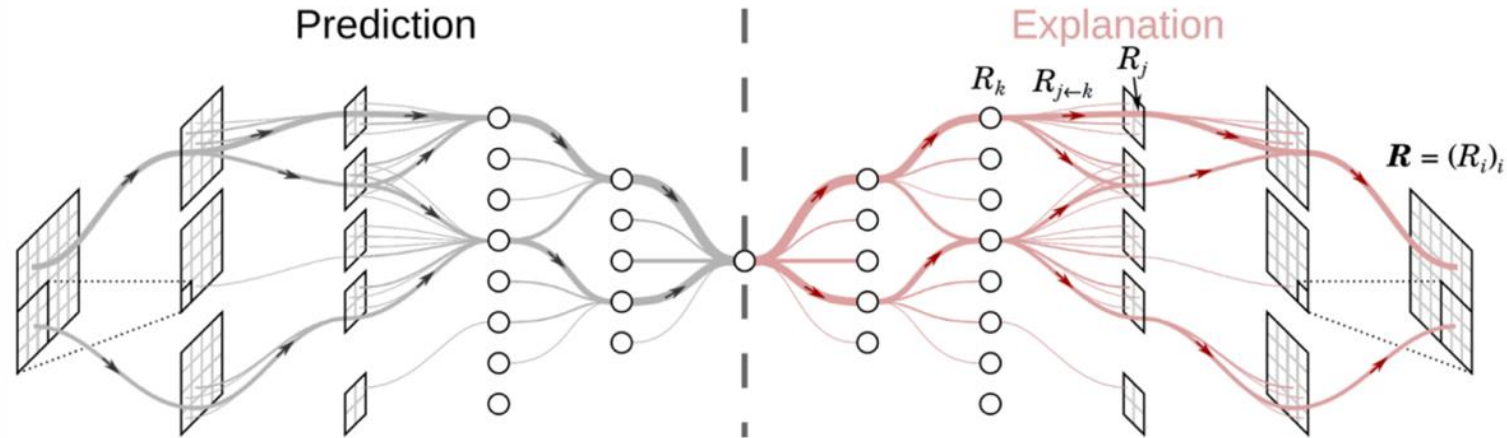


# From Feature Attributions to Next-Generation Explainable AI

Wojciech Samek

TU Berlin & Fraunhofer HHI



# Syllabus

## Part I: First Generation XAI

- What to explain
- Explaining by attribution
- DTD Framework
- Evaluating explanations

## Part II: New Developments

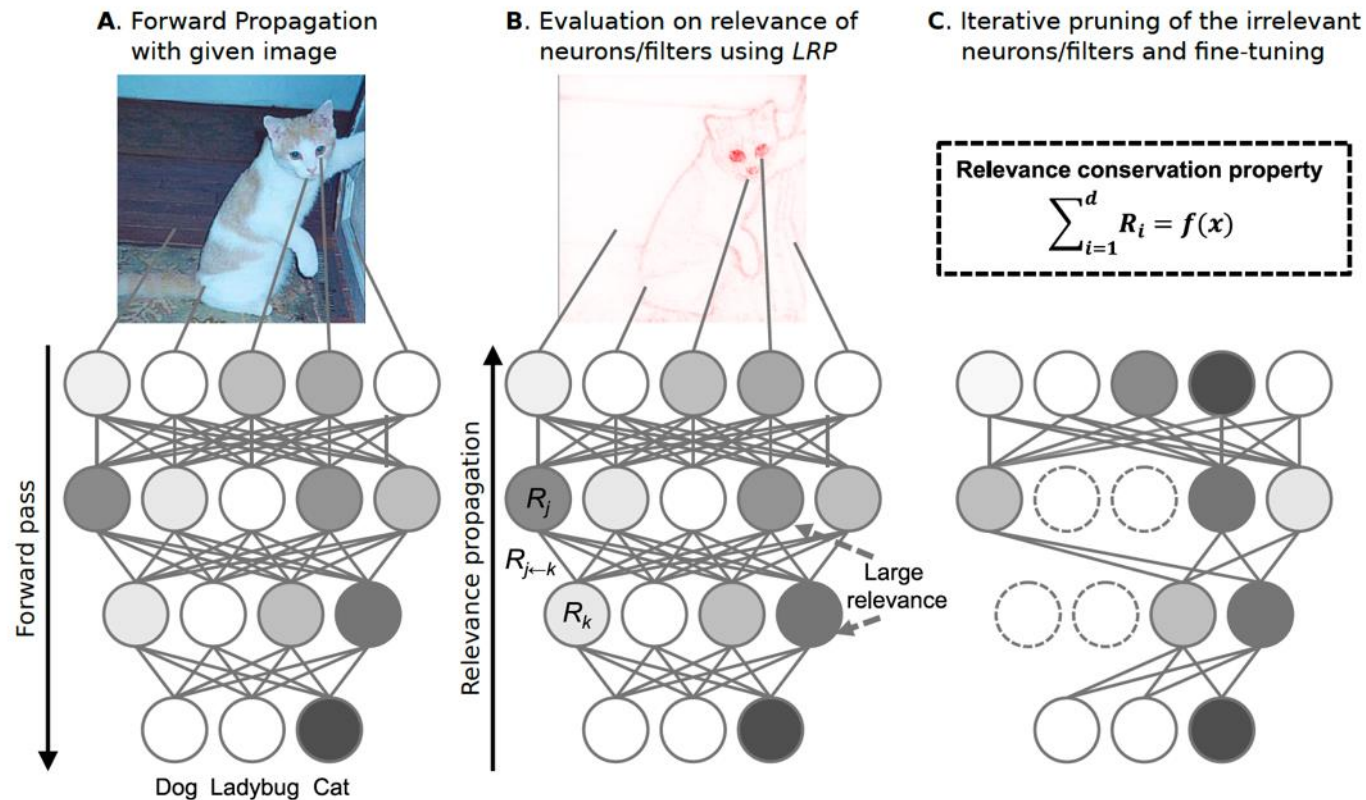
- Concepts and prototypes
- XAI for LLMs
- Non-Interpretable domains
- Beyond classification

## Part III: Beyond Explaining

- XAI-Based model surgery
- Reveal and revise
- Explanatory interactive ML
- Future of XAI

# **XAI-Based Model Surgery**

# XAI-Based Pruning



(Yeom et al. 2019)

# XAI-Based Pruning

VGG-16	Scene 15 @ 55%					Event 8 @ 55%					Cats & Dogs @ 60%				
	U	W	T	G	L	U	W	T	G	L	U	W	T	G	L
Loss	2.09	2.27	1.76	1.90	<b>1.62</b>	0.85	1.35	1.01	1.18	<b>0.83</b>	0.19	0.50	0.51	0.57	<b>0.44</b>
Accuracy	88.59	82.07	83.00	82.72	<b>83.99</b>	95.95	90.19	91.79	90.55	<b>93.29</b>	99.36	97.90	97.54	97.19	<b>98.24</b>
Params	119.61	56.17	53.10	53.01	<b>49.67</b>	119.58	56.78	48.48	50.25	<b>47.35</b>	119.55	47.47	51.19	57.27	<b>43.75</b>
FLOPs	15.50	8.03	<b>4.66</b>	4.81	6.94	15.50	8.10	5.21	<b>5.05</b>	7.57	15.50	7.02	3.86	<b>3.68</b>	6.49
	Oxford Flower 102 @ 70%					CIFAR-10 @ 30%									
	U	W	T	G	L	U	W	T	G	L					
Loss	3.69	3.83	3.27	3.54	<b>2.96</b>	1.57	1.83	1.76	1.80	<b>1.71</b>					
Accuracy	82.26	71.84	72.11	70.53	<b>74.59</b>	91.04	93.36	93.29	93.05	<b>93.42</b>					
Params	119.96	39.34	41.37	42.68	<b>37.54</b>	119.59	<b>74.55</b>	97.30	97.33	89.20					
FLOPs	15.50	5.48	<b>2.38</b>	2.45	4.50	15.50	11.70	<b>8.14</b>	8.24	9.93					

## Pruning criteria

U = Unpruned

W = Weight

T = Taylor

G = Gradient

L = LRP

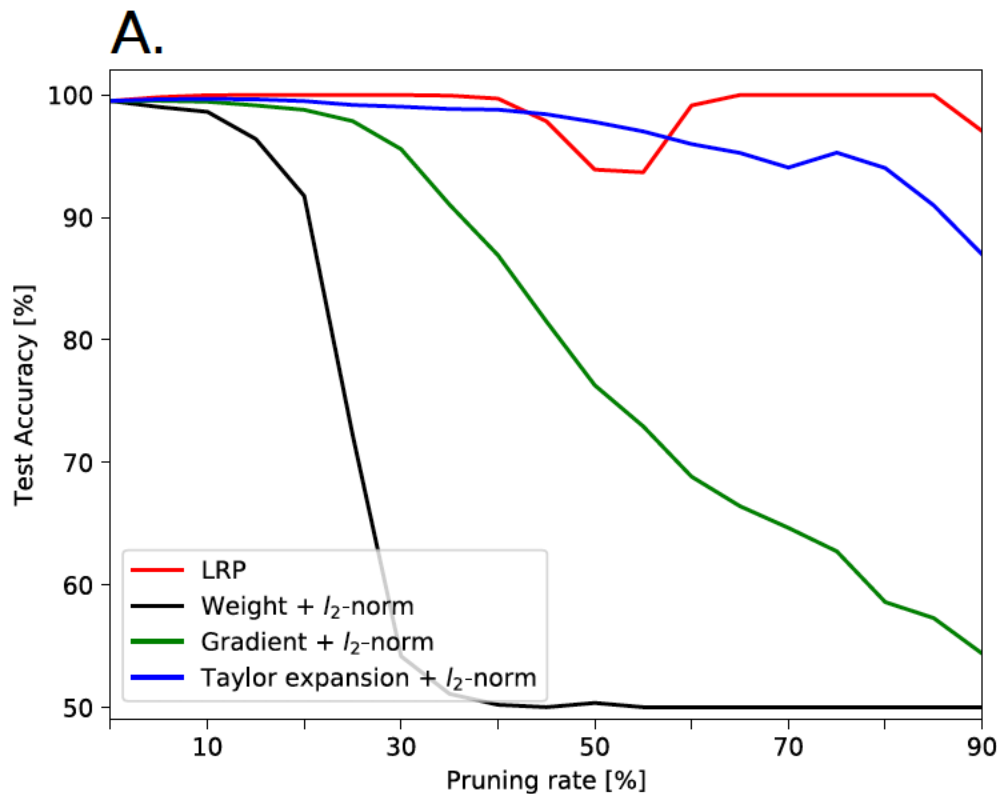
## Observations:

- small weights can be relevant
- LRP conservation property act as normalizer (better comparability across layers)

(Yeom et al. 2019)

With fine-tuning

# XAI-Based Pruning



Idea: Take generic model (VGG, 1000 classes) and make it smaller and more specialized (cats vs. dogs only) by pruning.

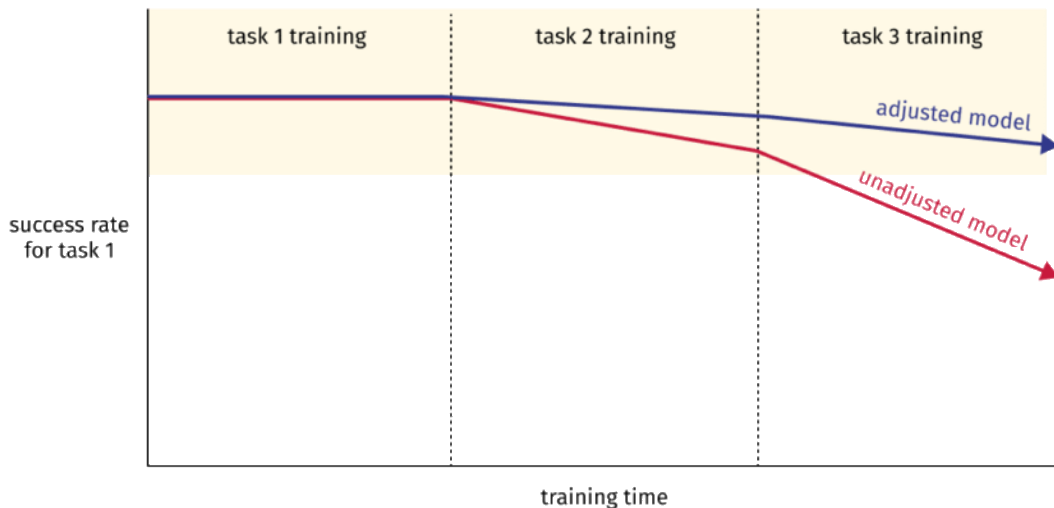
Only 10 samples per class (domain adaptation scenario)

(Yeom et al. 2019)

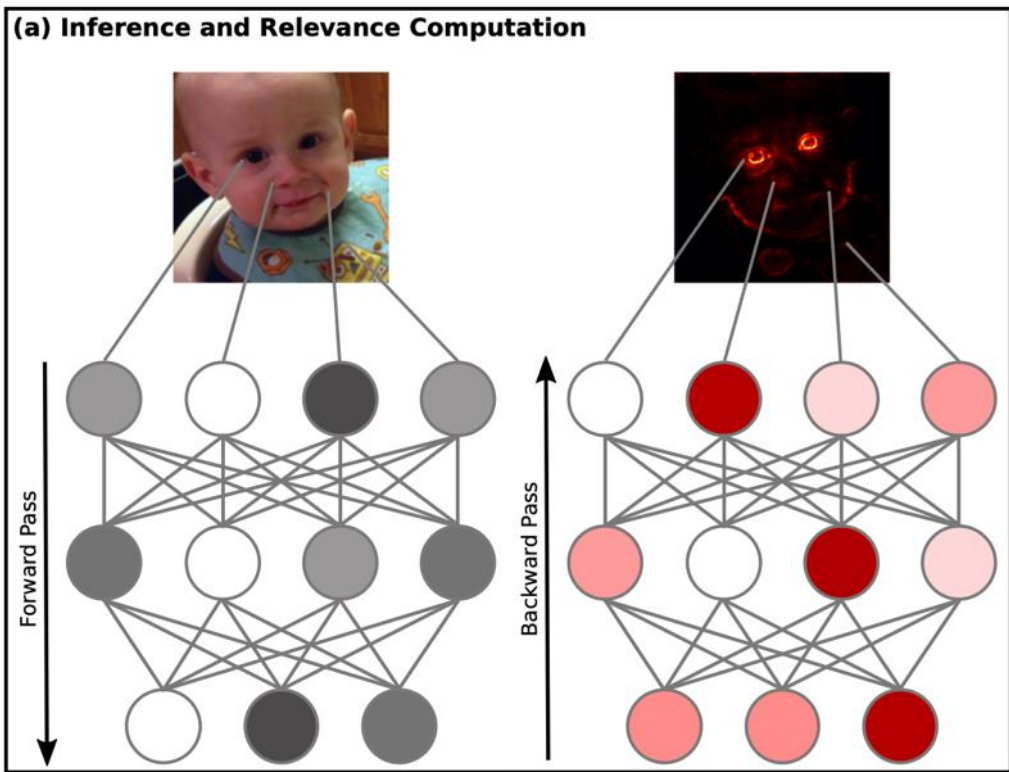
Without fine-tuning

# Catastrophic Forgetting

*Catastrophic forgetting is the tendency of an artificial neural network to abruptly and drastically forget previously learned information upon learning new information.*



