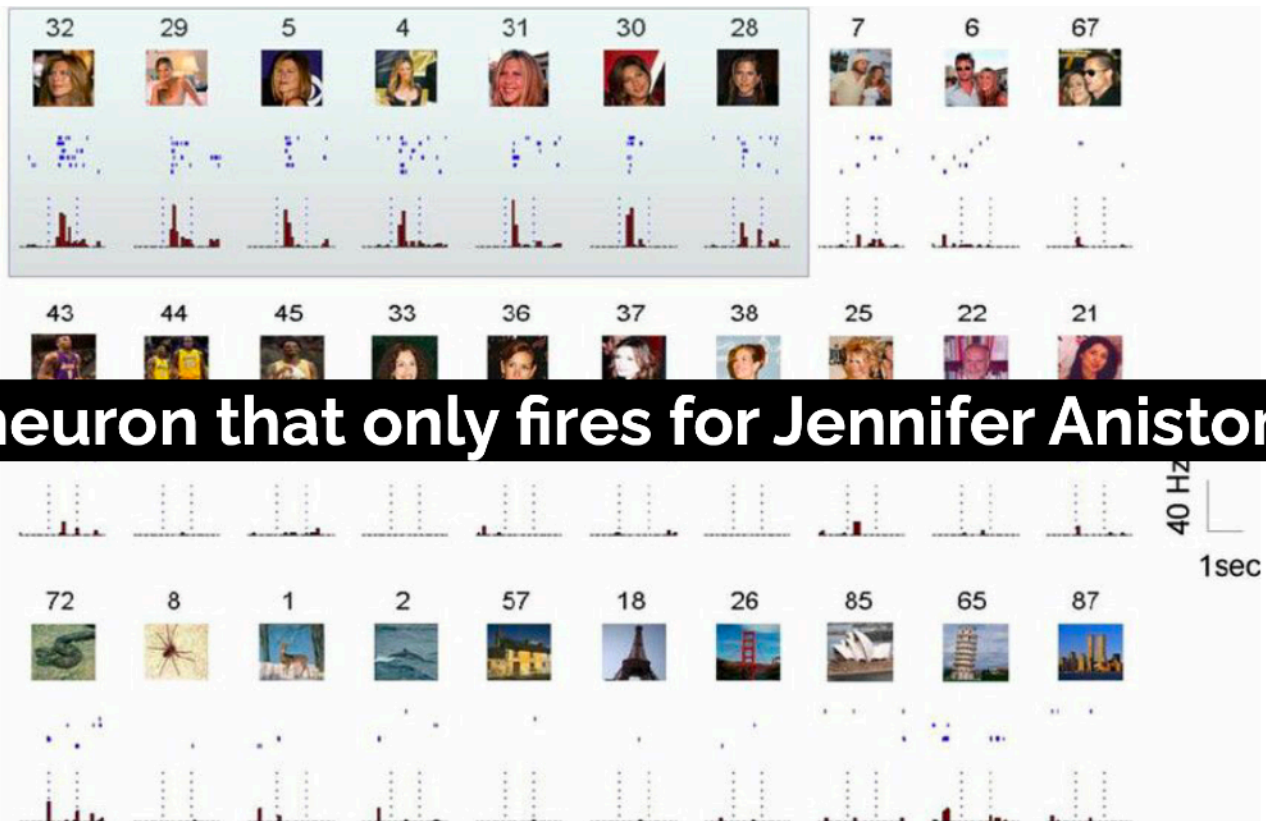
Answer: **France**

AttnLRP

The Normans (Norman: Nourmands; French: Normands; Latin: Normanni) were the people who in the 10th and 11th centuries gave their name to **Normandy**, a region in **France**. They were descended from Norse ("Norman" comes from "Norseman") raiders and pirates from Denmark, Iceland and Norway who, under their leader Rollo, agreed to swear fealty to King Charles III of West Francia. Through generations of assimilation and mixing with the native Frankish and Roman-Gaulish populations, their descendants would gradually merge with the Carolingian-based cultures of West Francia. The distinct cultural and ethnic identity of the Normans emerged initially in the first half of the 10th century, and it continued to evolve over the succeeding centuries.

Second Wave of XAI: "Understand Model"

Interpreting the Model



A neuron that only fires for Jennifer Aniston

Interpreting the Model



nature

Vol 435|23 June 2005|doi:10.1038/nature03687

LETTERS

Invariant visual representation by single neurons in the human brain

R. Quian Quiroga^{1,2,†}, L. Reddy¹, G. Kreiman³, C. Koch¹ & I. Fried^{2,4}



Interpreting the Model

32 29 5 4 31 30 28 7 6 67

nature Vol 435 | 23 June 2005 | doi:10.1038/nature03687

LETTERS

Invari
the human brain

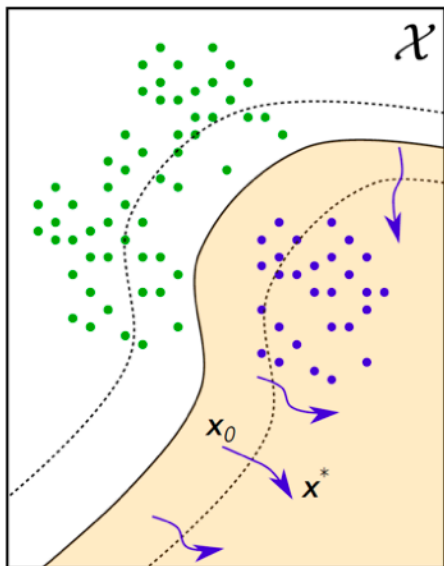
R. Quian Quiroga^{1,2,†}, L. Reddy¹, G. Kreiman³, C. Koch¹ & I. Fried^{2,4}

Do neural networks have a Jennifer Aniston neuron ?

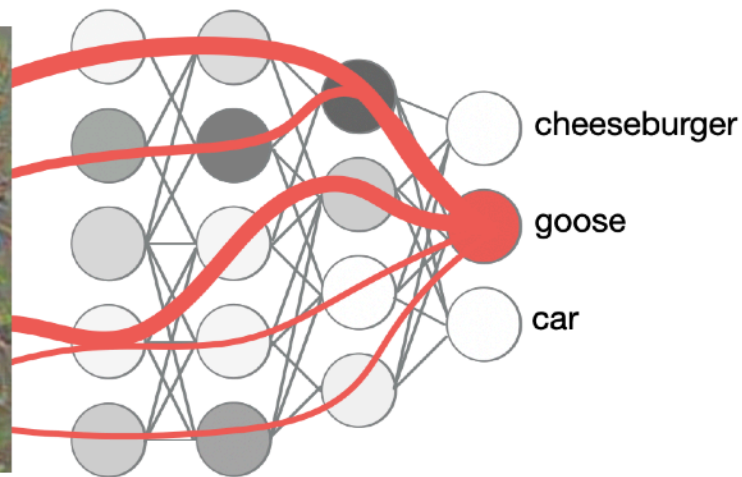
s in

Activation Maximization

Find the input pattern that maximizes class probability.



simple regularizer
(Simonyan et al. 2013)

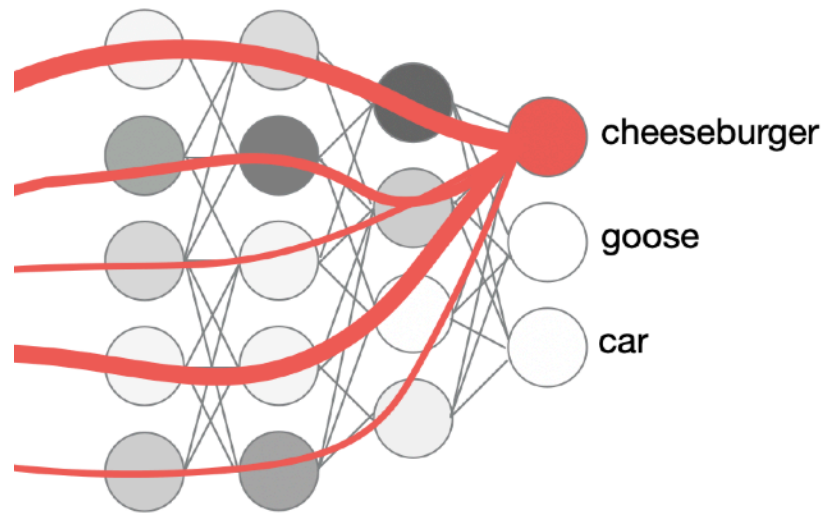


$$\max_{x \in \mathcal{X}} p_{\theta}(\omega_c | x) + \lambda \Omega(x)$$

Data-Based Activation Maximization

Find training samples, which maximally activate (output) neuron.

Most relevant training samples



(Chen et al., 2020) data-based activation maximization

Explainability 2.0: Where, What and How

nature machine intelligence



Article


<https://doi.org/10.1038/s42256-023-00711-8>

From attribution maps to human-understandable explanations through Concept Relevance Propagation

Received: 7 June 2022

Accepted: 31 July 2023

Published online: 20 September 2023

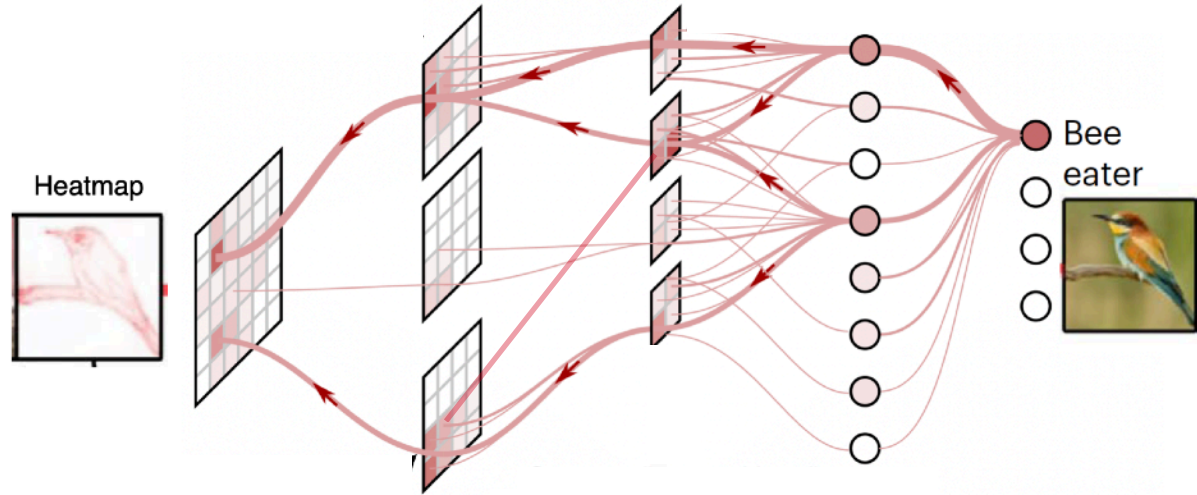
 Check for updates

Reduan Achtibat^{1,4}, Maximilian Dreyer^{1,4}, Ilona Eisenbraun¹, Sebastian Bosse¹,
Thomas Wiegand^{1,2,3}, Wojciech Samek^{1,2,3}✉ & Sebastian Lapuschkin¹✉

<https://doi.org/10.1038/s42256-023-00711-8>

Explainability 2.0: Where, What and How

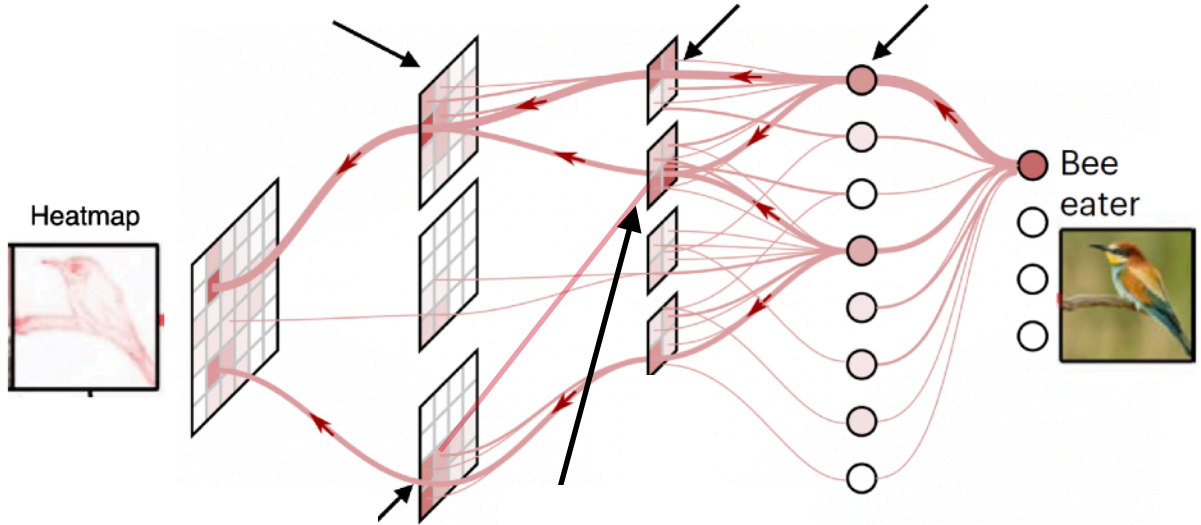
Known: Hidden layers encode semantic concepts.