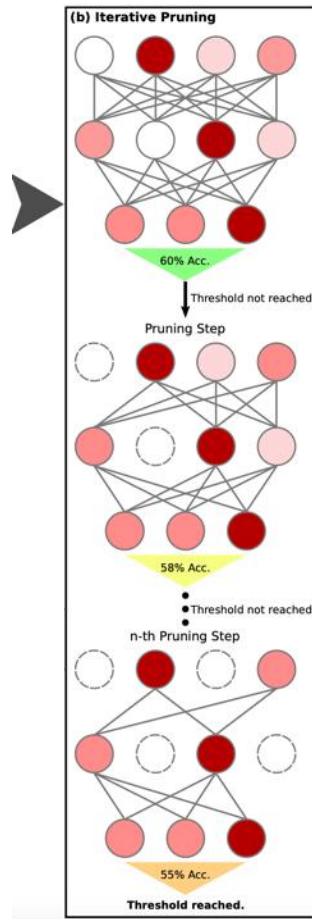
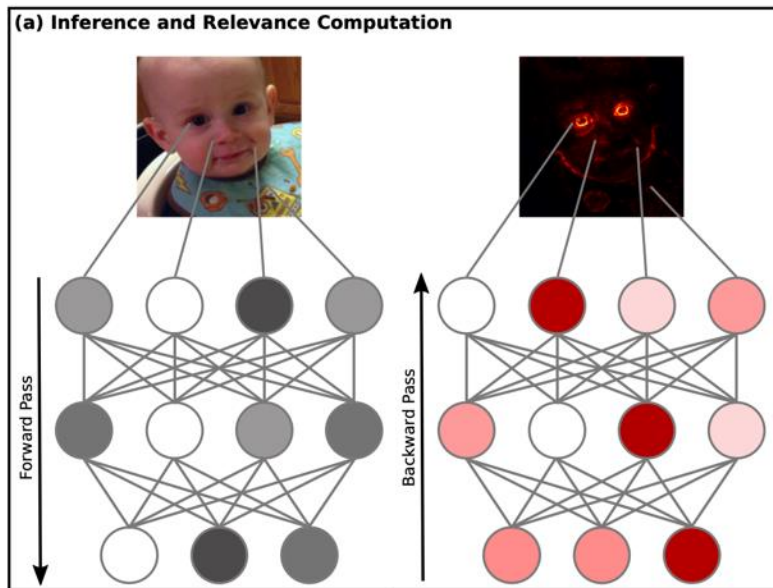
# Tackling Catastrophic Forgetting with XAI



1. Step: Identify relevant subnetworks

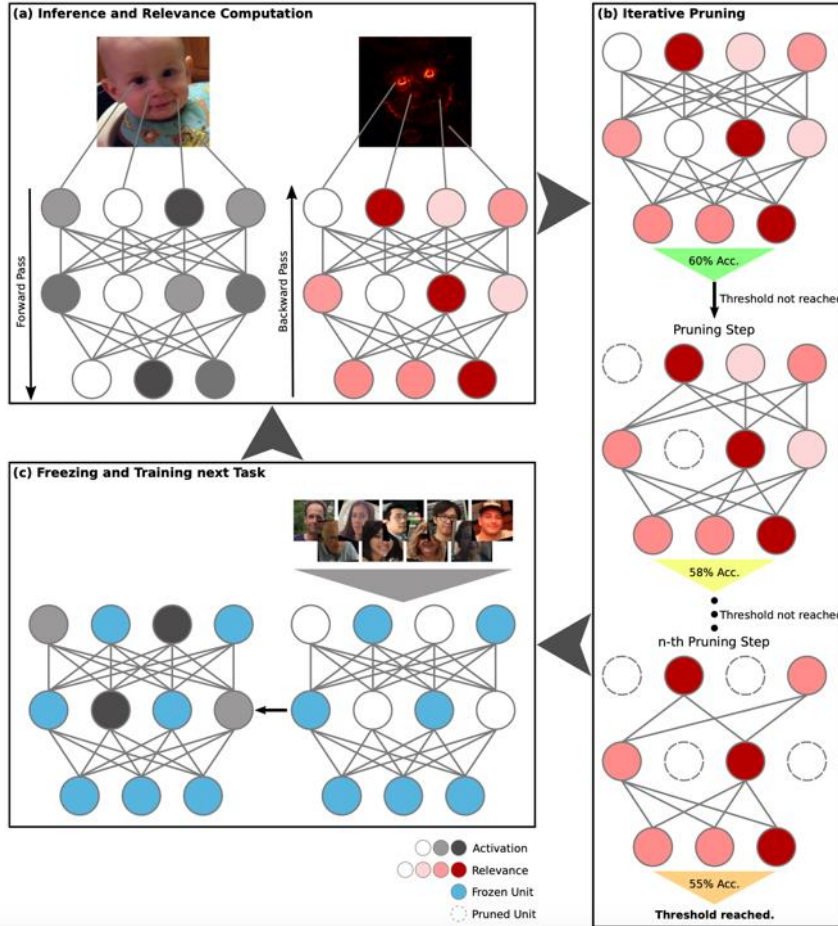


# Tackling Catastrophic Forgetting with XAI



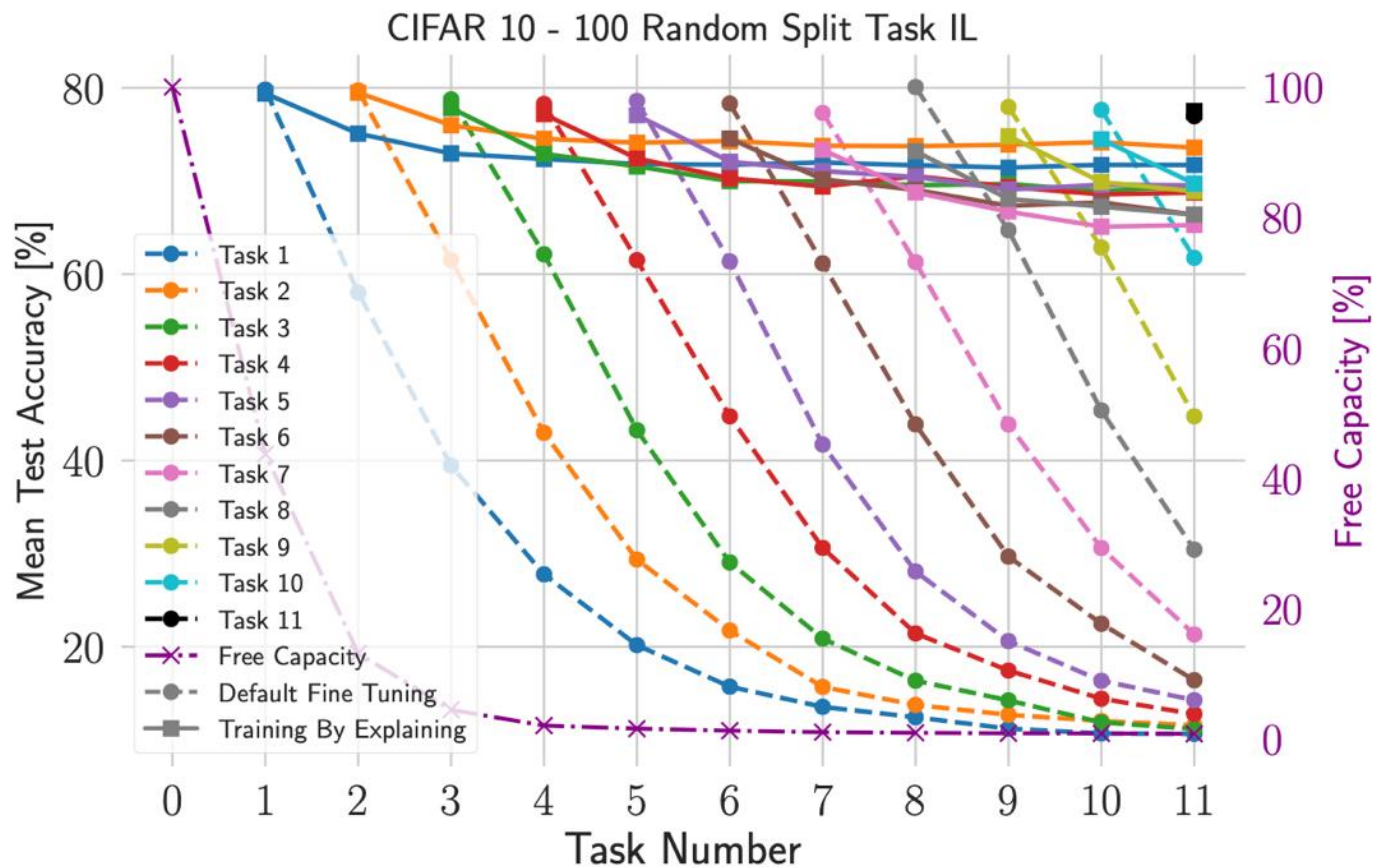
1. Step: Identify relevant subnetworks.
2. Prune model

# Tackling Catastrophic Forgetting with XAI



1. Step: Identify relevant subnetworks.
2. Prune model
3. Freeze and Train next Task

# Tackling Catastrophic Forgetting with XAI



(Ede et al. 2022)

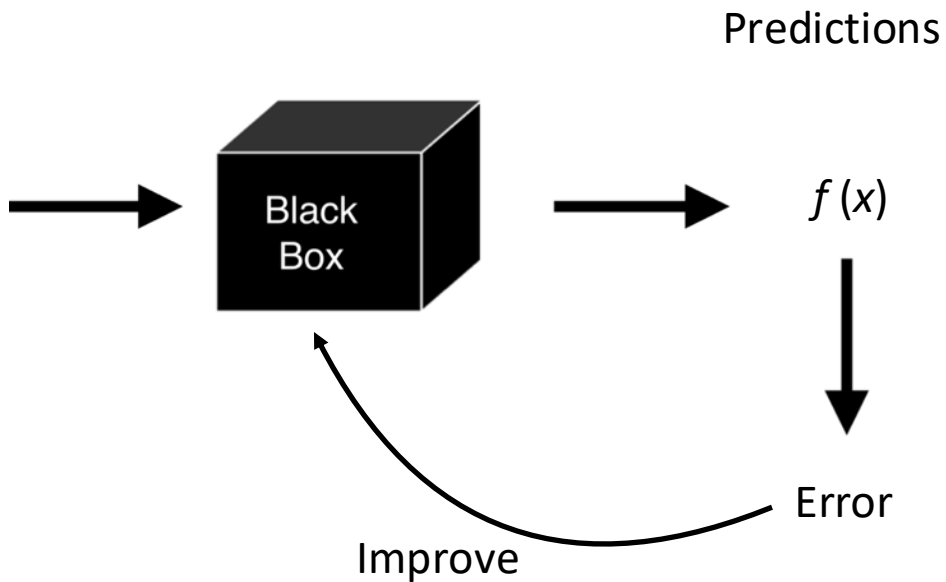
# 9th Take Home Message

*Understanding the model better opens up a lot of opportunities beyond "just explaining" its predictions.*

# **Reveal and Revise**

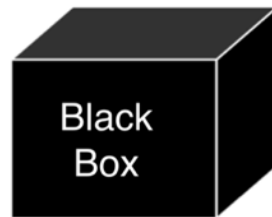
# Model Training

Training data  $x$  with labels  $y$



# Model Training

Training data  $x$  with labels  $y$



Predictions

$f(x)$



Error

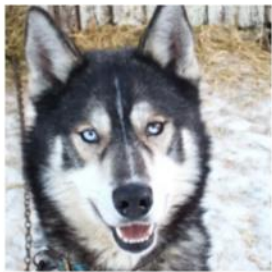
Improve

$$\min_{f \in \mathcal{F}} \int_{x,y} \|f(x) - y\|^2 dp(x,y)$$

Is minimizing the error a guarantee for the model to work well in practice?



# Clever Hans



(a) Husky classified as wolf

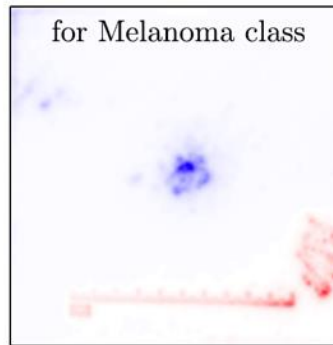
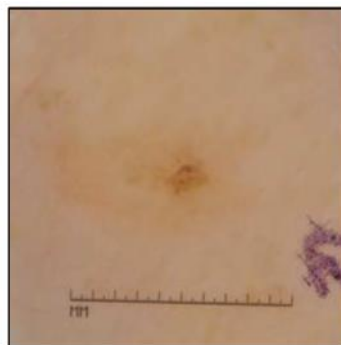
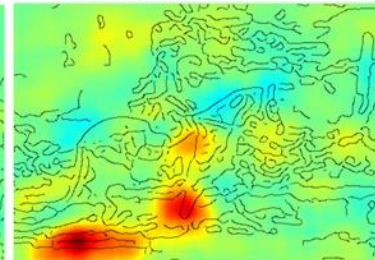
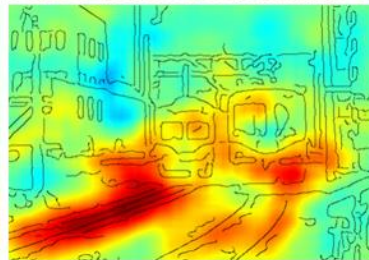
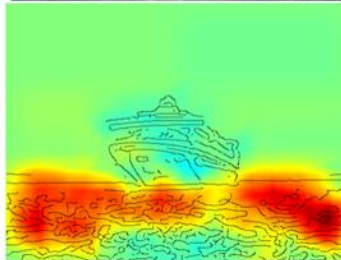
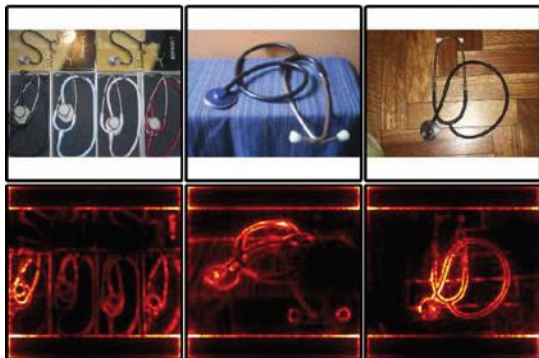


(b) Explanation

stethoscope  
1

stethoscope  
1

stethoscope  
1





# XAI-Based Model Improvement

XAI-Bases model improvement aims to improve the:

Reasoning: alignment with human, better generalization (no Clever Hans), no bias

# XAI-Based Model Improvement

XAI-Bases model improvement aims to improve the:

Reasoning: alignment with human, better generalization (no Clever Hans), no bias

Other Goals:

Performance: better test performance

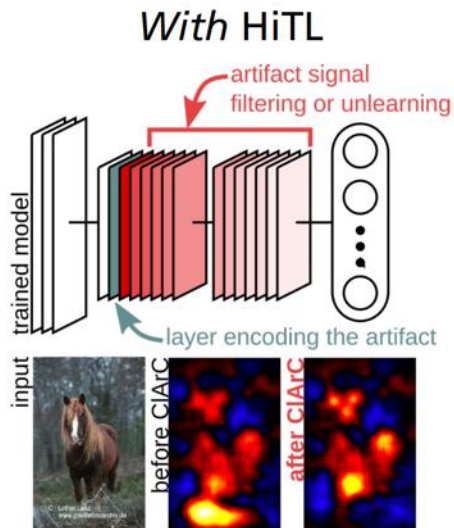
Convergence: faster learning

Robustness: robut against adversarial attackes and changes in distribution

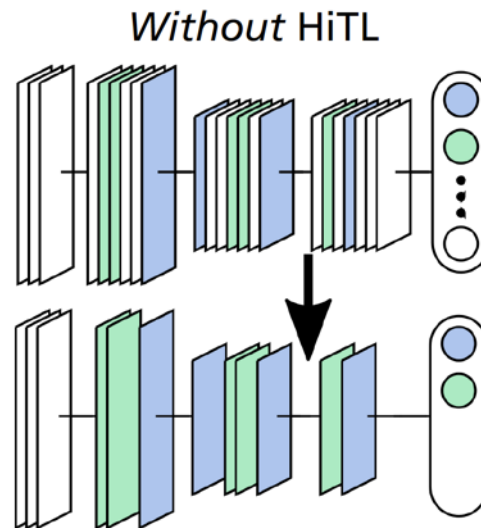
Efficiency: less (labeled) data, faster inference, less energy

Equality: better handle unbalanced data, minority classes

# XAI-Based Model Improvement

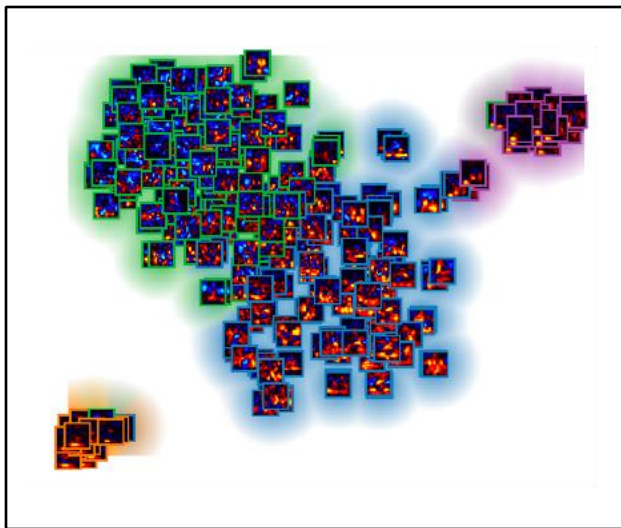


“semantic” model improvement, e.g.,  
[Schramowski *et al.* 2020; Anders, Weber, *et al.* 2022]

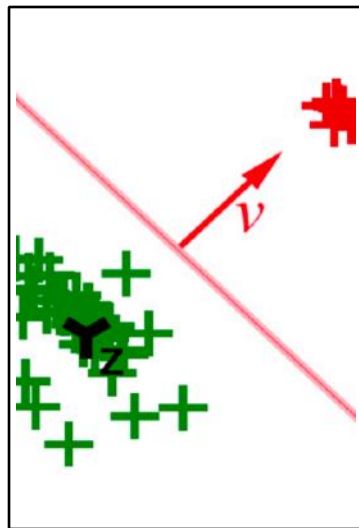


# Two Steps of Model Improvement

Analyse explanations



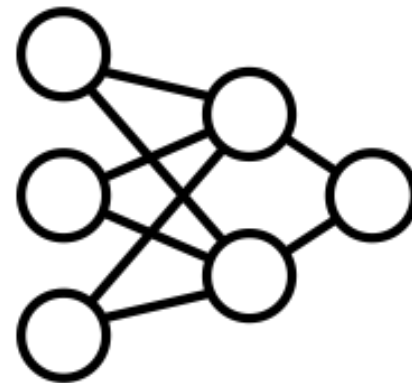
Artifact modelling



Clever Hans?



Correct model



**Step 1: Reveal**  
(or given a priori)

**Step 2: Revise**

# Two Steps of Model Improvement

Analyse explanations

Artifact modelling

Clever Hans?

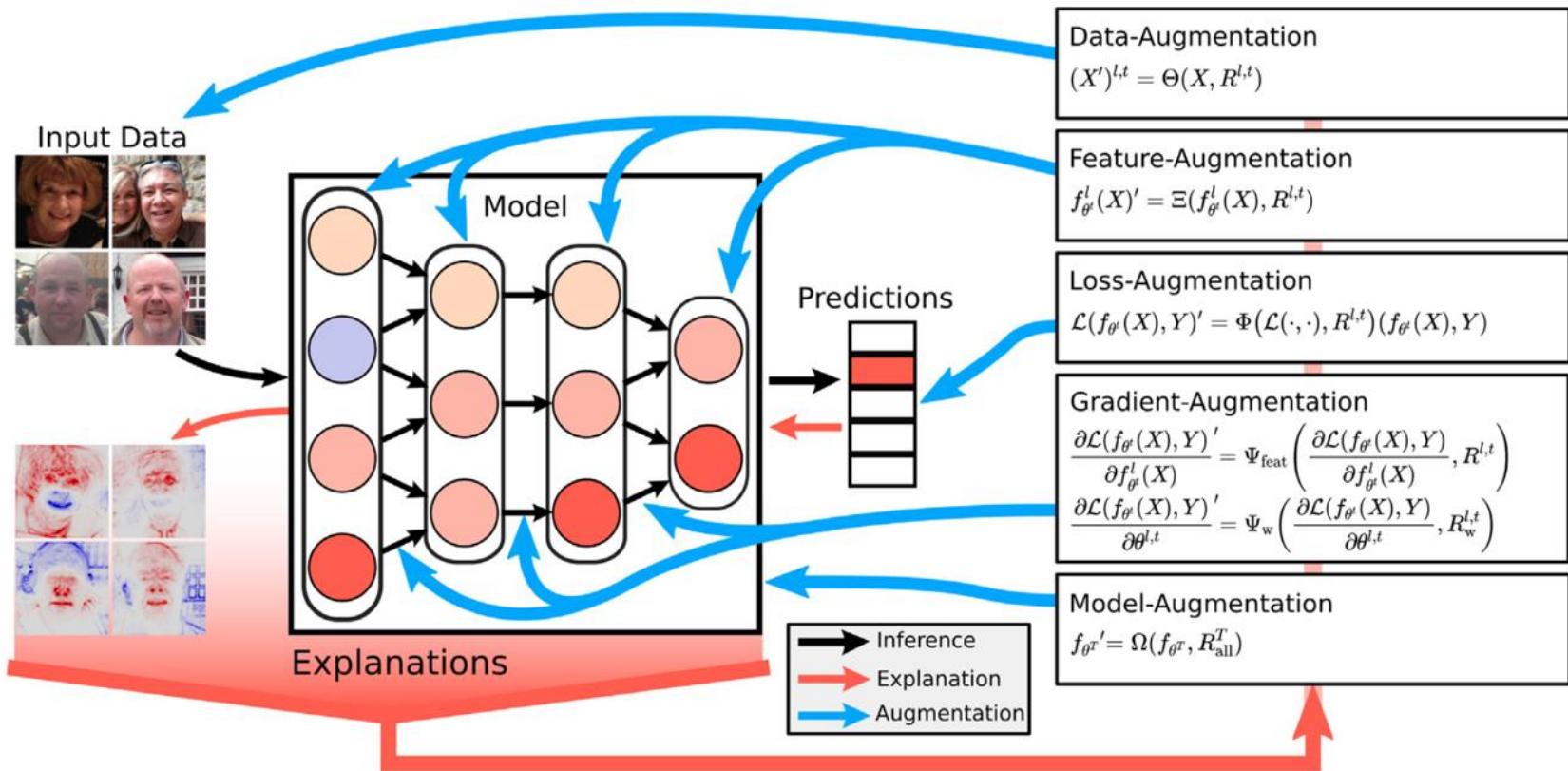
Correct model



**Step 1: Reveal**  
(or given a priori)

**Step 2: Revise**

# Where to Integrate Explanations?



Step 1a: Detect the Artifact  
with Explanations (Reveal)

# Reveal Clever Hans

(a) directly on data



**infeasible:**  
does not reveal  
model behavior





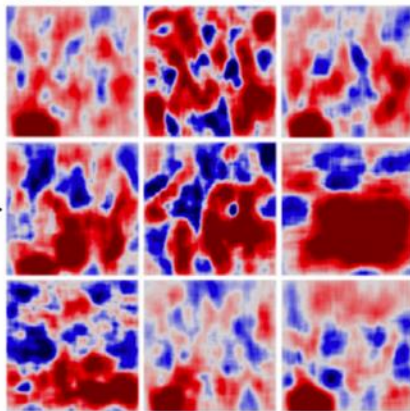
# Reveal Clever Hans

(a) directly on data



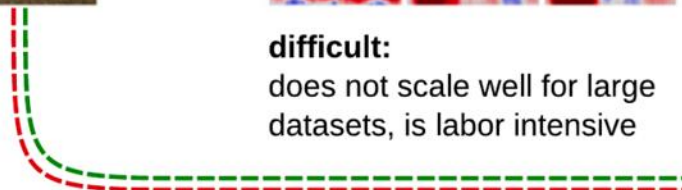
**infeasible:**  
does not reveal  
model behavior