Goal: Explain in terms of these concepts.

Explainability 2.0: Where, What and How

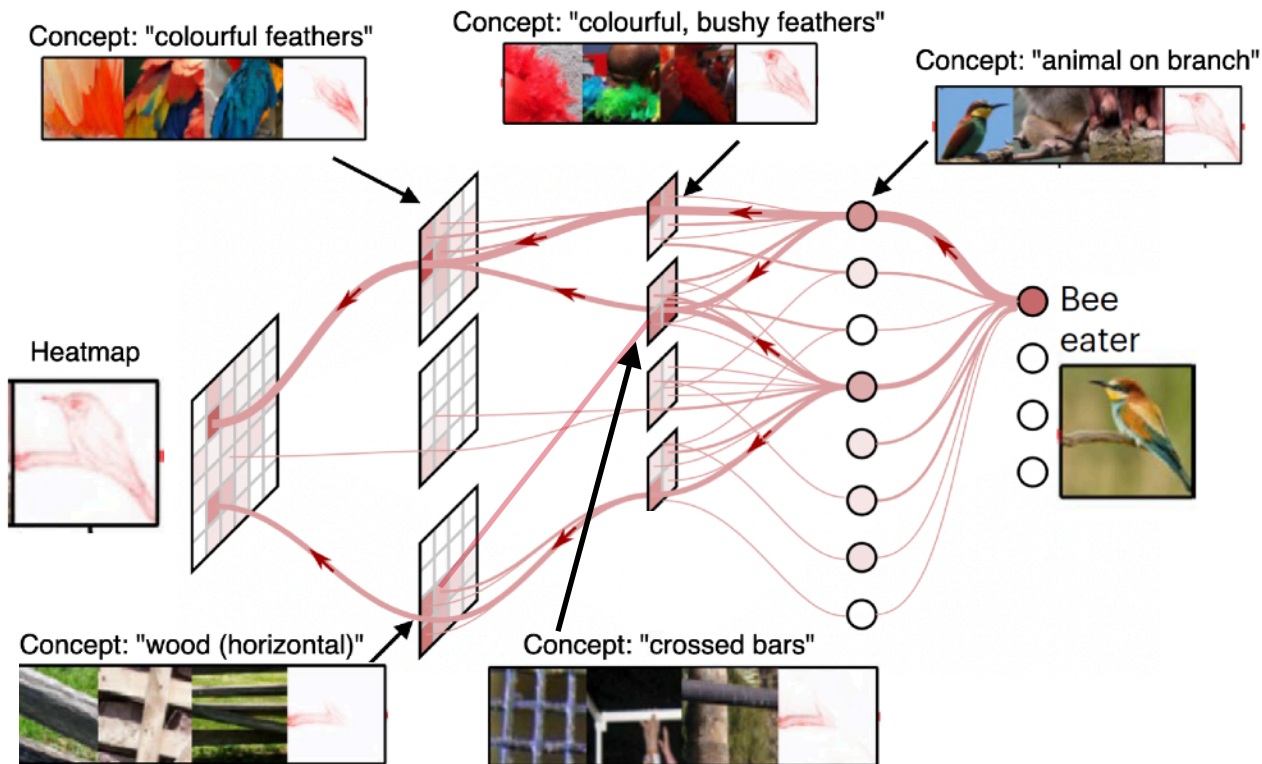
Known: Hidden layers encode semantic concepts.



Which neurons are relevant ?
—> LRP

Explainability 2.0: Where, What and How

Known: Hidden layers encode semantic concepts.



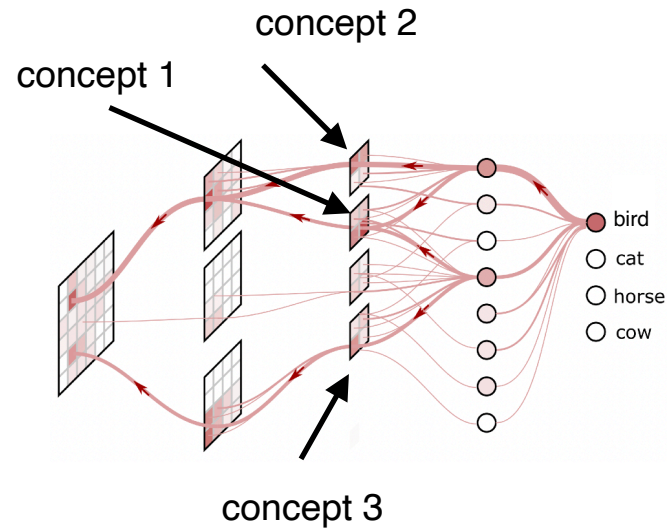
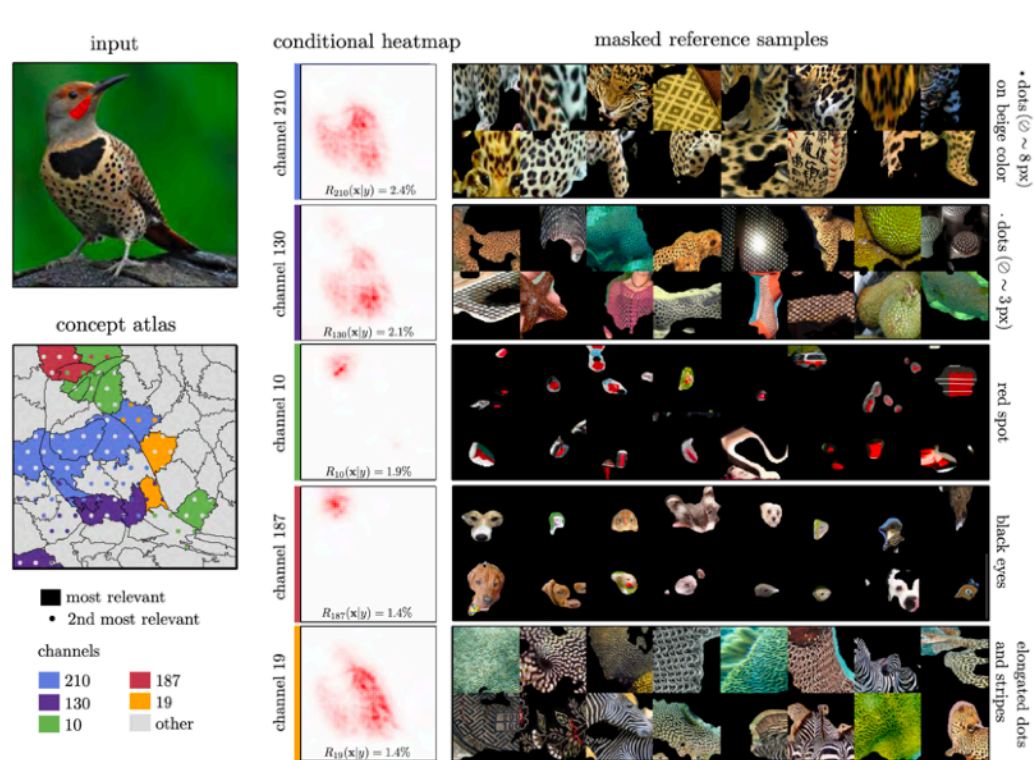
Which neurons are relevant ?

—> LRP

What are they encoding ?

—> Activation Maximization

Concept Relevance Propagation (CRP)

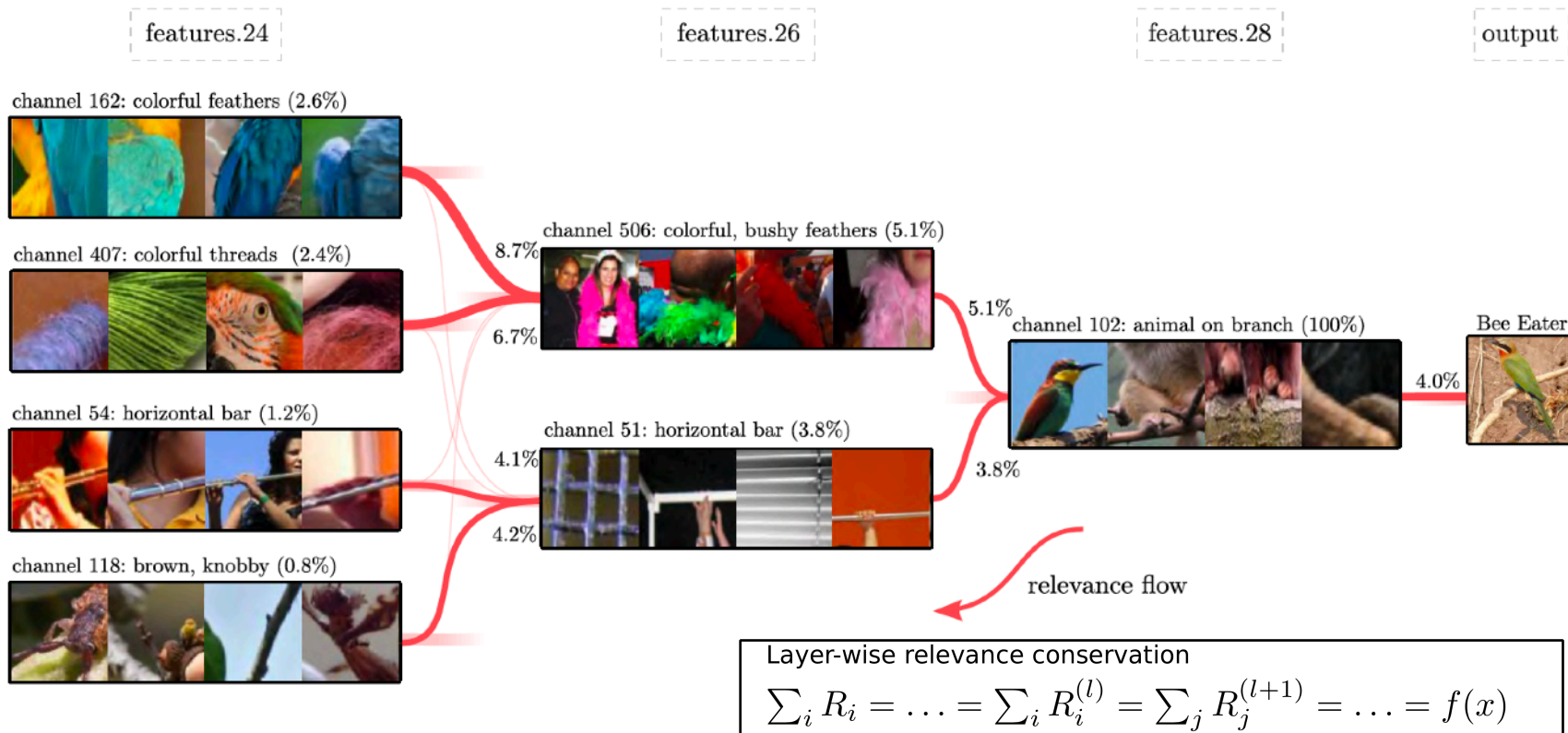


Step 1: Find relevant concepts

Step 2: Compute conditional explanation (*where*)

Step 3: Visualize relevant samples (*what*)