(b) per-sample XAI



input  
for

**difficult:**  
does not scale well for large  
datasets, is labor intensive



# Reveal Clever Hans

Idea: Clustering of explanations

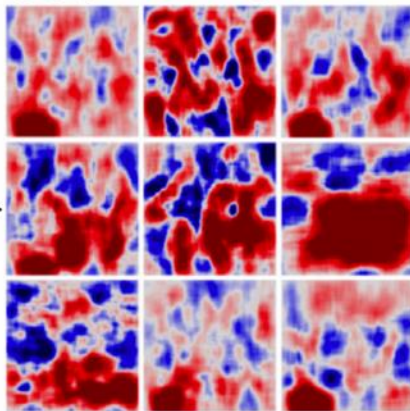
(a) directly on data



**infeasible:**  
does not reveal  
model behavior

input  
for

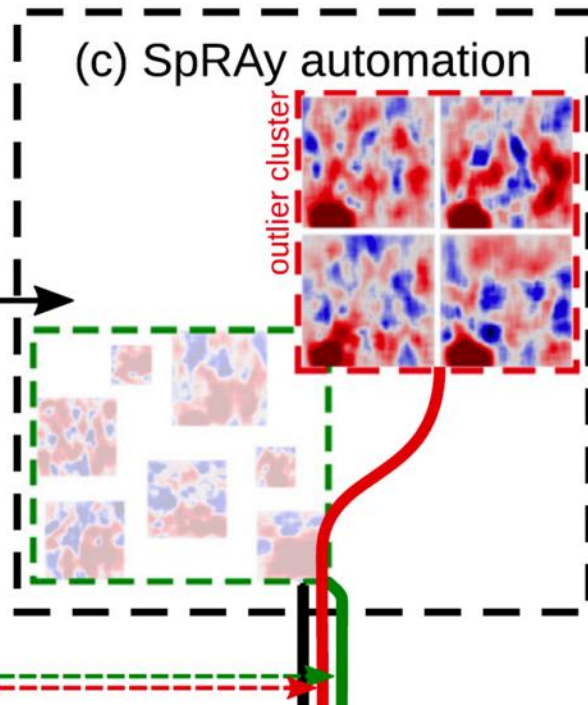
(b) per-sample XAI



**difficult:**  
does not scale well for large  
datasets, is labor intensive

input  
for

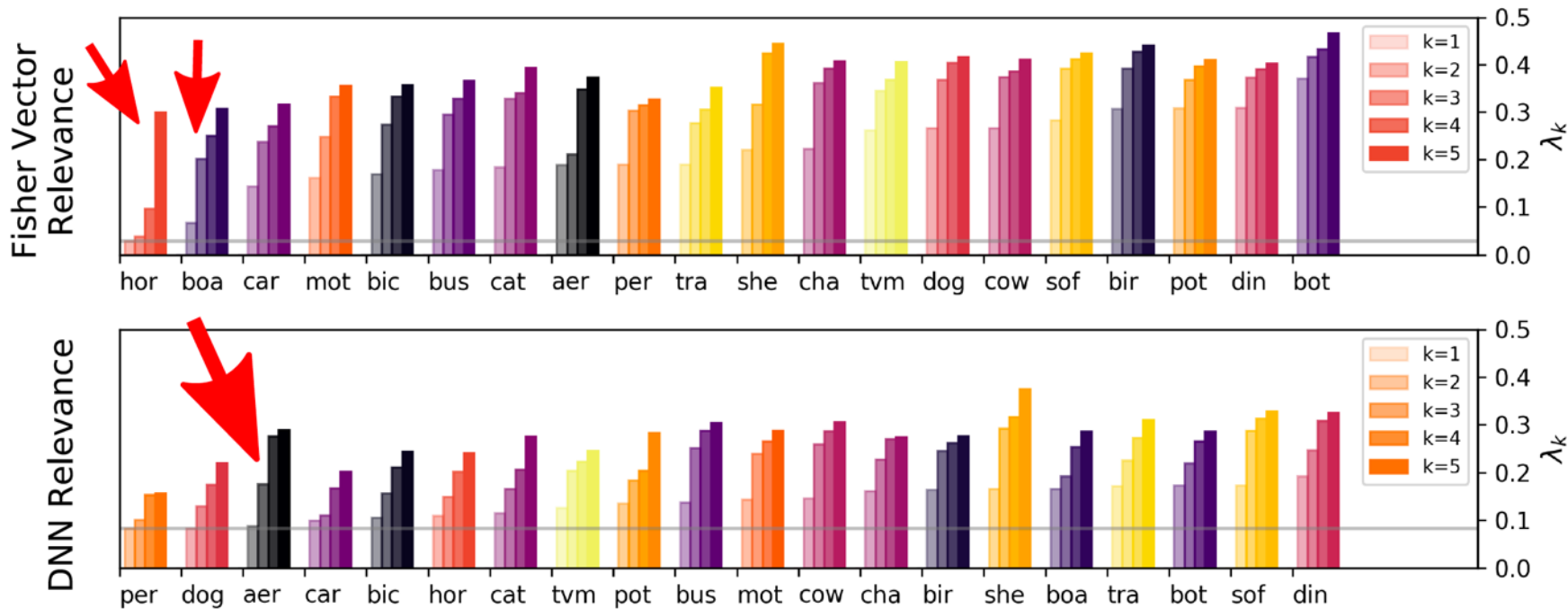
(c) SpRAY automation



SpRAY = Spectral Relevance Analysis

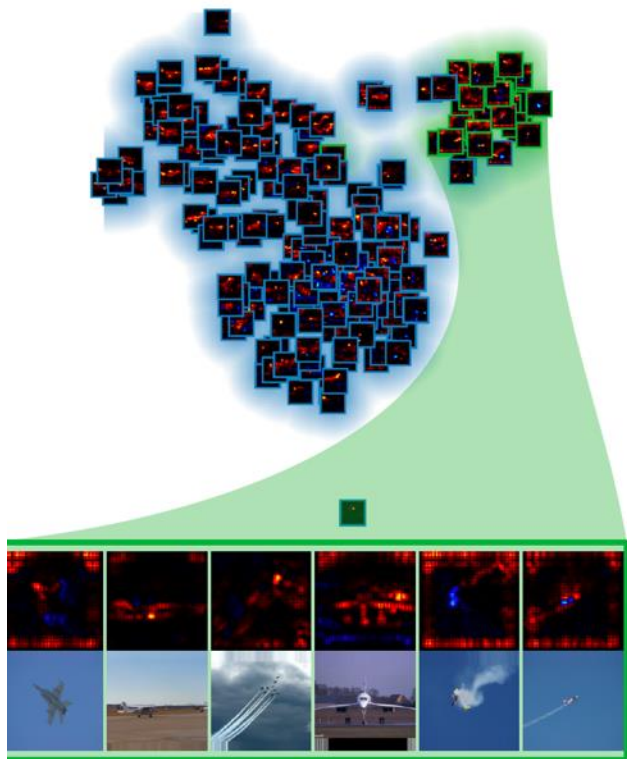
# Spectral Clustering

A strong increase in the difference between two successive eigenvalues (**eigengap**) indicates w



# Aeroplane Cluster

DNN Heatmaps

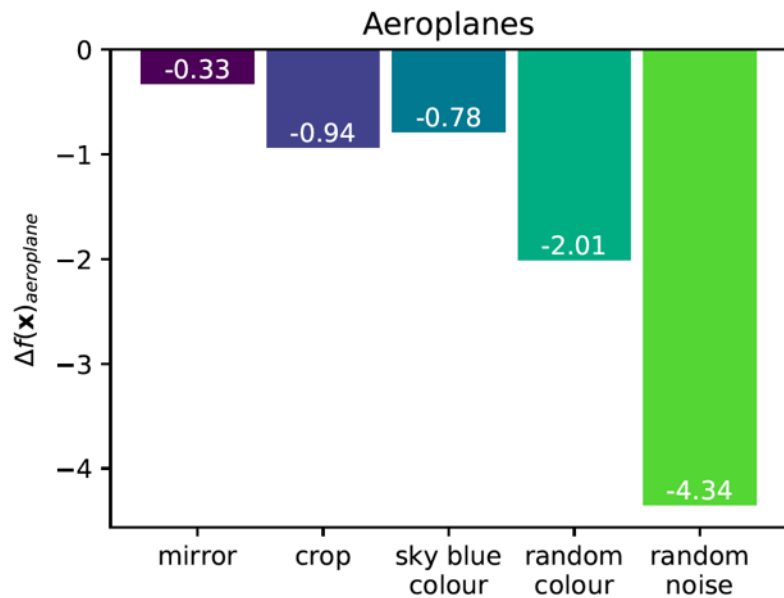
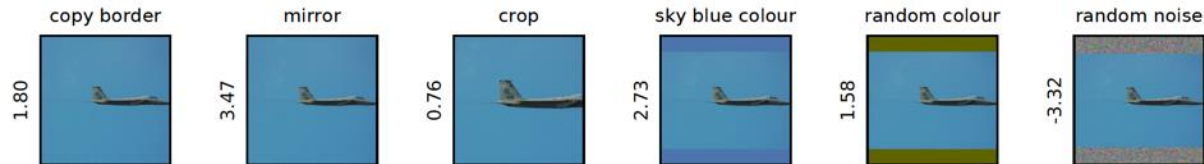
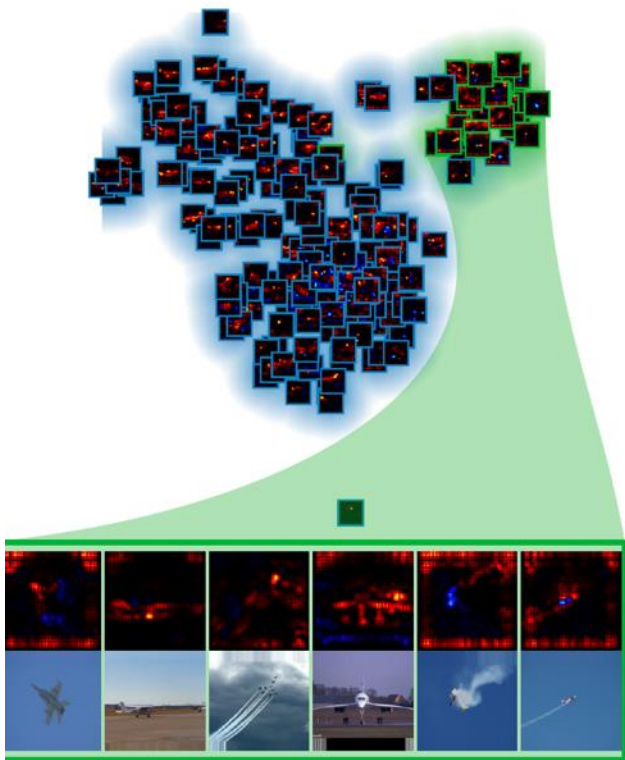


Observation: Artifact at the border.

(Lapuschkin et al. 2019)

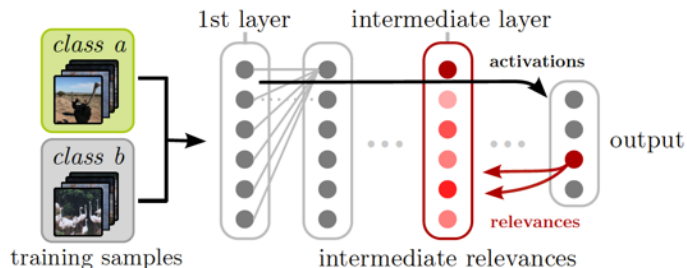
# Aeroplane Cluster

DNN Heatmaps



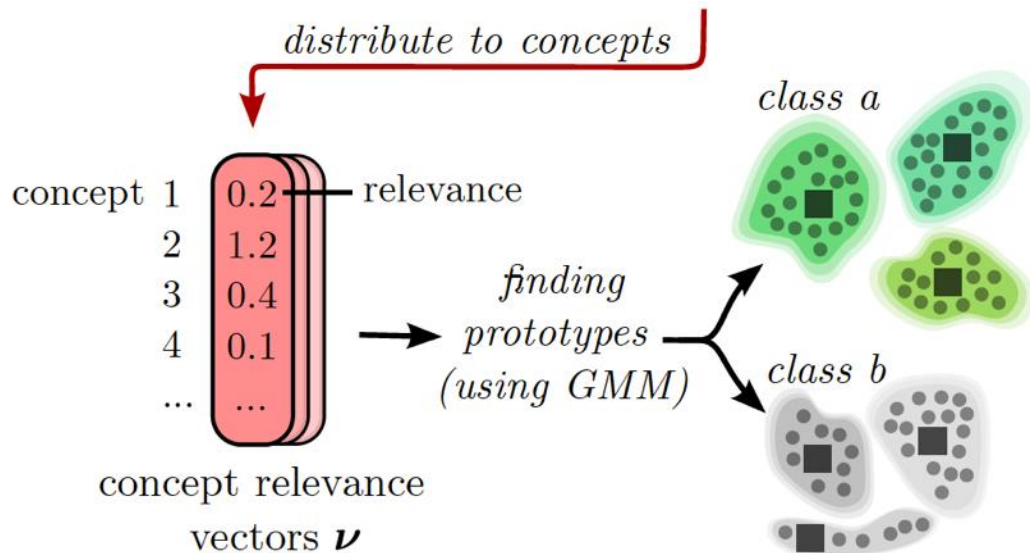
**Limitation of Input Space Analysis: Natural invariances (rotation, translation, scaling) are difficult to capture.**

# Prototypical Concept-based Explanations



Quantifying the (Extra-)Ordinary:

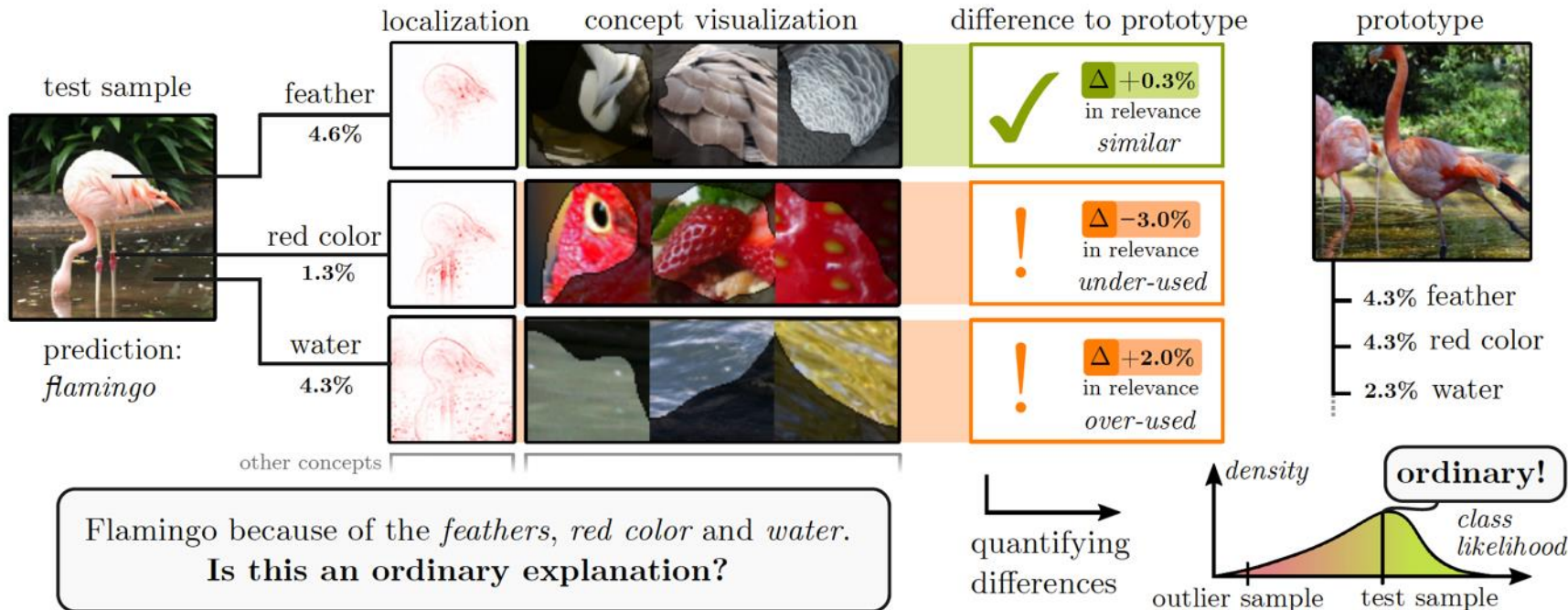
$$L^k(\boldsymbol{\nu}) = \log p^k(\boldsymbol{\nu})$$



$$p_i^k(\boldsymbol{\nu}) = \frac{1}{(2\pi)^{\frac{n}{2}} \det(\boldsymbol{\Sigma}_i^k)^{\frac{1}{2}}} e^{-\frac{1}{2}(\boldsymbol{\nu} - \boldsymbol{\mu}_i^k)^\top (\boldsymbol{\Sigma}_i^k)^{-1} (\boldsymbol{\nu} - \boldsymbol{\mu}_i^k)}$$



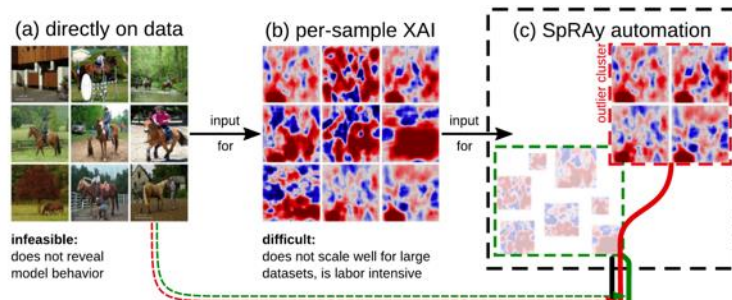
# Prototypical Concept-based Explanations





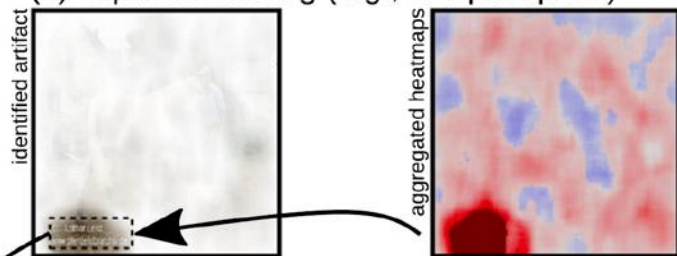
Step 1b: Model the Artifact

# Estimate Artifact Model



## II: Estimate Artifact Model

### (a) explicit modeling (e.g., in input space)

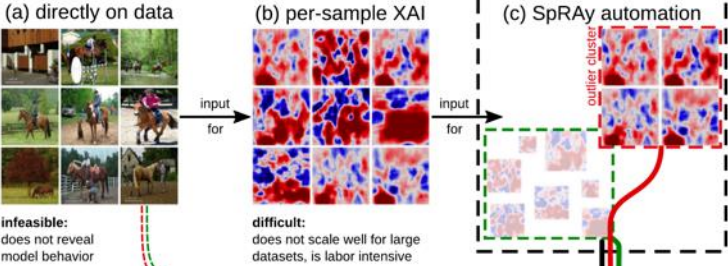


extracted artifact

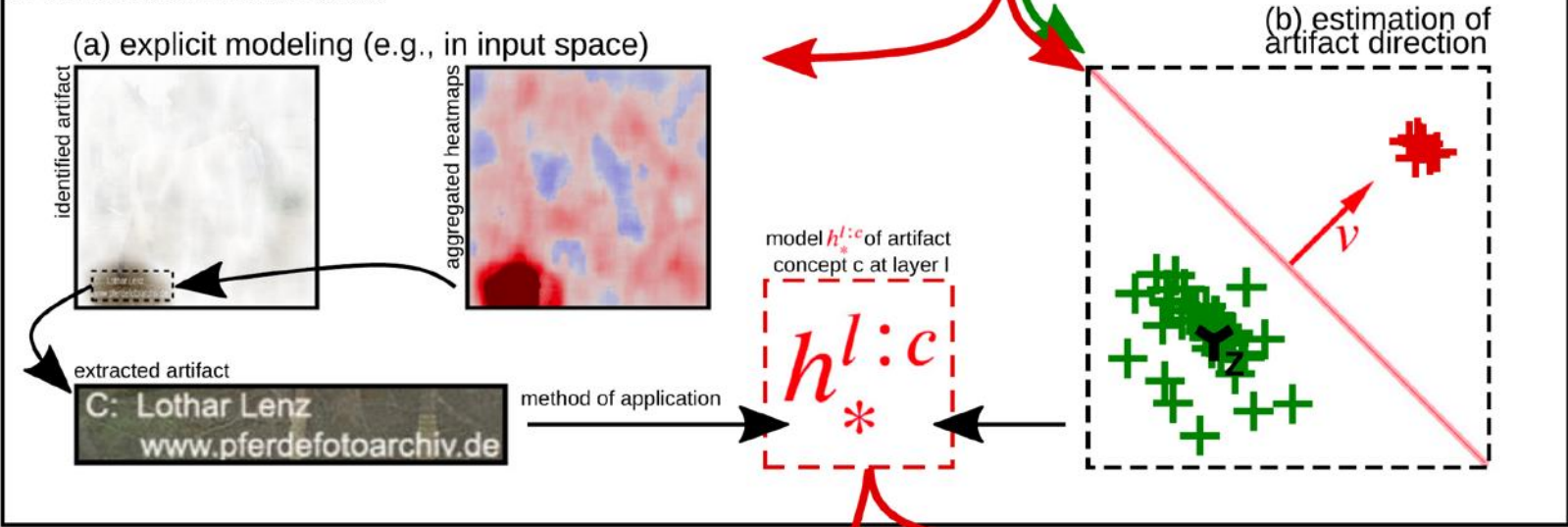
C: Lothar Lenz  
www.pferdefotoarchiv.de



# Estimate Artifact Model

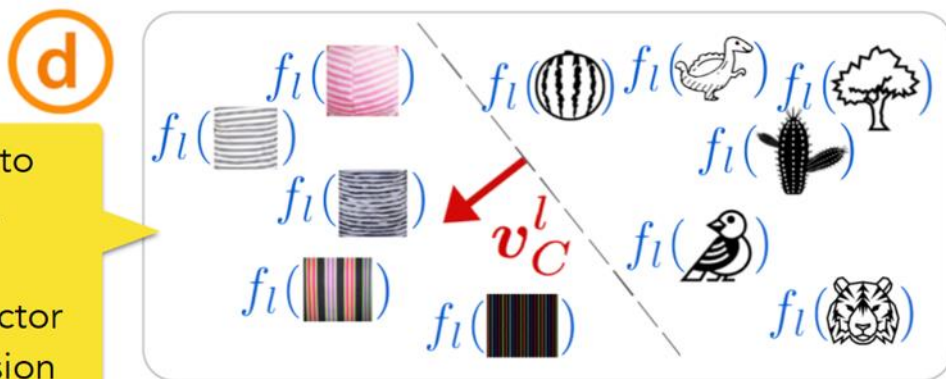
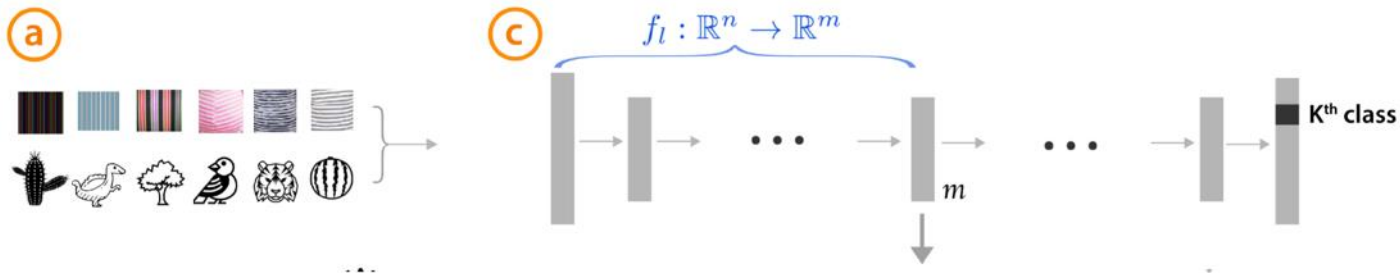


## II: Estimate Artifact Model



# Revisited: Concept Activation Vector (CAV)

## Inputs:



Train a linear classifier to separate activations.

CAV ( $v_C^l$ ) is the vector **orthogonal** to the decision boundary.

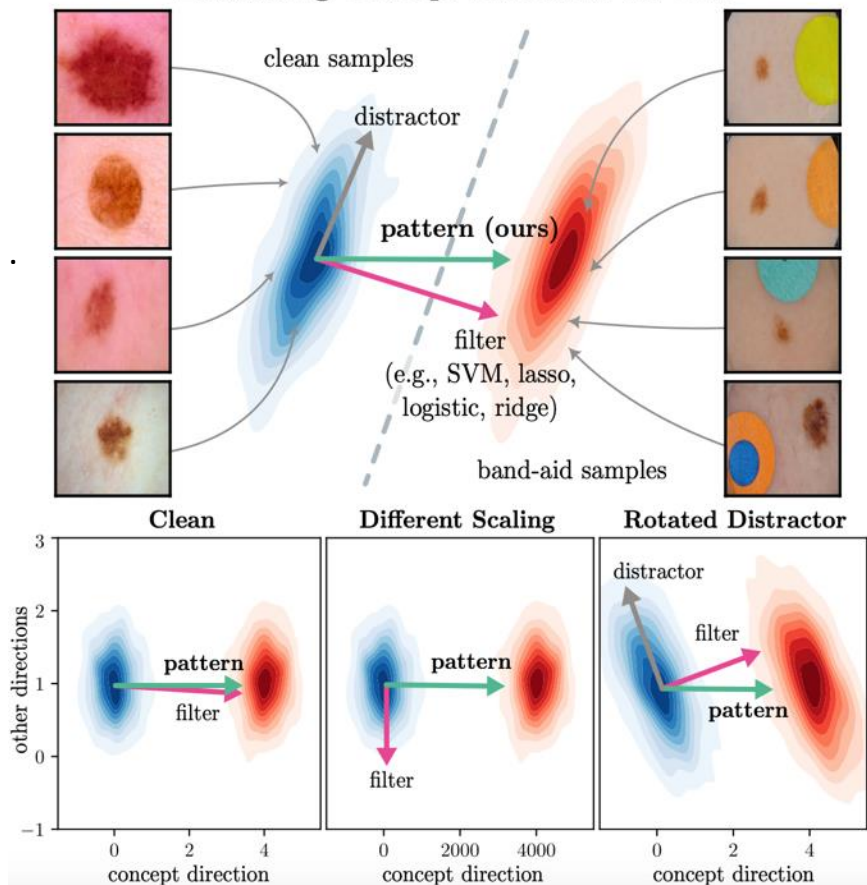
[Smilkov '17, Bolukbasi '16, Schmidt '15]

# Pattern-CAV: New Noise Resilient CAV

## Filter-Pattern Problem (Haufe et al. 2014):

Significant influence of distractor (i.e., non-signal) directions contained in the data, which are picked up by filters (i.e., weights) of linear models to optimize class-separability

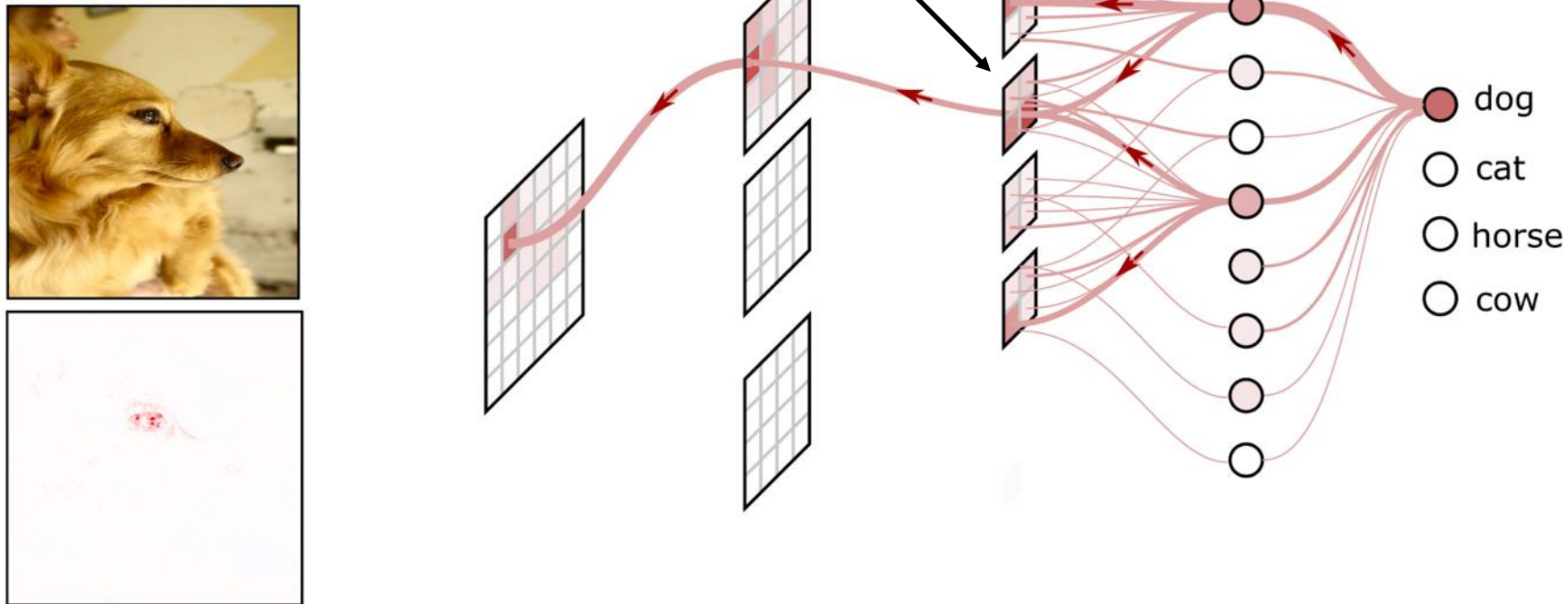
## Estimating Concept Direction via CAV



## Navigating Neural Space: Revisiting Concept Activation Vectors to Overcome Directional Divergence

Frederik Pahde<sup>1</sup>, Maximilian Dreyer<sup>1</sup>, Leander Weber<sup>1</sup>, Moritz Weckbecker<sup>1</sup>,  
Christopher J. Anders<sup>2,3</sup>, Thomas Wiegand<sup>1,2,3</sup>, Wojciech Samek<sup>1,2,3,†</sup>, Sebastian Lapuschkin<sup>1,†</sup>

# Back to Input Space



LRP / CRP can be used to localize CAV in input space.