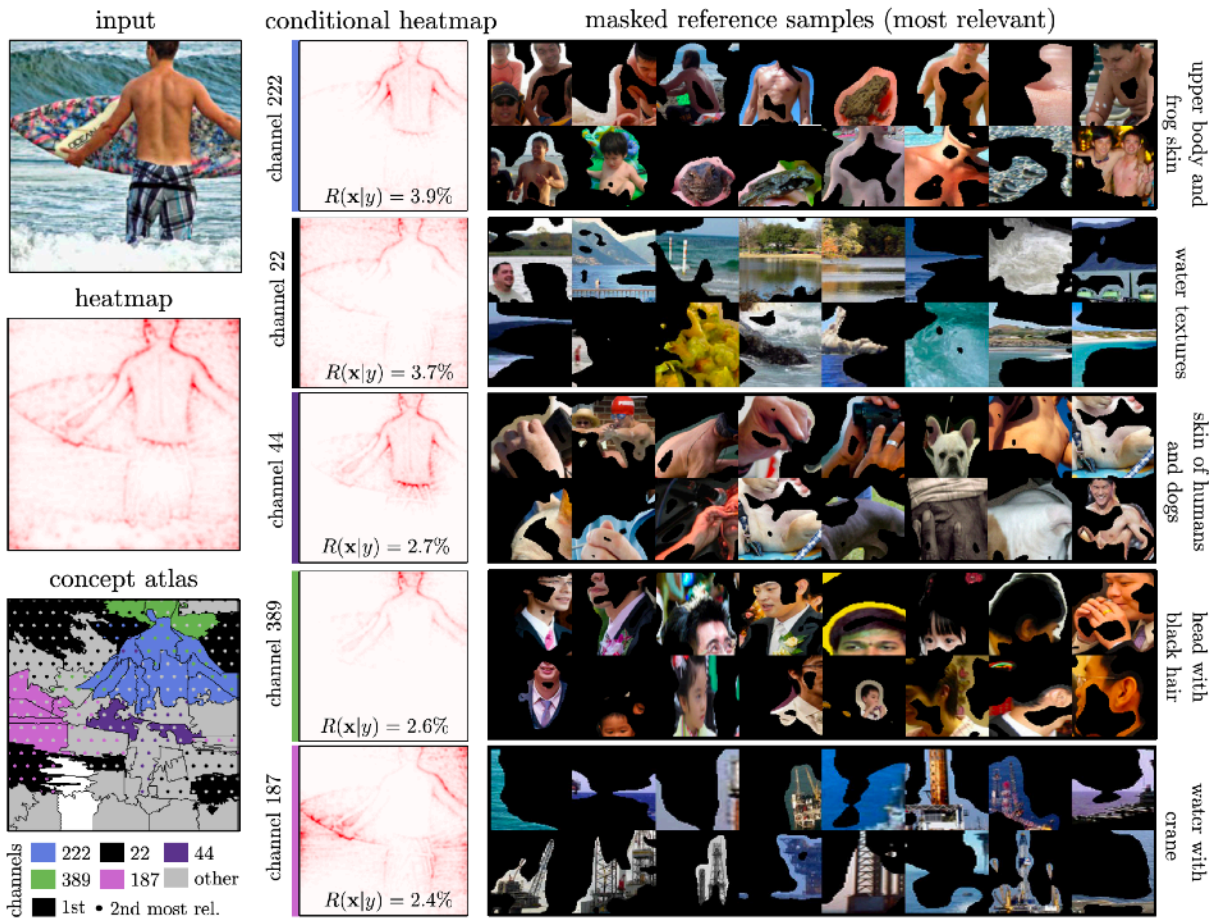
Concept Composition



Identifying Clever Hans

Prediction: swimming trunk

Relevant concepts: skin,
body, hair, water



Third Wave of XAI: "Understand Everything"

Mechanistic understanding and validation of large AI models with SemanticLens

Maximilian Dreyer^{1*} Jim Berend^{1*} Tobias Labarta¹ Johanna Vielhaben¹

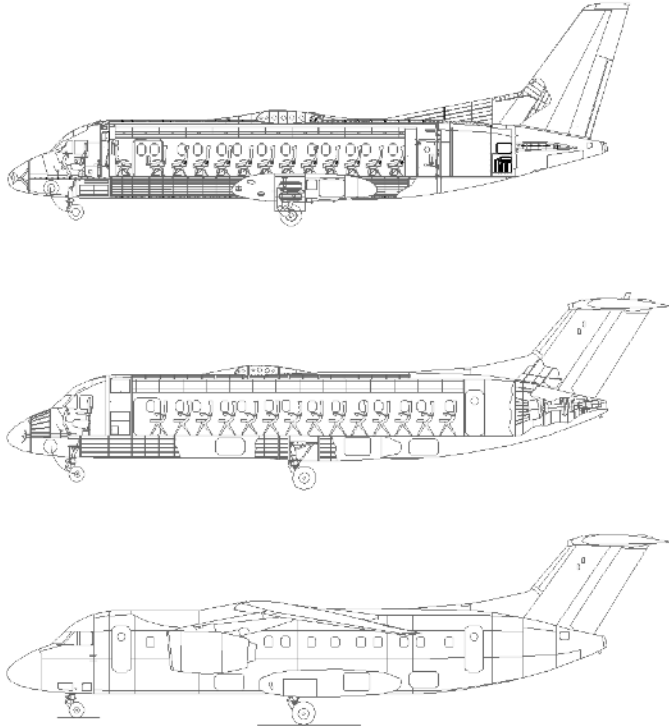
Thomas Wiegand^{1,2,3} Sebastian Lapuschkin¹ Wojciech Samek^{1,2,3}

¹Fraunhofer Heinrich Hertz Institute ²Technische Universität Berlin

³BIFOLD – Berlin Institute for the Foundations of Learning and Data

`{wojciech.samek,sebastian.lapuschkin}@hhi.fraunhofer.de`

Technical systems designed by humans

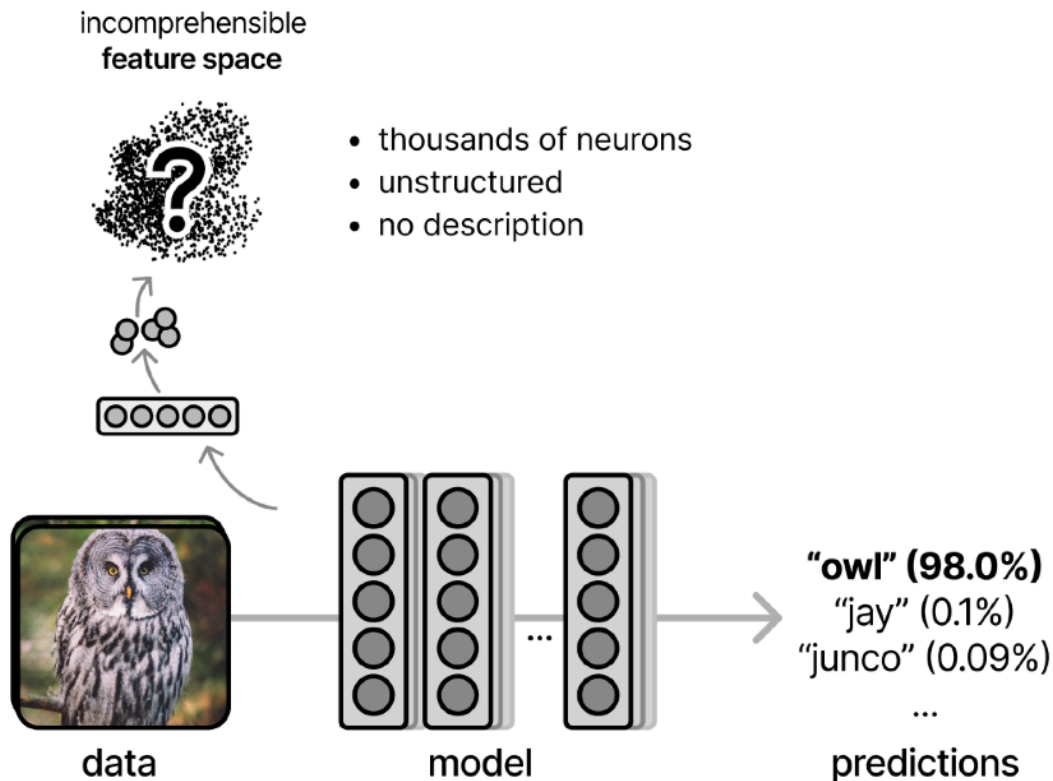


Technical systems designed by humans

- constructed step by step
- modular
- each component serving a specific, well-understood function
- can be validated and certified

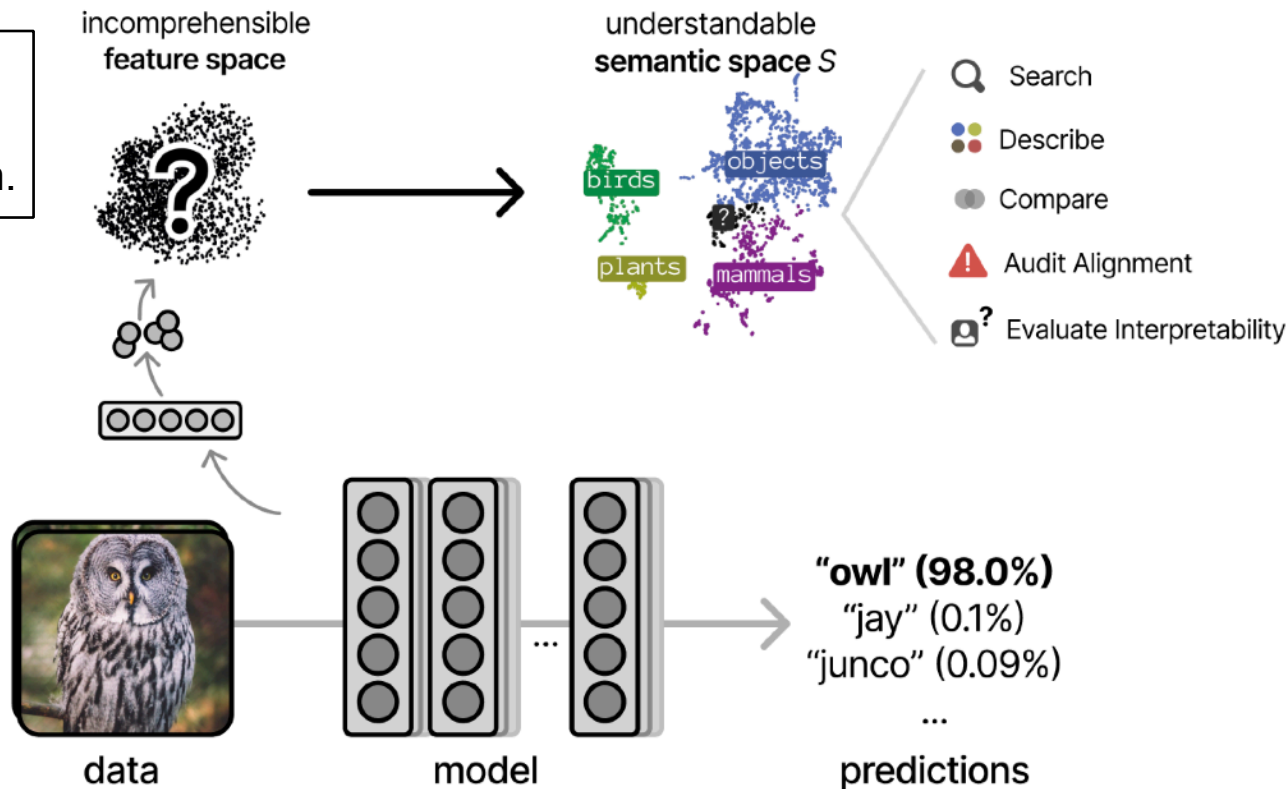
...

What Happens Inside the Model?





What Happens Inside the Model?

Idea: Transform model into comprehensible form.



SemanticLens

1 Describing the Role of Components

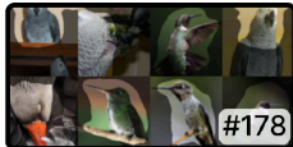
component  \longrightarrow *concept examples* 

By collecting highly activating data samples + CRP.

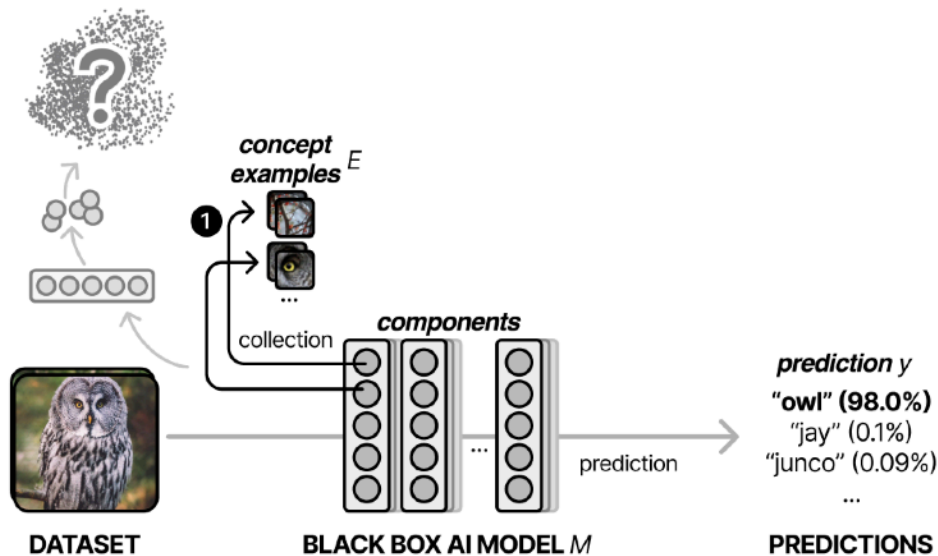
concept examples



concept examples





incomprehensible
feature space



SemanticLens

① Describing the Role of Components

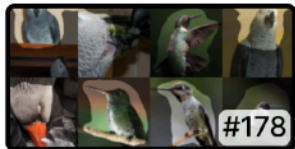
component  \longrightarrow *concept examples* 

By collecting highly activating data samples + CRP.