Step 2: Correct the Model  
Behaviour (Revise)



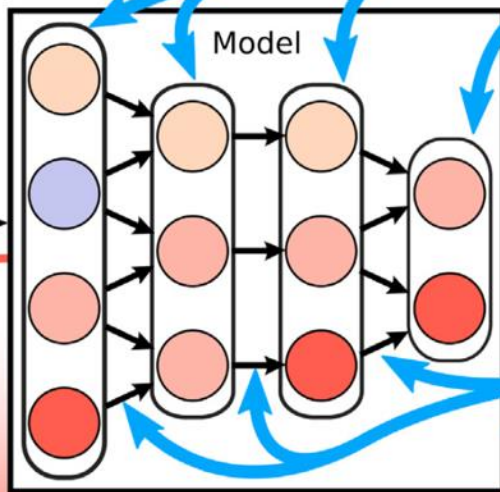
# XAI-Based Model Improvement

Data-Augmentation

$$(X')^{l,t} = \Theta(X, R^{l,t})$$

*Generate new samples depending on explanations or alter the distribution of existing data*

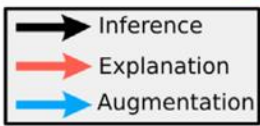
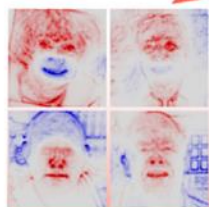
Input Data



Predictions



Explanations





# Data Augmentation



Isolate artefact

# Data Augmentation



Isolate artefact, add to *other/all* classes

# Data Augmentation



unmodified fine-tuning



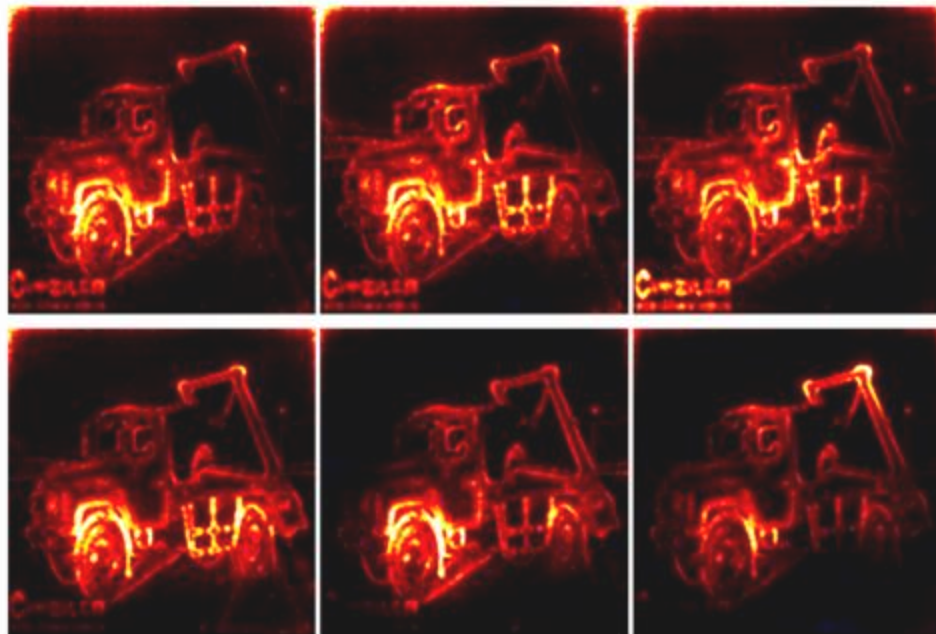
CIArC fine-tuning



1 epoch

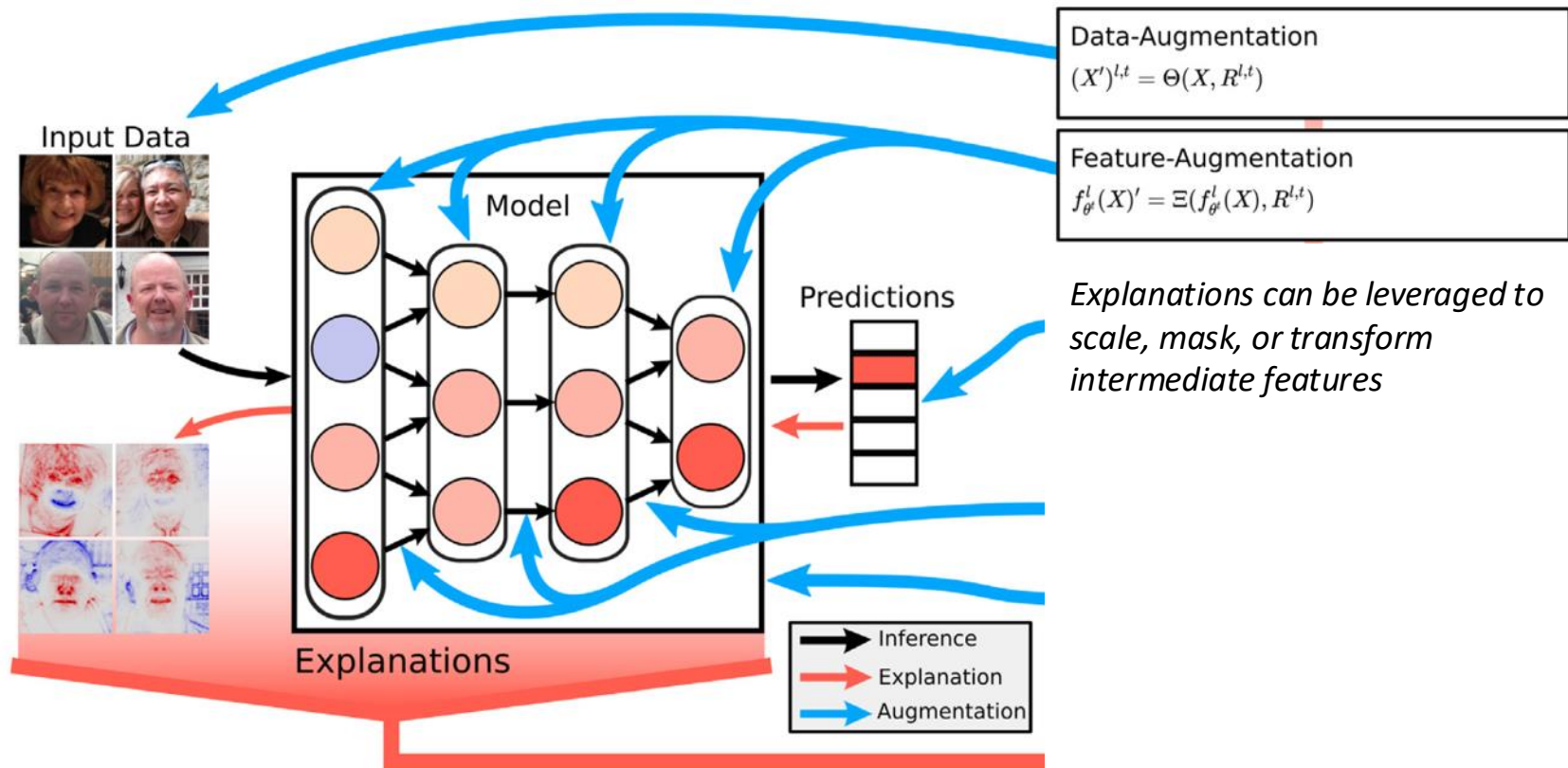
5 epochs

10 epochs



Isolate artefact, add to *other/all* classes, re-train model.

# XAI-Based Model Improvement

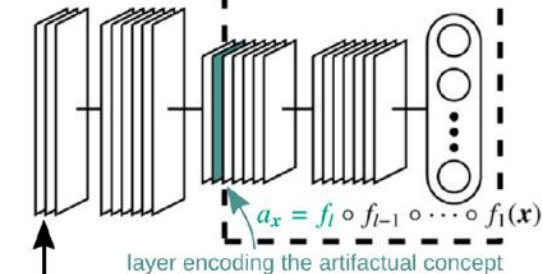


# Feature Augmentation

## III: Update Model to Reduce Impact of Artifact

original model

layers affected by CIArC

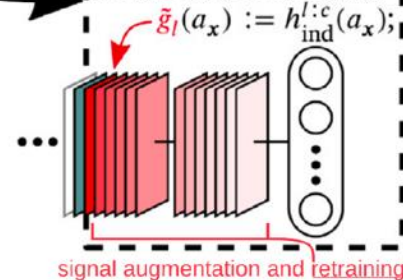


### Augmentative Class Artifact Compensation (A-CIArC)

Augment samples in such a way that the SGD-trained classifier becomes insensitive to an artifact.

--> Signal augmentation and retraining

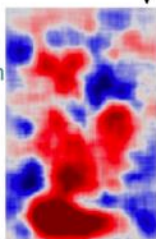
A-CIArC: Artifact Unlearning



input

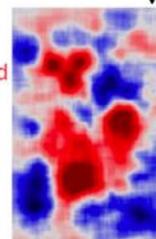


$f(x) = 0.94$   
"class 'horse' is recognized from a watermark in the image, and some supporting information [...]"



original explanation

$f(x) = 0.63$   
"class 'horse' is recognized from the expected horse-like information amplified during augmentative training [...]"

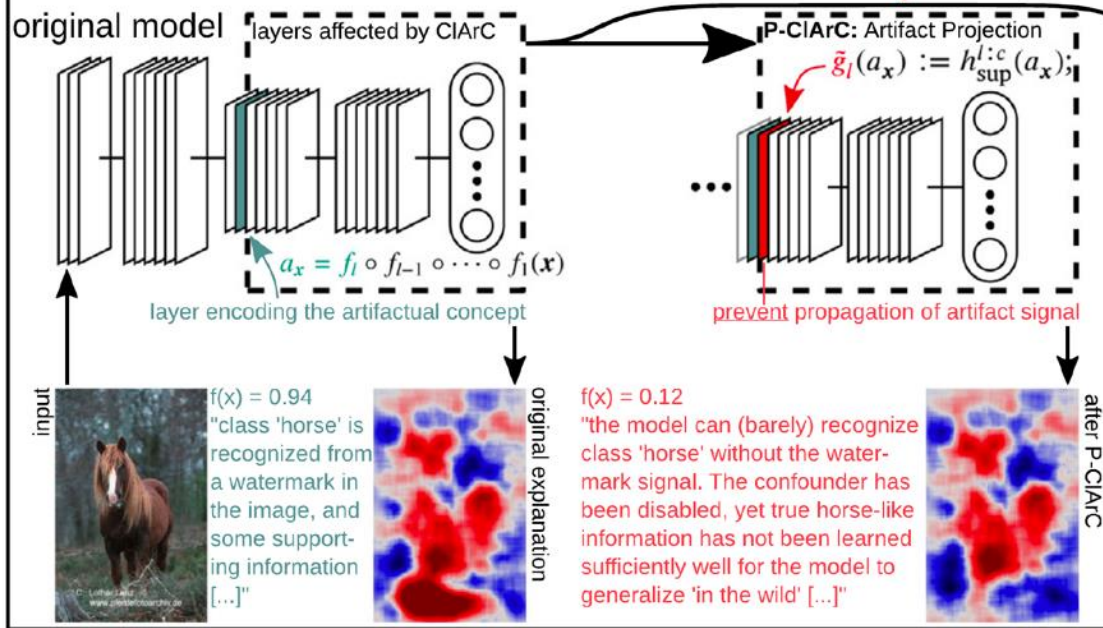


after A-CIArC

*inductive* removal with A-CIArC.

# Feature Augmentation

## III: Update Model to Reduce Impact of Artifact



## Projective Class Artifact Compensation

Correct the model without retraining by incorporating a suppressive artifact model directly into the prediction model.

--> Prevent propagation of artifact signal.

*inductive* removal with **A-ClArC**. *suppressive* removal with **P-ClArC**.



# Augmentative Class Artifact Compensation

Pushing samples around, to artifact reference point  $z$

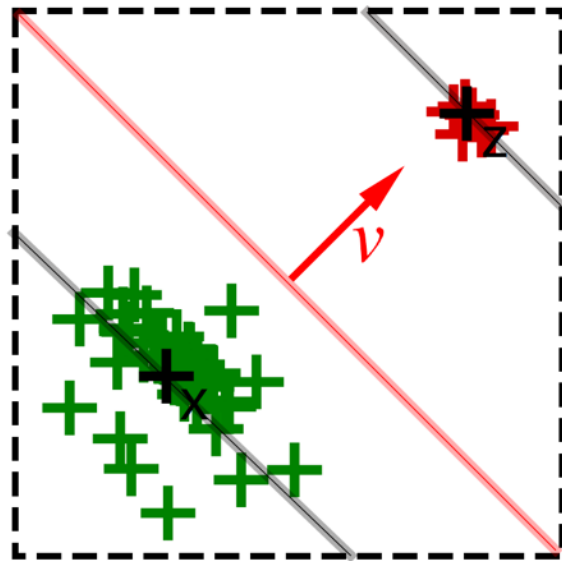
activation vector

artifact CAV

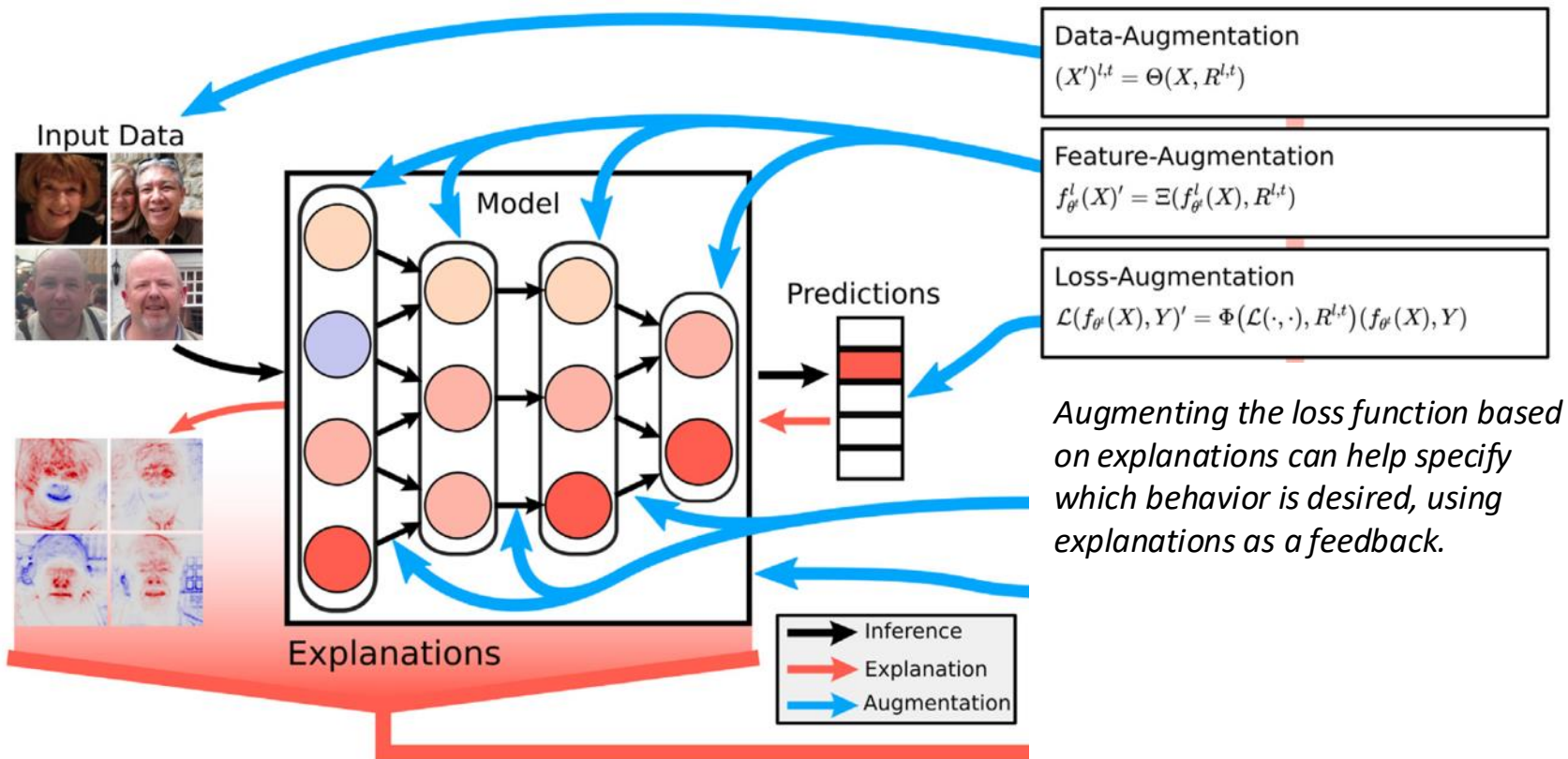
$$h_{\text{ind}}(\mathbf{x}) = \mathbf{x} + \mathbf{v} \cdot (-\mathbf{v}^\top \mathbf{x} + \mathbf{v}^\top \mathbf{z})$$

$$\mathbf{z} = \frac{1}{|X^+|} \sum_{\mathbf{x}^+ \in X^+} \mathbf{x}^+$$

pushing samples over the  
decision boundary (e.g.  
SVM CAV)

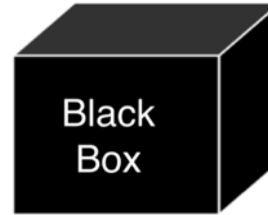


# XAI-Based Model Improvement





# Right for the Right Reason

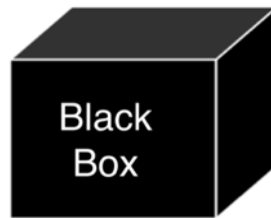


*"Look at the horse!"*

Otherwise you will be penalized



# Right for the Right Reason



*"Look at the horse!"*

Otherwise you will be penalized

"explanation" (sensitivity)

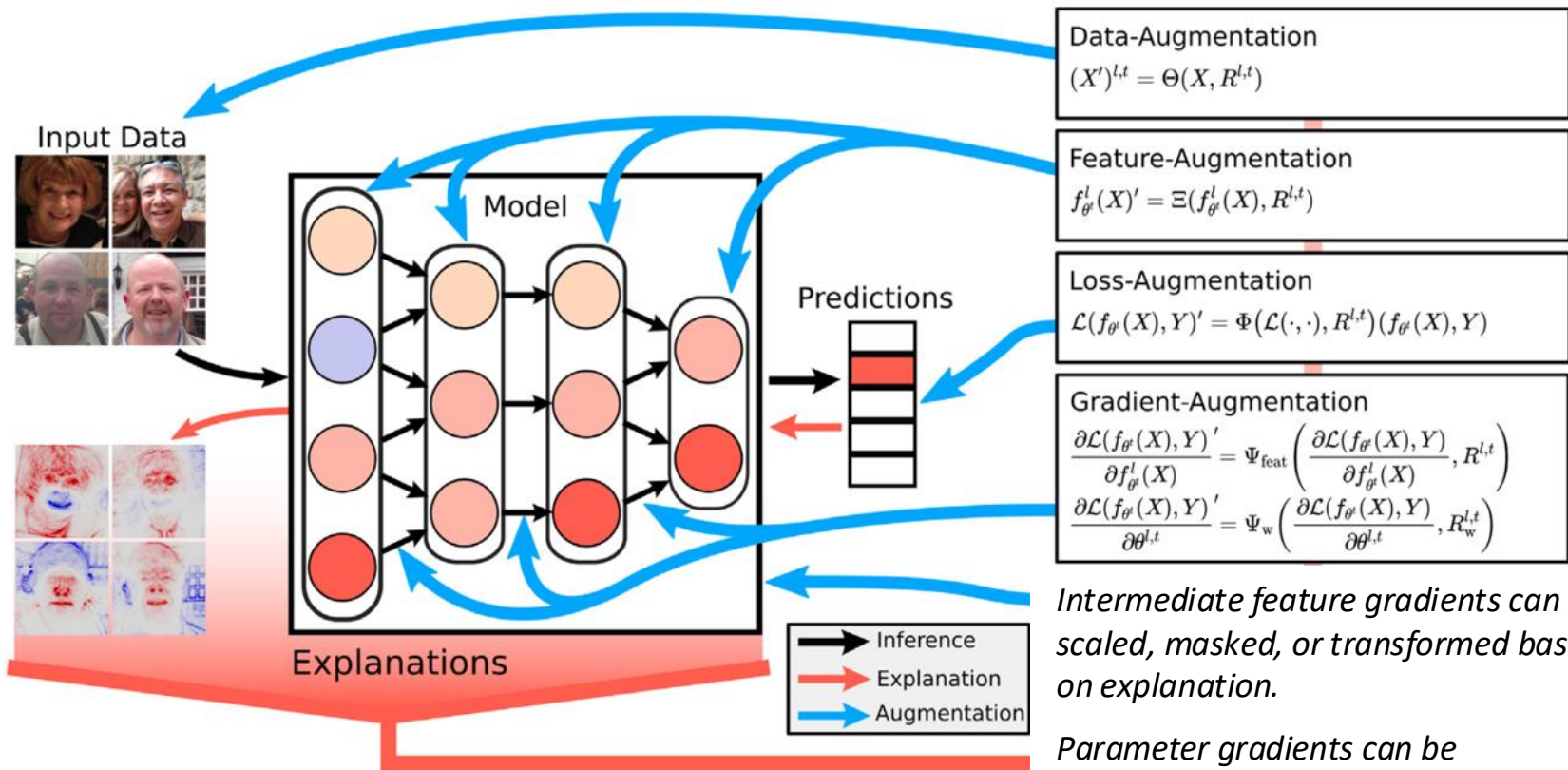
$$L(\theta, X, y, A) = \underbrace{\sum_{n=1}^N \sum_{k=1}^K -y_{nk} \log(\hat{y}_{nk})}_{\text{Right answers}}$$

binary mask

$$+ \lambda_1 \underbrace{\sum_{n=1}^N \sum_{d=1}^D \left( A_{nd} \frac{\partial}{\partial x_{nd}} \sum_{k=1}^K \log(\hat{y}_{nk}) \right)^2}_{\text{Right reasons}} + \lambda_2 \underbrace{\sum_i \theta_i^2}_{\text{Regular}}$$

(Ross et al. 2017)

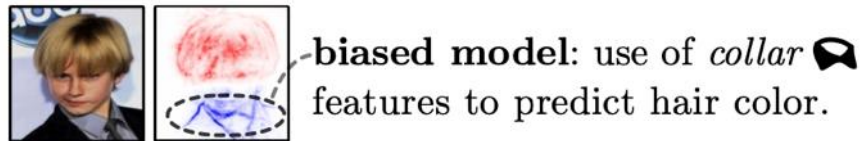
# XAI-Based Model Improvement



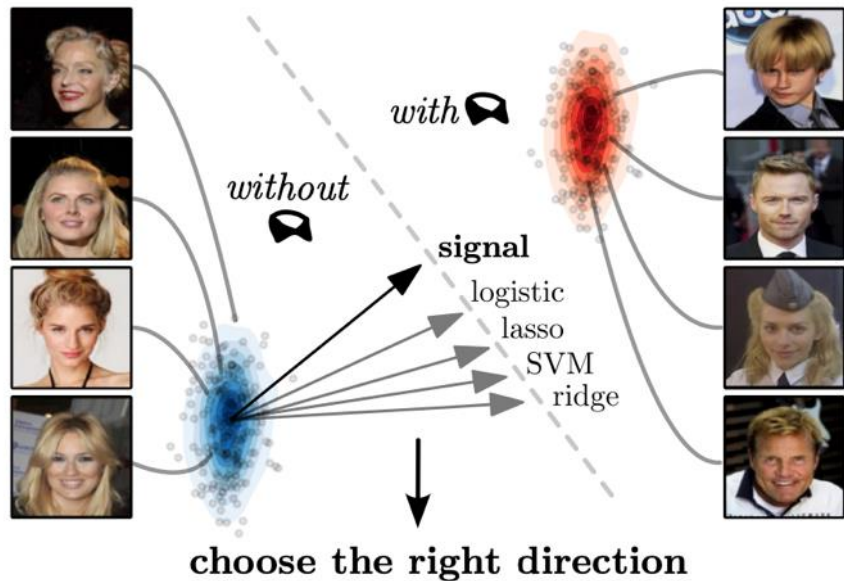
*Intermediate feature gradients can be scaled, masked, or transformed based on explanation.*

*Parameter gradients can be augmented directly by computing parameter-wise relevance scores.*

# RR-CIArC



a modeling bias via CAV



b correcting bias

bias samples



Vanilla  
(without correction)



A-CIArC



activation space

RR-CIArC (ours)



gradient space

ensure the right reasons

(Dreyer et al. 2024)

# RR-CIArC

$$L_{\text{RR}}(\mathbf{x}) = \left( \nabla_{\mathbf{a}} \tilde{f}(\mathbf{a}(\mathbf{x})) \cdot \mathbf{h} \right)^2$$

sensitivity w.r.t. latent features

bias direction (artifact CAV)

Interpretation: Loss penalizes the model for changing the output when slightly adding or removing activations along the bias direction  $\mathbf{h}$ .

$$\lim_{\epsilon \rightarrow 0} \frac{\tilde{f}(\mathbf{a}(\mathbf{x}) + \epsilon \mathbf{h}) - \tilde{f}(\mathbf{a}(\mathbf{x}))}{\epsilon} \approx \nabla_{\mathbf{a}} \tilde{f}(\mathbf{a}(\mathbf{x})) \cdot \mathbf{h} \stackrel{!}{=} 0$$

# RR-CIArC

architecture	method	Bone Age			ISIC			ImageNet			CelebA		
		<i>clean</i>	<i>biased</i>	TCAV	<i>clean</i>	<i>biased</i>	TCAV	<i>clean</i>	<i>biased</i>	TCAV	<i>clean</i>	<i>biased</i>	TCAV
VGG-16	<i>Vanilla</i>	78.8	49.8	0.86	76.2	34.9	0.84	68.7	43.5	0.63	93.7	82.8	0.37
	RRR	78.8	49.8	0.86	76.7	42.8	0.72	68.6	49.6	0.55	93.7	91.2	0.43
	P-CIArC	78.9	77.4	0.66	75.1	49.0	0.77	68.3	<b>62.6</b>	0.37	56.6	60.8	0.19
	A-CIArC	77.8	69.0	0.66	75.2	49.5	0.65	67.7	60.9	<b>0.49</b>	93.0	90.4	0.44
	RR-CIArC (ours)	78.8	<b>77.7</b>	<b>0.52</b>	74.3	<b>57.0</b>	<b>0.49</b>	68.5	<b>62.6</b>	<b>0.49</b>	93.6	<b>92.6</b>	<b>0.54</b>
ResNet-18	<i>Vanilla</i>	75.1	46.3	1.00	81.8	56.8	1.00	66.7	52.9	1.00	96.8	58.3	0.21
	RRR	74.5	47.9	1.00	78.7	61.1	1.00	66.4	59.1	0.08	95.5	74.7	0.92
	P-CIArC	75.0	70.7	<b>0.60</b>	60.8	59.9	1.00	67.0	61.7	0.80	96.5	64.4	0.06
	A-CIArC	74.8	57.4	0.34	77.1	65.0	0.98	65.0	63.3	0.88	96.1	62.9	0.38
	RR-CIArC (ours)	71.1	<b>74.2</b>	0.39	78.5	<b>71.2</b>	<b>0.76</b>	66.5	<b>64.0</b>	<b>0.55</b>	95.8	<b>75.3</b>	<b>0.61</b>
Efficient Net-B0	<i>Vanilla</i>	78.2	44.3	0.90	84.2	62.9	1.00	73.9	53.2	0.99	96.6	58.3	0.25
	RRR	78.4	49.6	0.79	83.1	68.7	0.85	73.9	59.1	0.66	95.4	75.6	<b>0.50</b>
	P-CIArC	65.2	35.1	0.02	19.7	29.6	1.00	74.1	54.6	0.21	96.8	55.0	0.05
	A-CIArC	78.0	54.2	0.64	77.7	72.8	0.68	71.4	69.9	0.90	96.7	60.6	0.24
	RR-CIArC (ours)	77.6	<b>70.3</b>	<b>0.53</b>	78.7	<b>75.6</b>	<b>0.54</b>	73.9	<b>70.8</b>	<b>0.56</b>	92.0	<b>77.6</b>	0.43