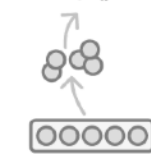
concept examples



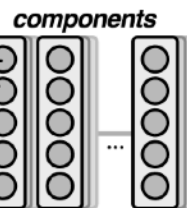
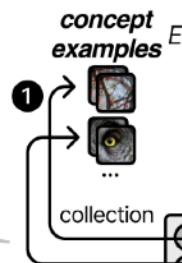
concept examples



incomprehensible
feature space



DATASET

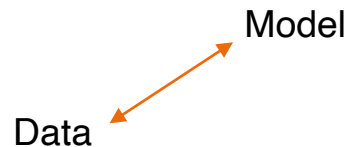


BLACK BOX AI MODEL M

prediction



prediction y
"owl" (98.0%)
"jay" (0.1%)
"junco" (0.09%)
...

PREDICTIONS


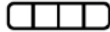


SemanticLens

1 Describing the Role of Components

component  \longrightarrow *concept examples* 

2 Semantic Embedding

concept examples  \xrightarrow{F} *semantic embeddings* 

Via Foundation Model as *Semantic Domain Expert*:

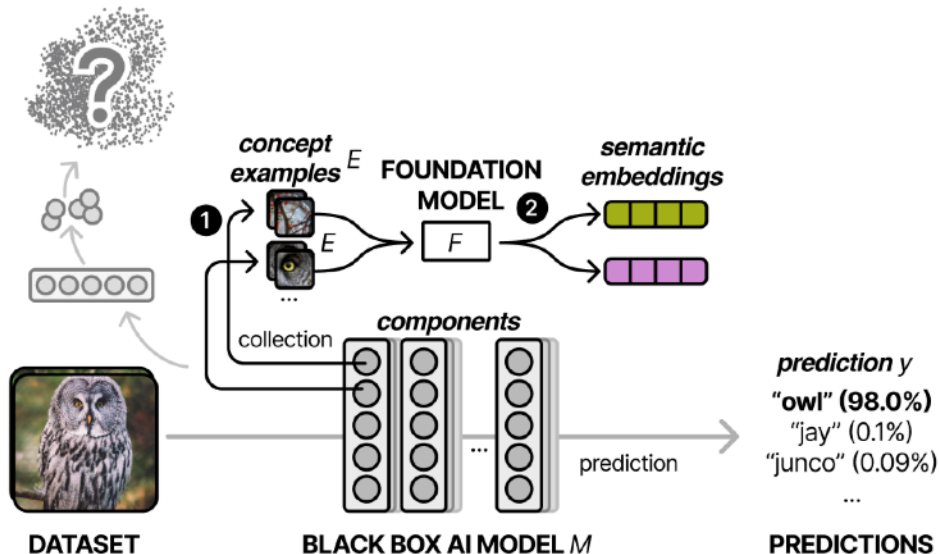
General Knowledge:

- CLIP
- Florence
- ...

Medical Domain:

- WhyLesion-CLIP
- CXR-CLIP
- ...



incomprehensible
feature space




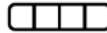
No need for human in the loop anymore

SemanticLens

1 Describing the Role of Components

component  \longrightarrow concept examples 

2 Semantic Embedding

concept examples  \xrightarrow{F} semantic embeddings 

Via Foundation Model as *Semantic Domain Expert*:

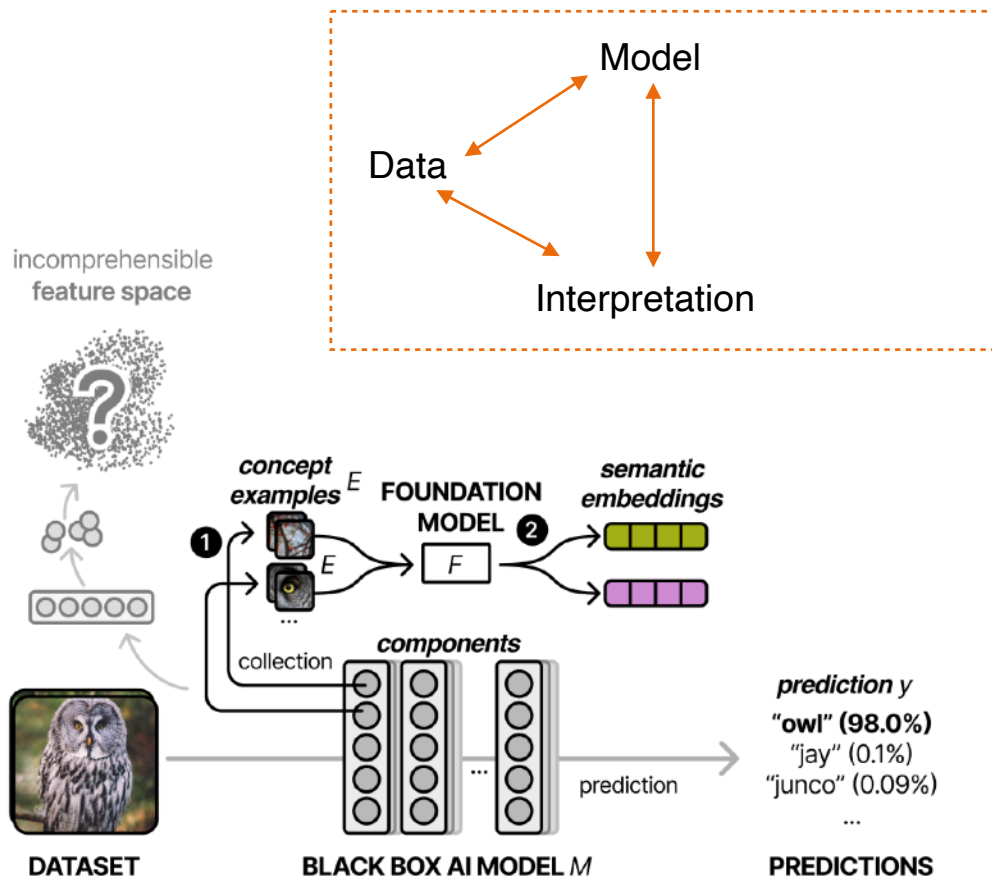
General Knowledge:

- CLIP
- Florence
- ...

Medical Domain:

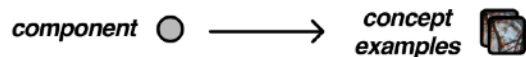
- WhyLesion-CLIP
- CXR-CLIP
- ...

No need for human in the loop anymore

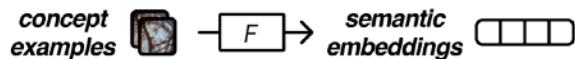


SemanticLens

1 Describing the Role of Components



2 Semantic Embedding



3 Connect with Concept Relevance



Retrieve component-level relevance scores with CRP for:

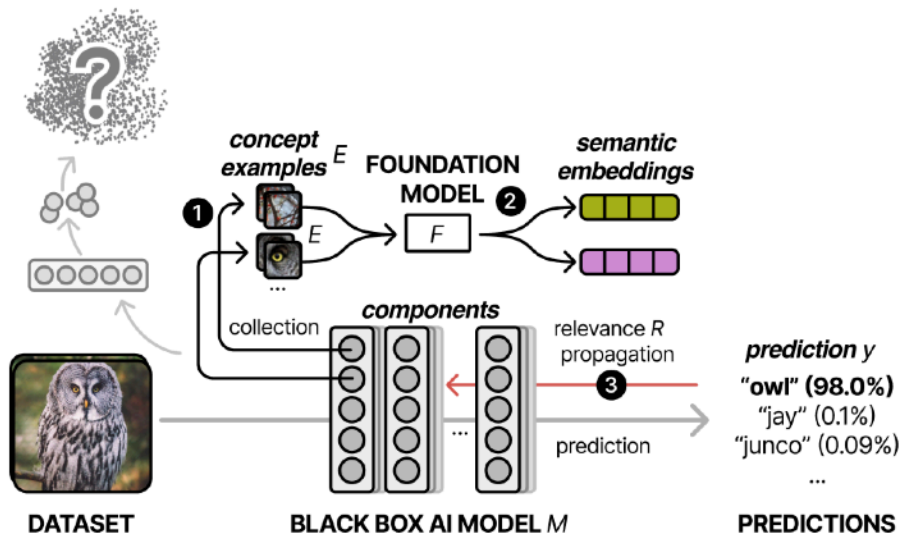
- output predictions



- upper-level components

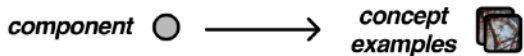


incomprehensible
feature space

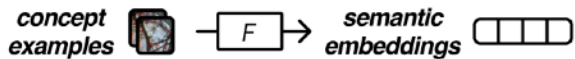


SemanticLens

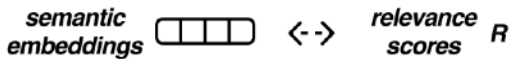
1 Describing the Role of Components



② Semantic Embedding



3 Connect with Concept Relevance

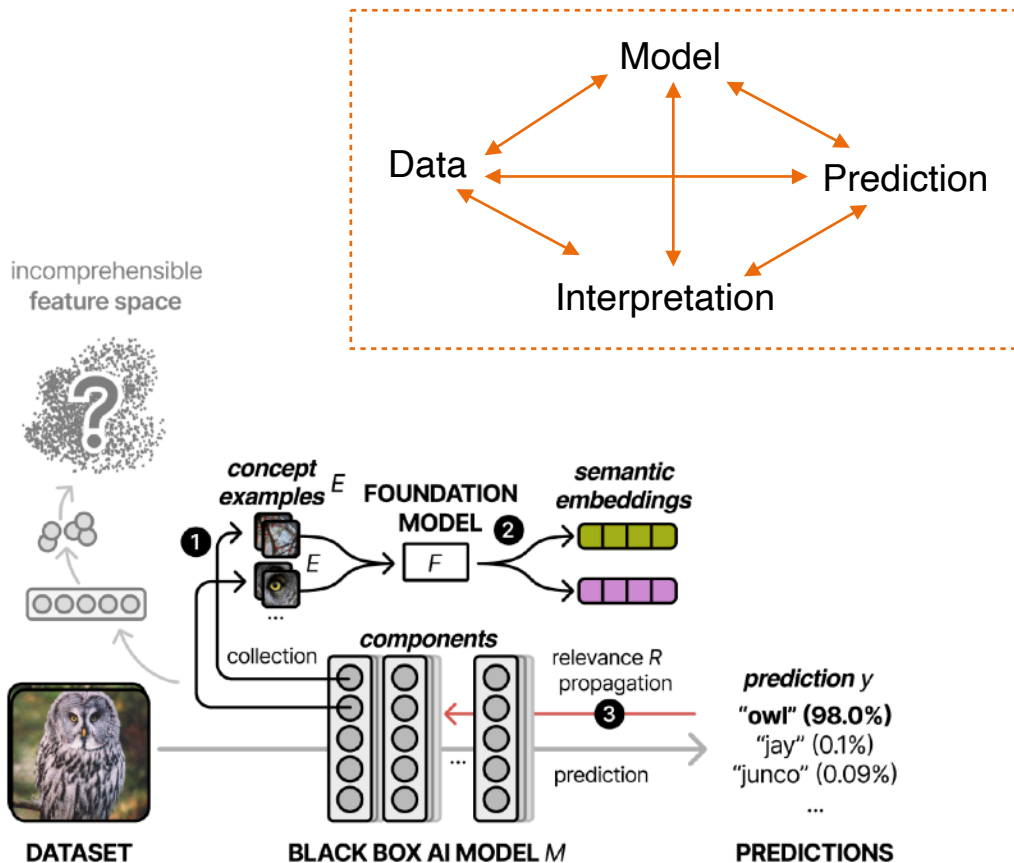


Retrieve component-level relevance scores with *CRP* for:

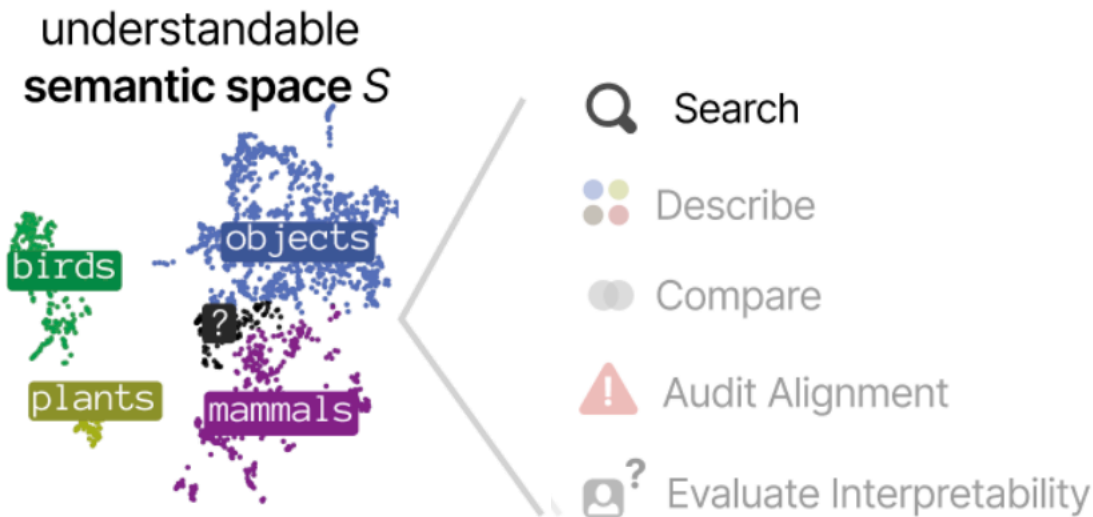
- output predictions



- upper-level components



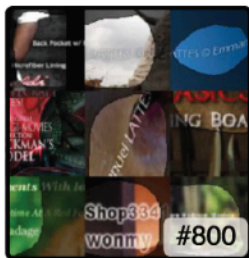
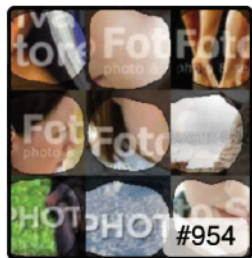
SemanticLens: What Can We do ?



Search: Finding the Needle in the Haystack

find artefact-related neurons

Q “watermark”

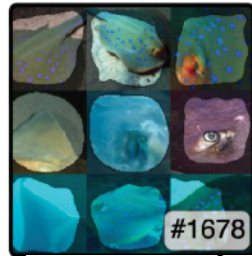
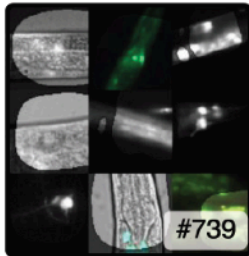


...

concept examples \mathcal{E}
of most aligned
semantic embeddings

find specific knowledge-related neurons

Q “bioluminescence”

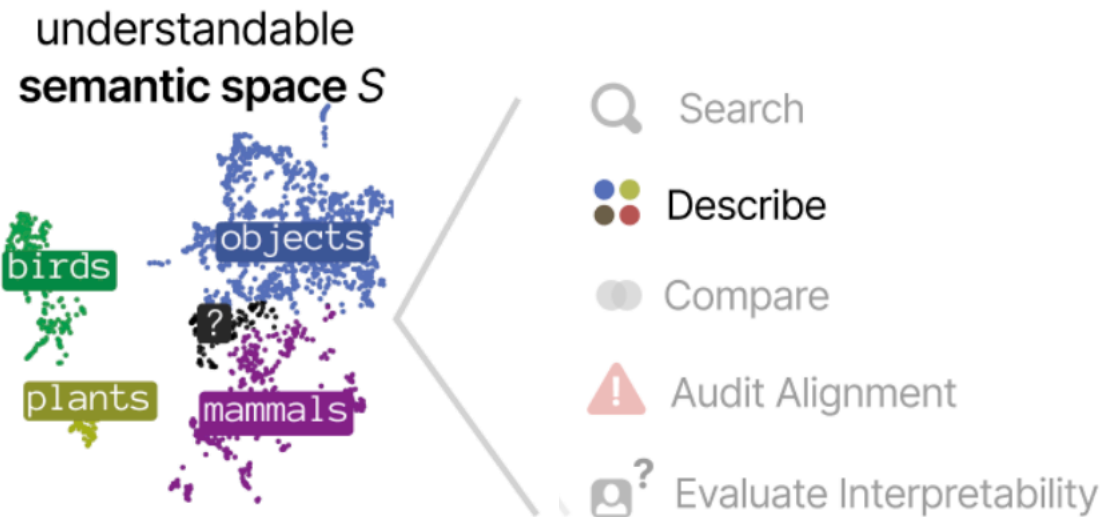


...

neuron number

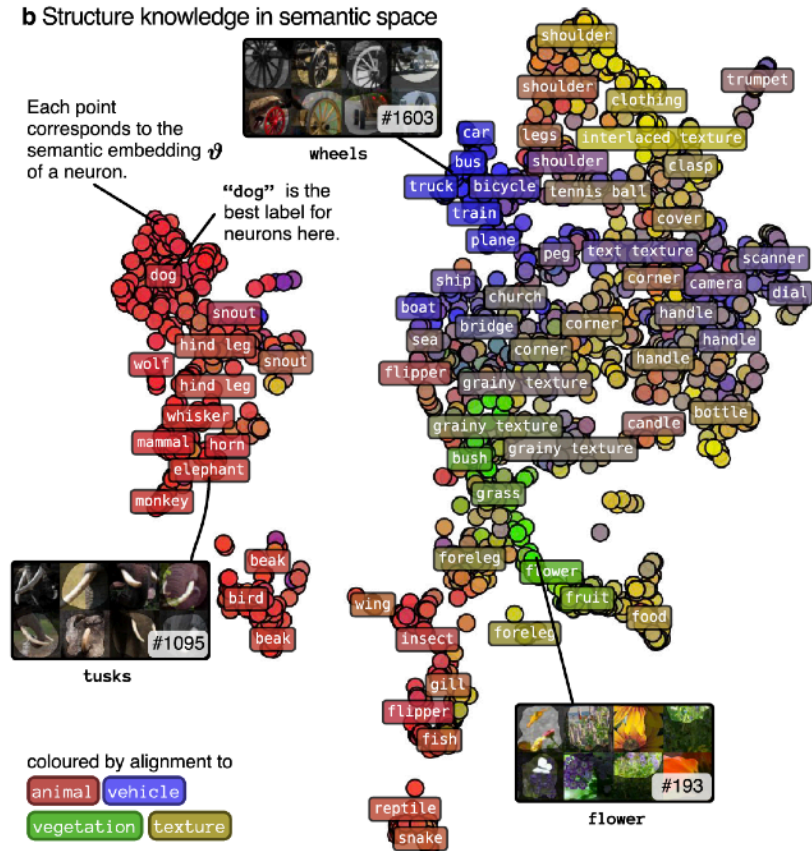
Note: CLIP models allow to measure similarity between image embeddings (here: neuron) and text embeddings (here: query).

SemanticLens: What Can We do ?

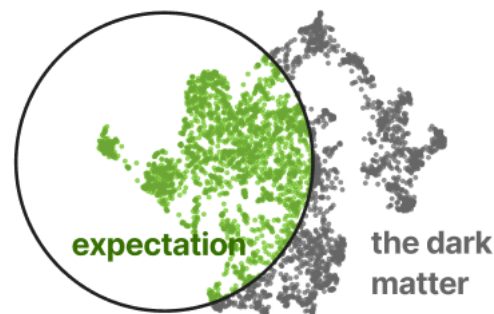


Describe: What Knowledge (does not) Exists ?

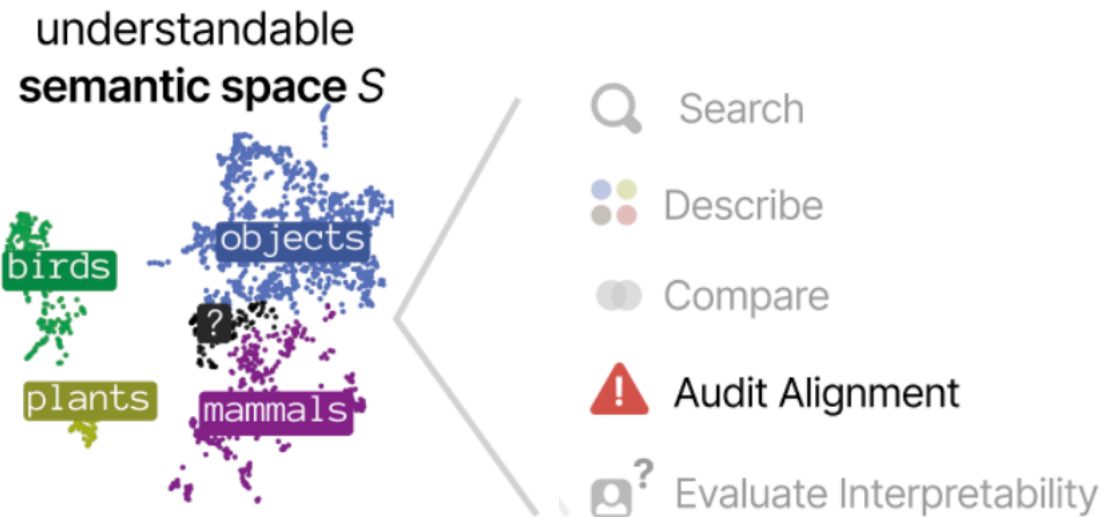
b Structure knowledge in semantic space



Also here we measure similarity between image embeddings (here: neuron) and text embeddings (here: label from a vocabulary of labels).



SemanticLens: What Can We do ?



A Tool for Auditing

1 define concepts for detecting Ox

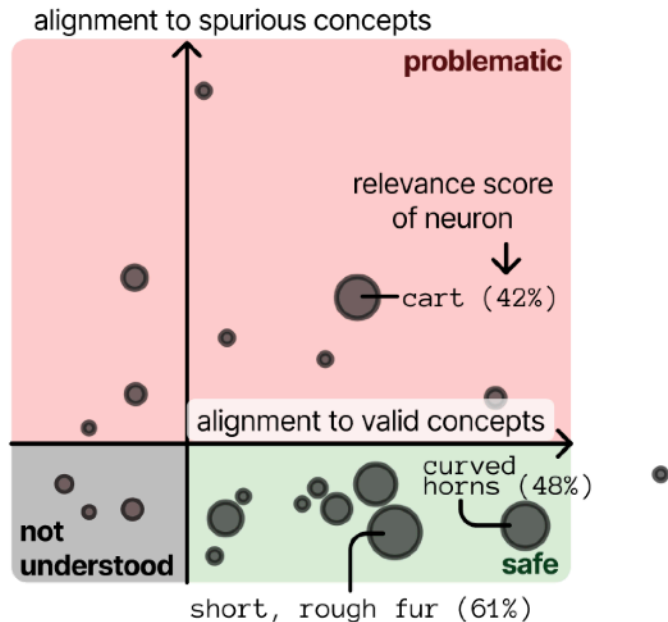
valid concepts:

large muscular body
curved horns
hooves
thick neck
short, rough fur
soft fur
long fur
brown coat
black coat
white coat
strong legs
long tail
wide muzzle

spurious concepts:

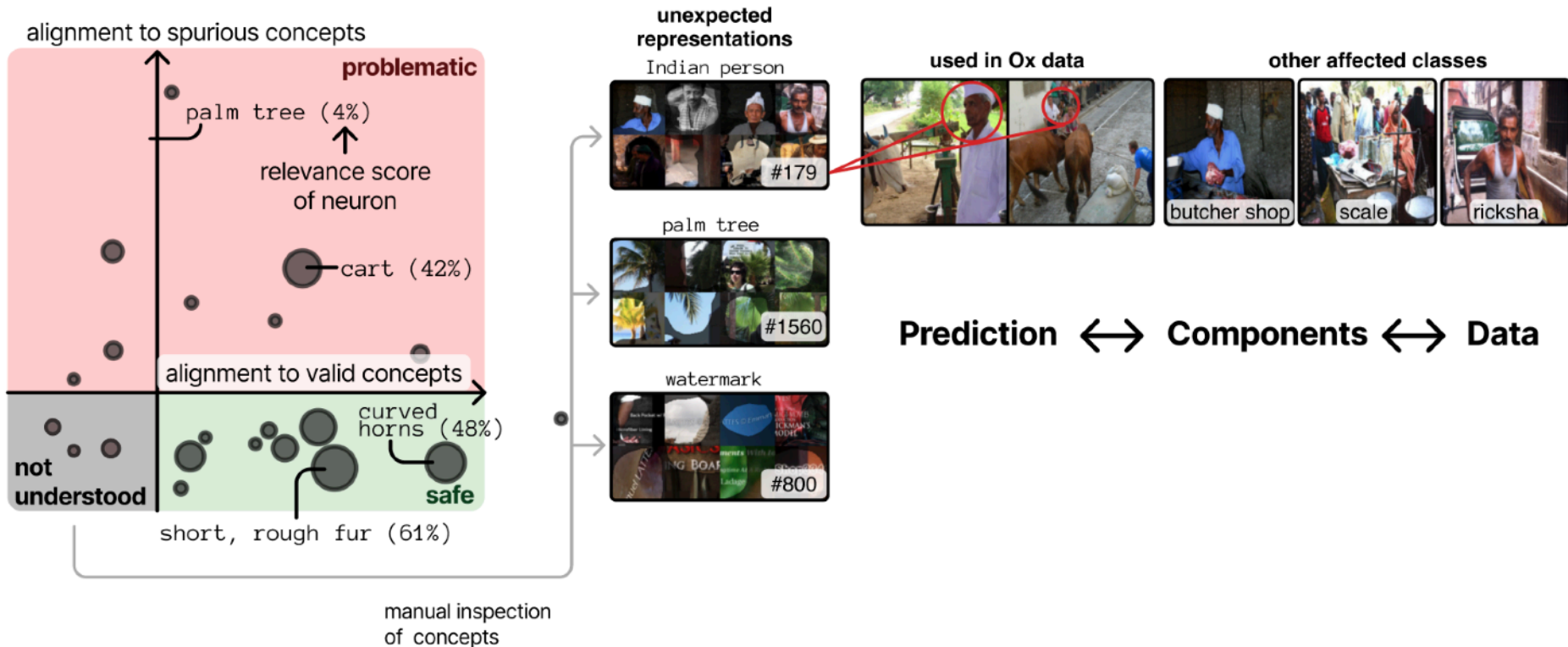
grassland
sky
tree
water
grain, straw
cart
wheel
mud, dirt
person
wooden

2 evaluate alignment of neurons



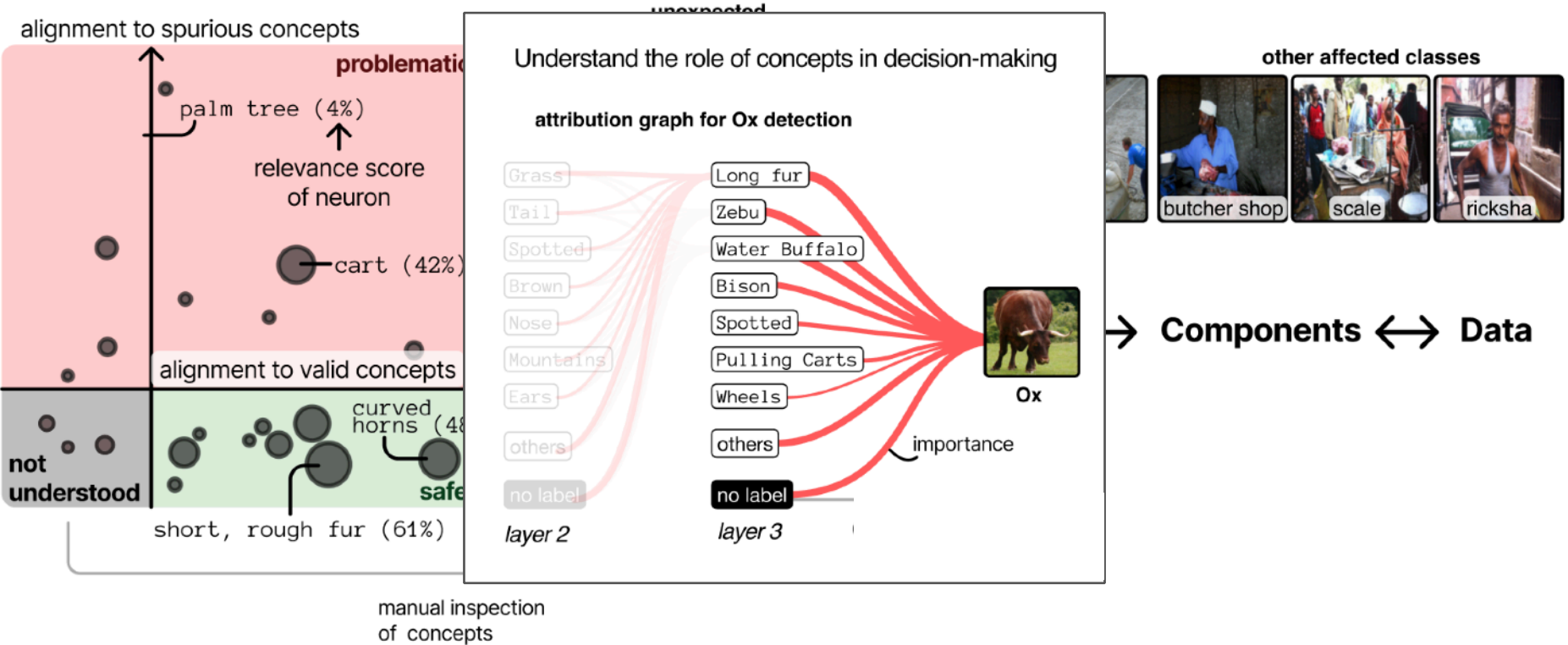
Note: Size of circle shows the “relevance” of this concept.

A Tool for Auditing

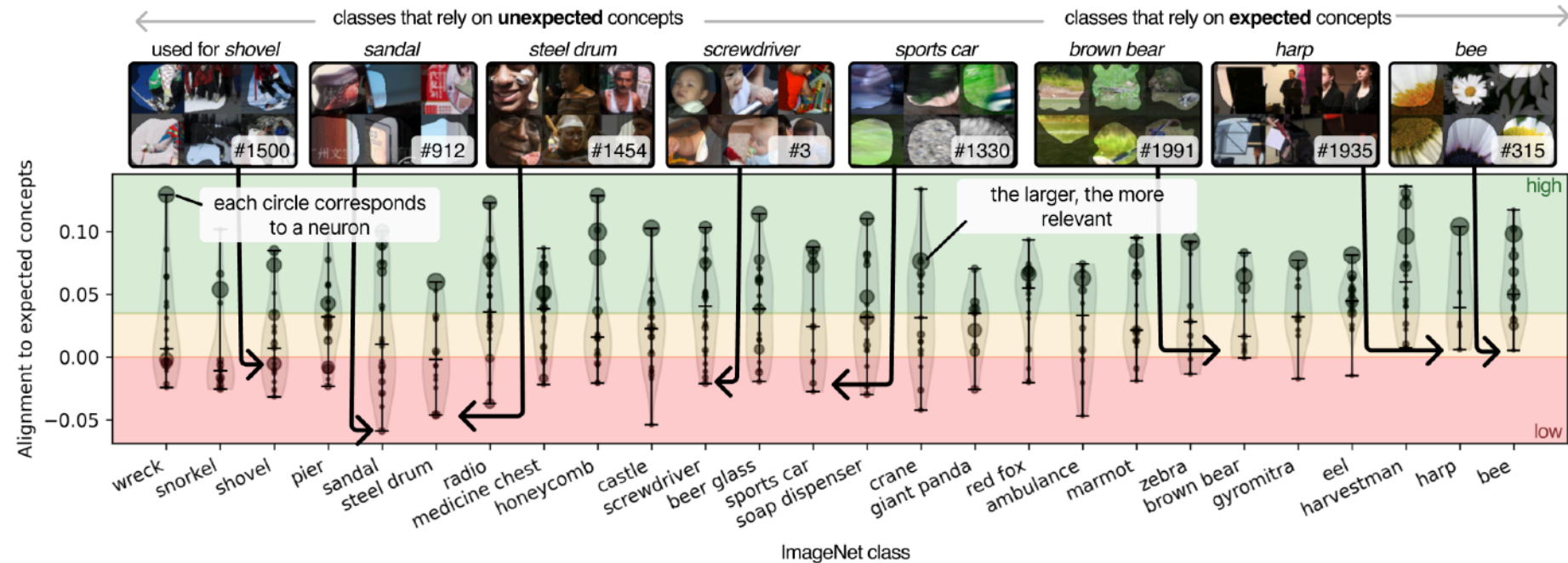


For concepts which we do not understand (i.e., dark matter) we can go back to data for manual inspection.

A Tool for Auditing

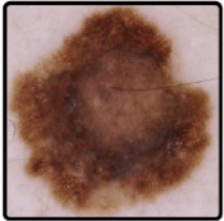


A Tool for Auditing

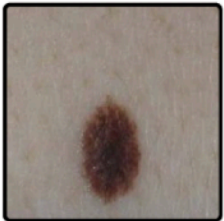


Audit Alignment: Medical Case

a Defining concepts for melanoma detection: ABCDE rule



Melanoma



Other (Regular)

valid concepts:

A. Asymmetry

- asymmetric lesion, ...

B. Border

- ragged border, ...

C. Color

- blue-white veil, ...

D. Diameter

- large lesion, ...

E. Evolving

- crusty surface, ...

Other (Regular)

- even border, ...

Other (Irregular)

- white or yellowish structures, ...

spurious concepts:

hairs

hairy

band-aid

blue-coloured band-aid

red skin

measurement scale bar

ruler

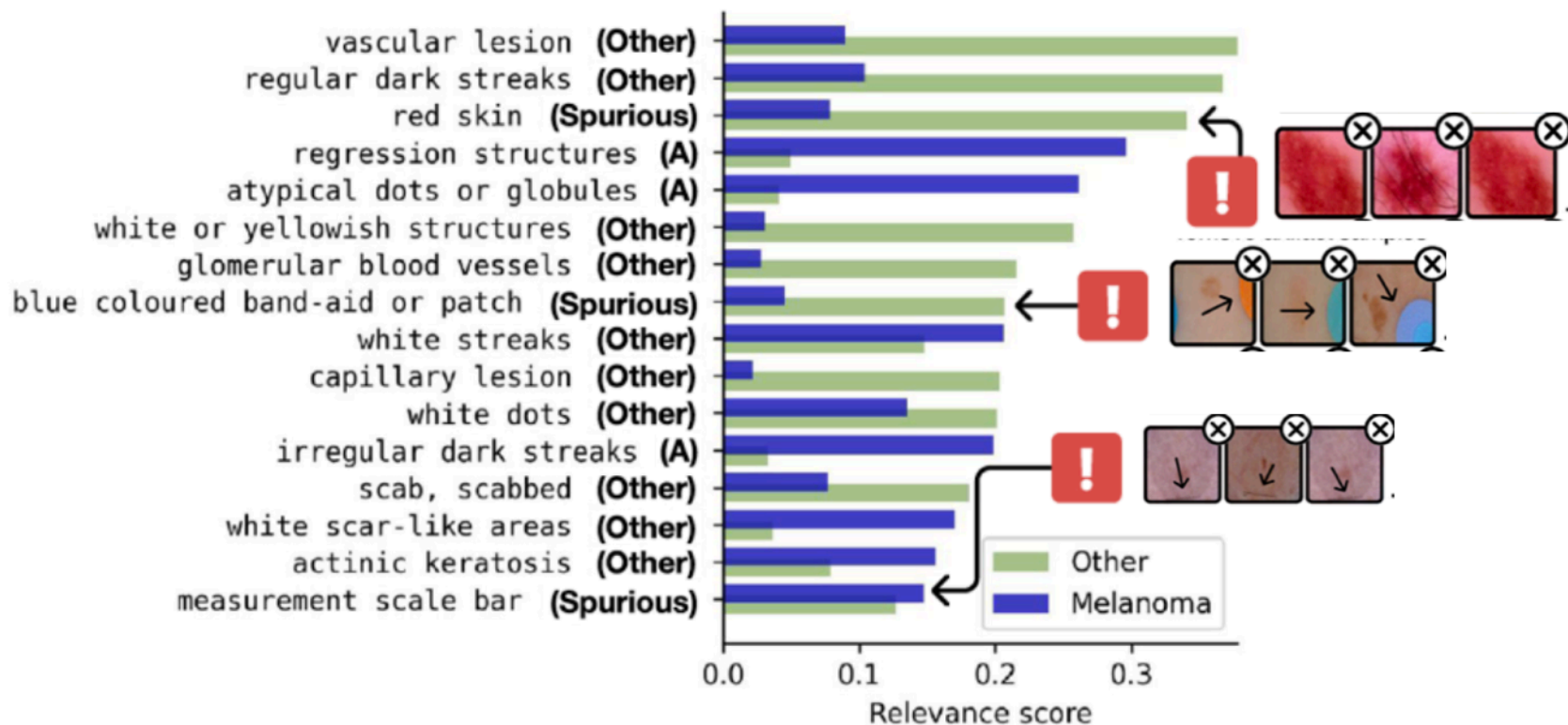
vignetting

purple ink

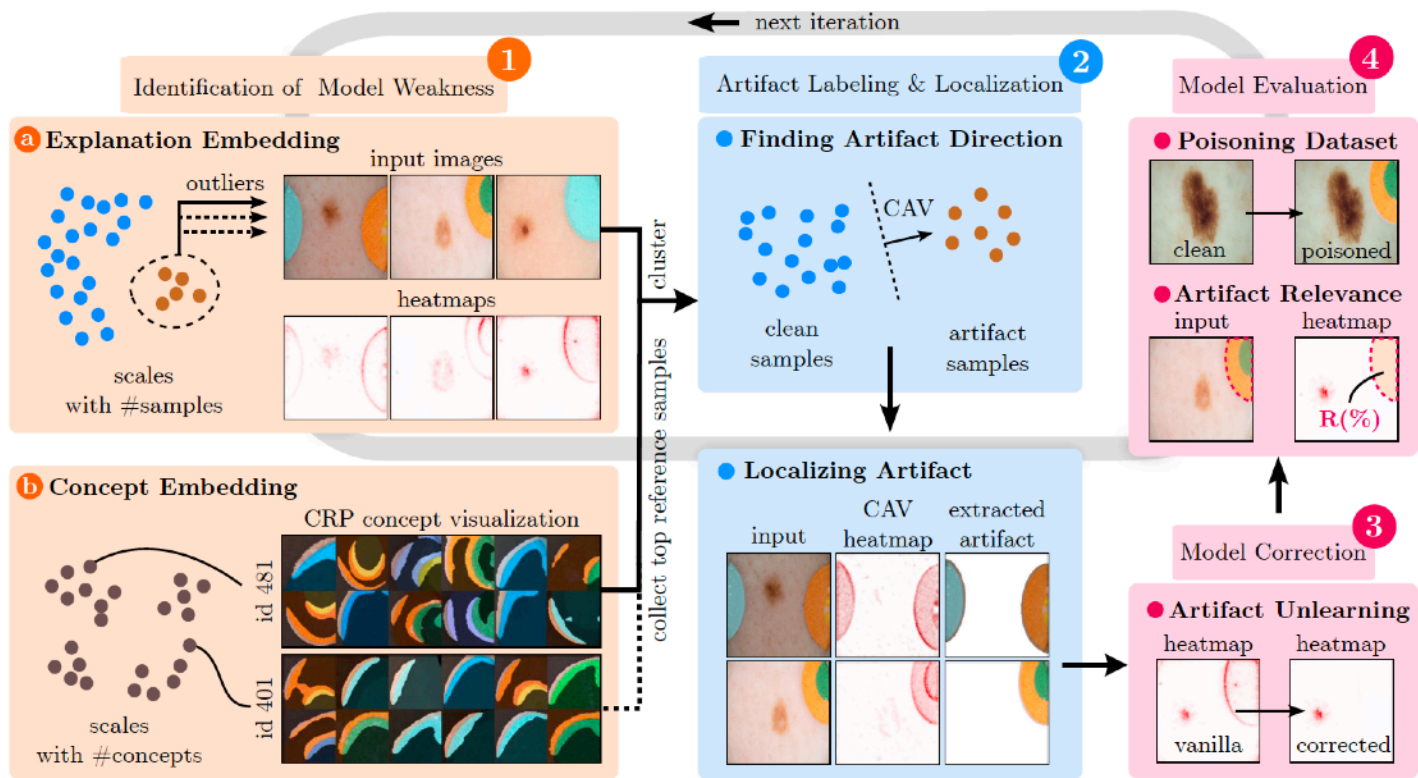
skin marker

...

Audit Alignment: Medical Case

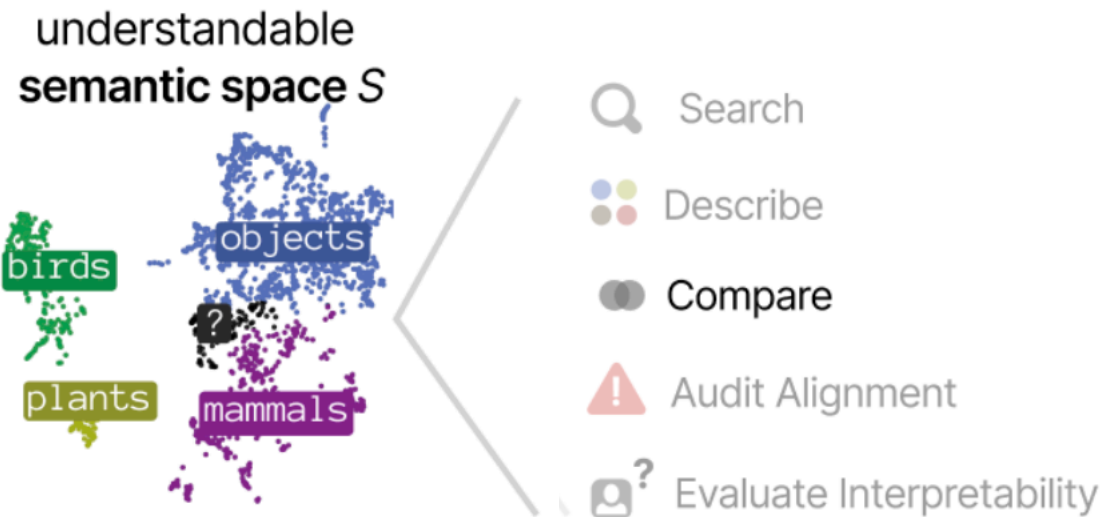


From Inspecting to Debugging



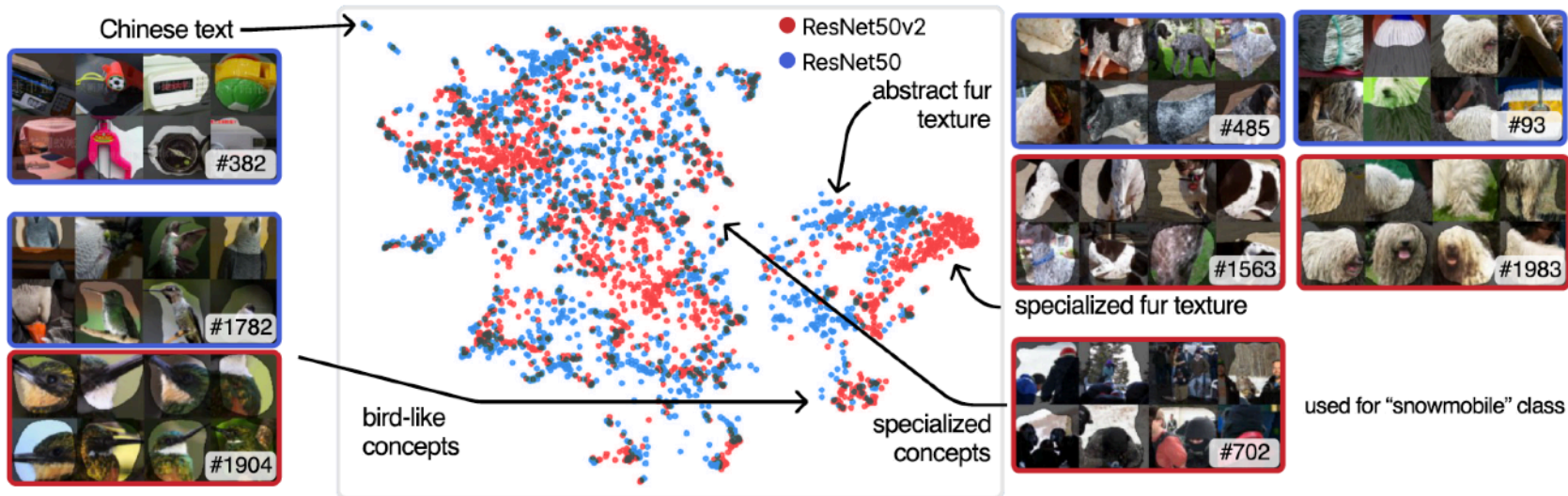
[Pahde et al. 2023]

SemanticLens: What Can We do ?



Compare: Identify Common and Unique Knowledge

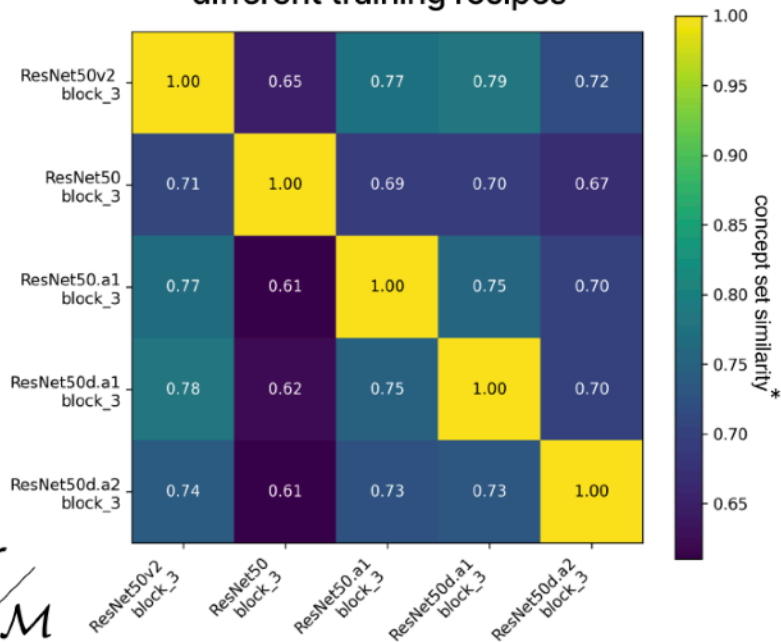
What concepts are shared between two models, and which are unique to each one?



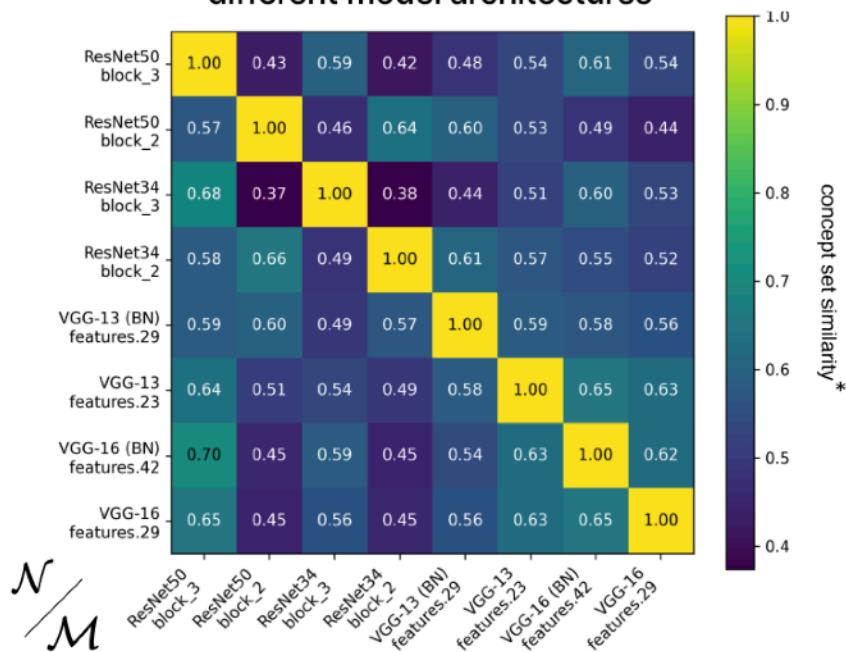
Note: Comparison can be done because components of different models maps into the same semantic space.

Compare: Identify Common and Unique Knowledge

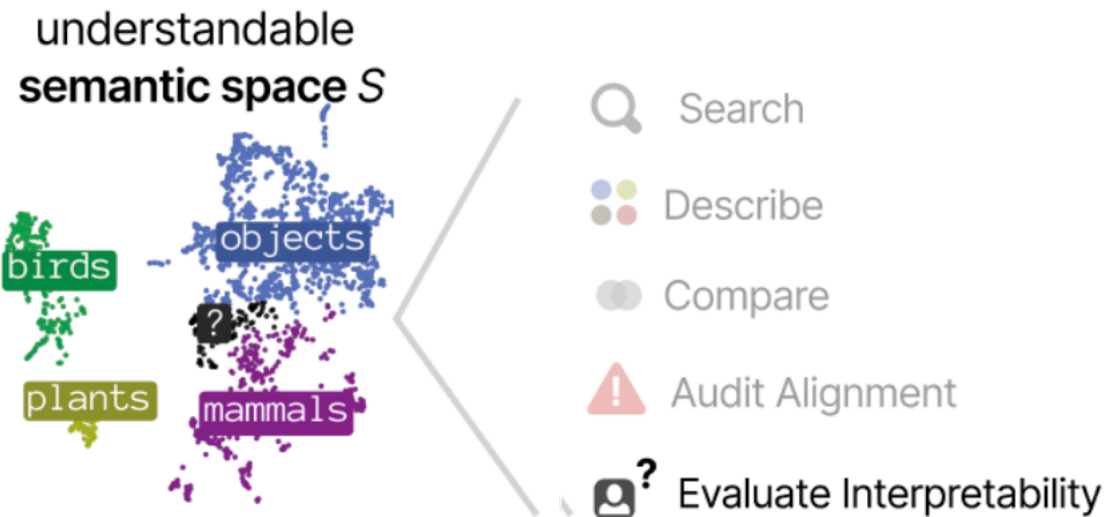
different training recipes



different model architectures



SemanticLens: What Can We do ?

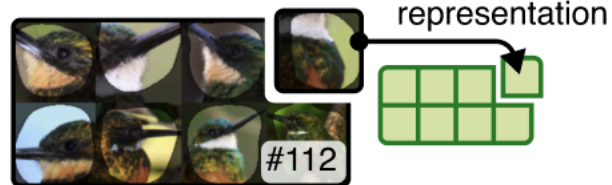


Evaluating Component Interpretability

clarity

per concept ●

how clear is a concept?



polysemanticity

per concept ●

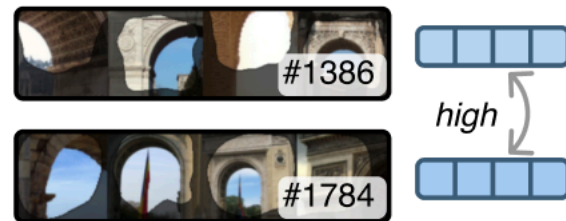
how polysemantic is a concept?



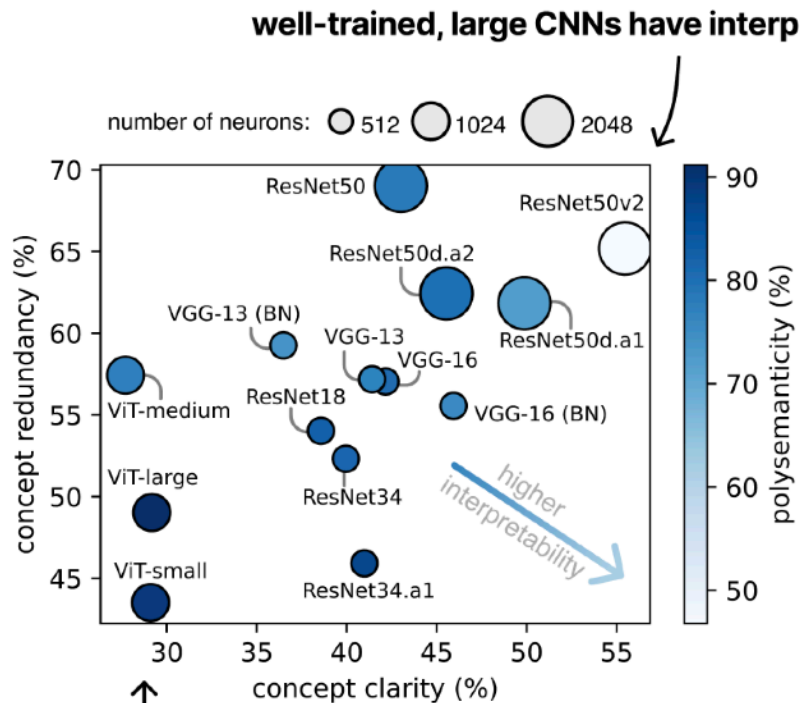
similarity

between two concepts ●●

how similar are concepts?

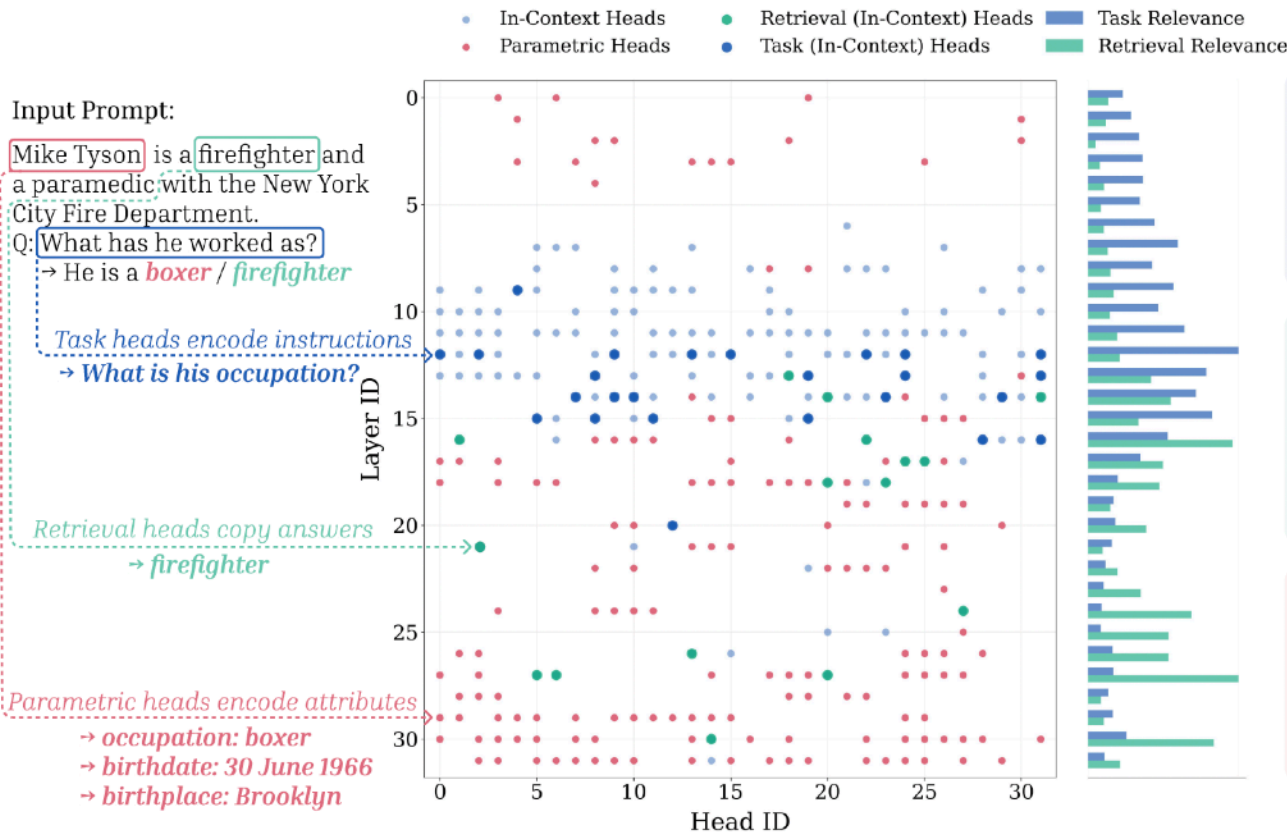


Evaluating Component Interpretability



ViTs have no inherently interpretable components

Next Steps: Component-Level Understanding of LLMs



a) Execute instructions in another prompt

What has he worked as?

patch instruction
into last token

Margaret Mitchell was born in Georgia.
 → She is a novelist

b) Change retrieved answer object



modify attention weights



Mike Tyson is a firefighter and a paramedic [...]
 → He is a paramedic

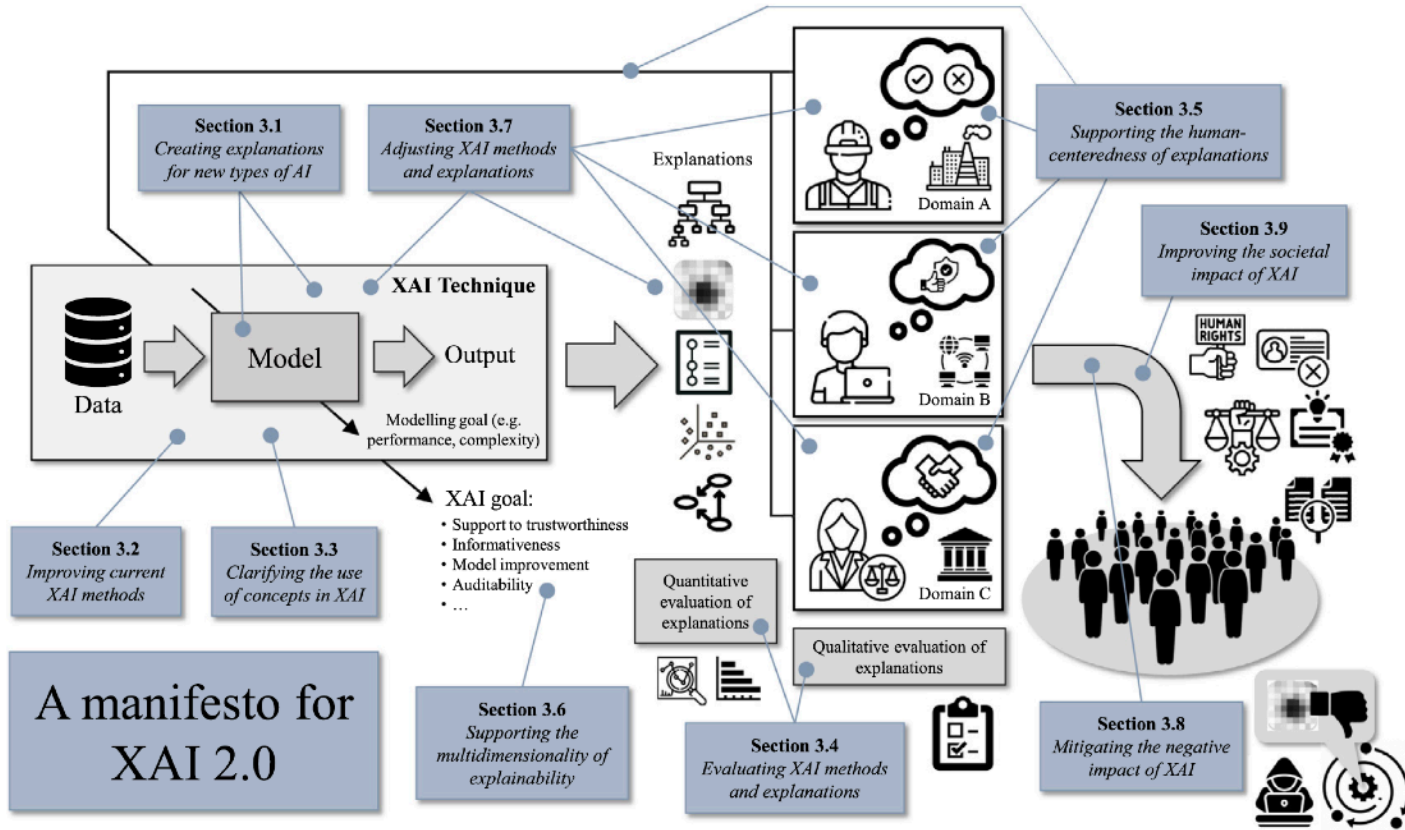
c) Overwrite entities' attributes

Mike Tyson

patch attributes
into last token

Q: What is the occupation of Albert Einstein?
 → He is a boxer

Future Work



(Longo et al. 2024)

<https://doi.org/10.1016/j.inffus.2024.102301>

Toolboxes

Benchmarking:

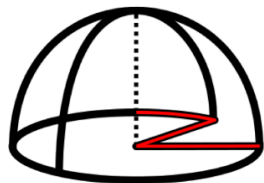
QUANTUS

<https://github.com/understandable-machine-intelligence-lab/Quantus>



SemanticLens

<https://github.com/jim-berend/semanticlens>



zennit

<https://github.com/chr5tphr/zennit>

Benchmarking:
CLEVR-XAI



<https://github.com/ahmedmagdiosman/clevr-xai>



zennit-crp

<https://github.com/rachtibat/zennit-crp>

iNNvestigate

<https://github.com/albermax/innvestigate>



quanda

<https://github.com/dilyabareeva/quanda>



**THANK YOU
FOR YOUR
ATTENTION**