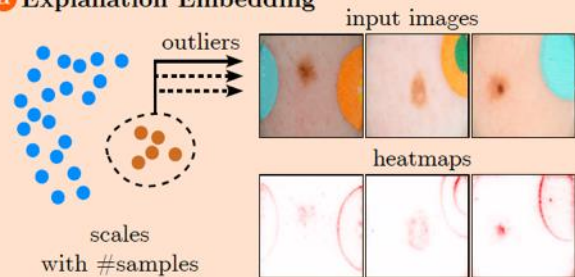
Reveal & Revise Example



# Explain and Improve

## 1 Identification of Model Weakness

### a Explanation Embedding



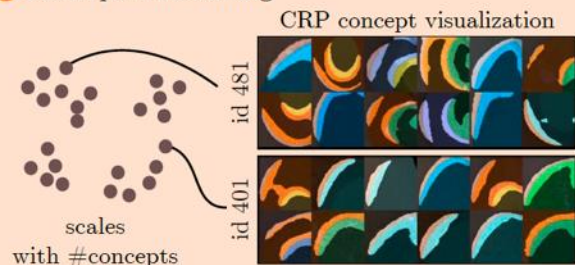
### Analysis in input space

Cluster heatmaps with SpRAY (Lapusckin et al. 2019)

Present cluster representative to human

Identify artifact clusters

### b Concept Embedding



### Analysis in concept space

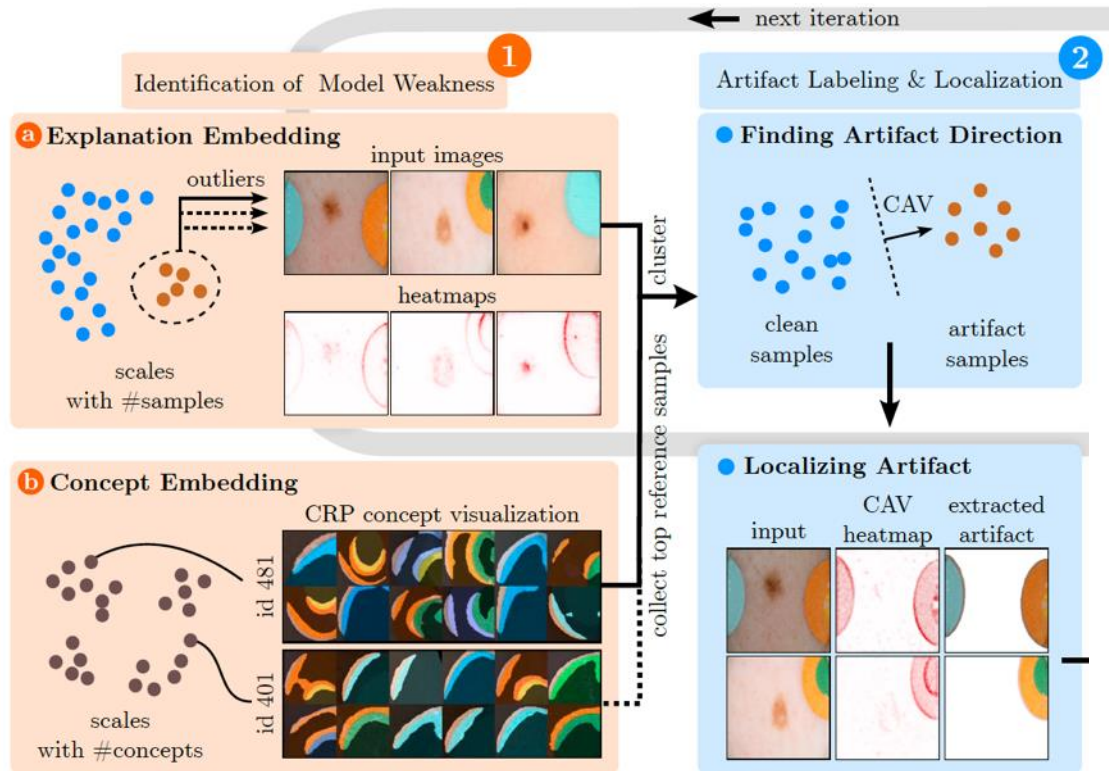
Present concepts (or concept clusters) to human

Identify artifact concepts / clusters

[Pahde et al. 2023]



# Explain and Improve



Compute "artifact CAV" h separating clean and outlier samples in latent space.

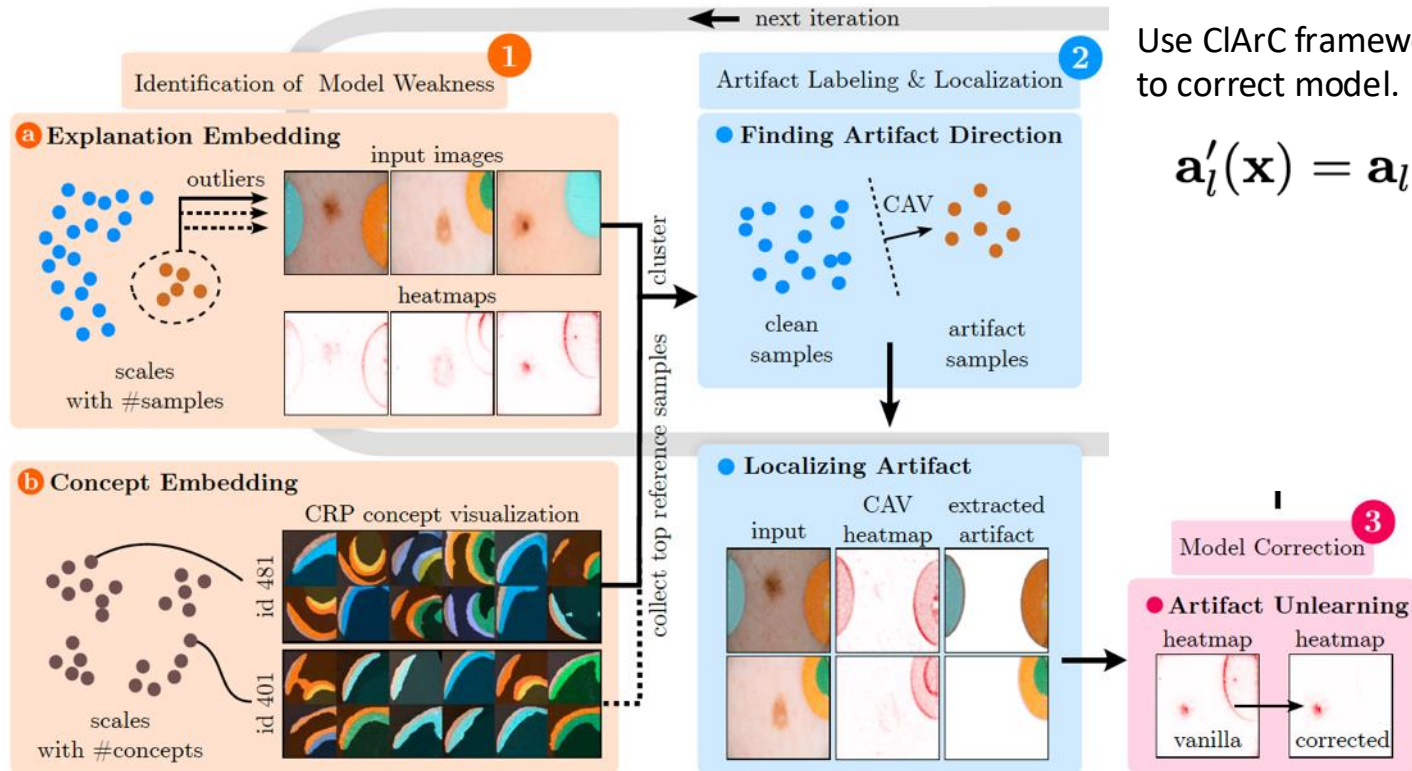
Localize the "artifact CAV" in input space (use LRP rule).

$$\mathbf{R}_l(\mathbf{x}) = \mathbf{a}_l(\mathbf{x}) \circ \mathbf{h}_l$$

→ explain artifact in input space using LRP.

[Pahde et al. 2023]

# Explain and Improve



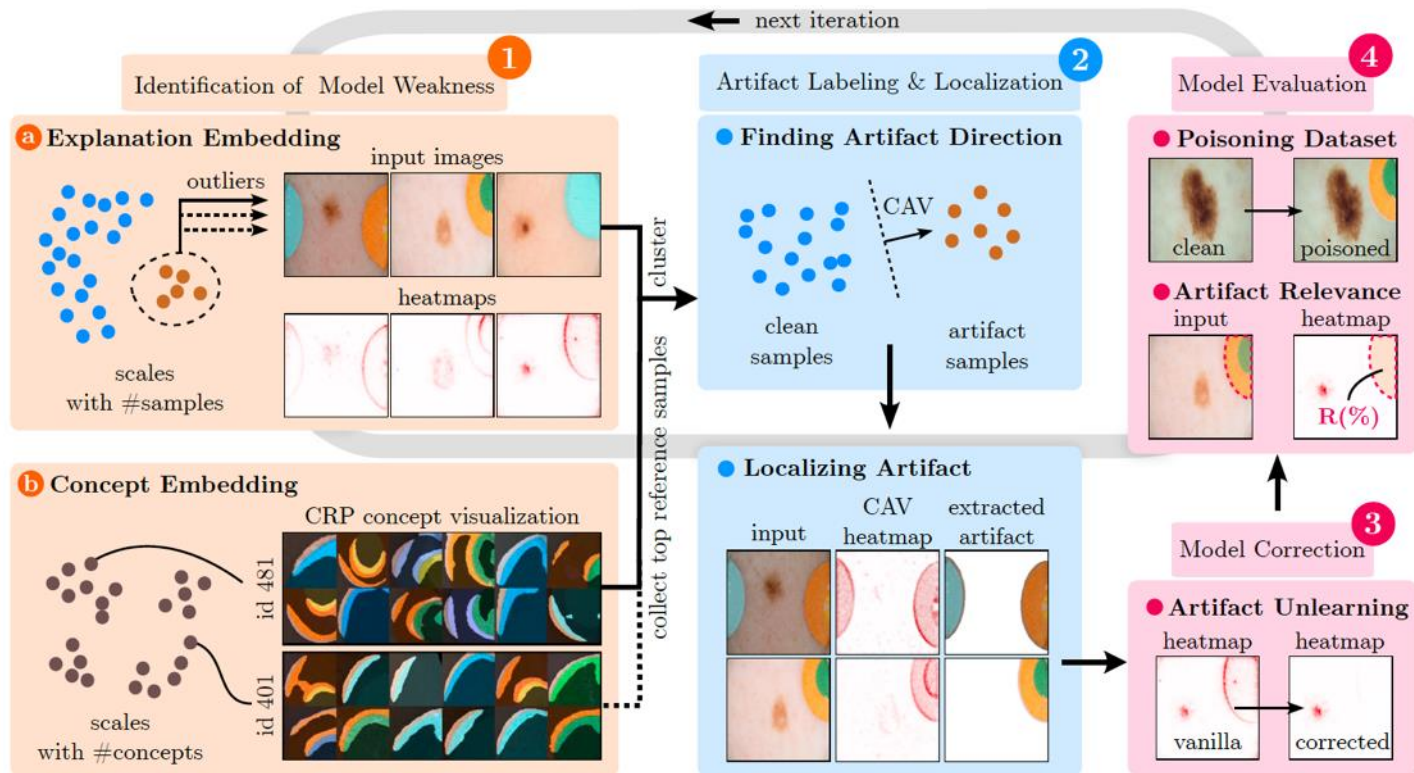
## Latent space correction

Use CIaRC framework (Anders et al. 2022) to correct model.

$$\mathbf{a}'_l(\mathbf{x}) = \mathbf{a}_l(\mathbf{x}) + \gamma(\mathbf{x})\mathbf{h}_l$$

[Pahde et al. 2023]

# Explain and Improve

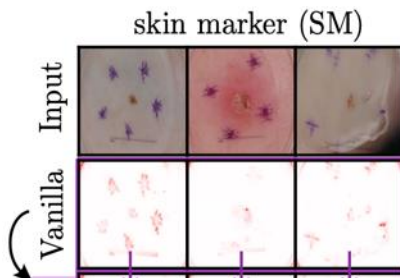


[Pahde et al. 2023]

# Example

ISIC Dataset Artifacts:

- Skin Marker



--- Reveal Step ---

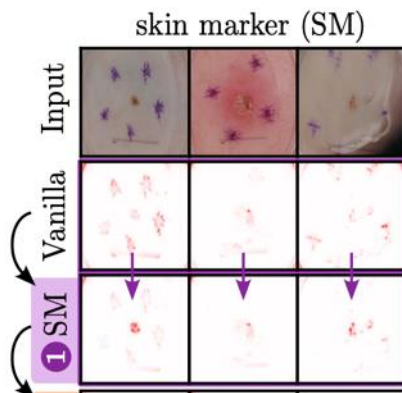
R2R iteration	corrected artifacts	↓ artifact relevance (%)			↑ F1 (%)			↑ accuracy (%)				
					<i>poisoned</i>	<i>original</i>		<i>poisoned</i>	<i>original</i>			
0	-	18.4	45.5	24.2	61.3	59.7	60.5	73.9	71.8	71.5	68.7	<b>80.1</b>

[Pahde et al. 2023]

# Example

ISIC Dataset Artifacts:

- Skin Marker



--- Revise Step ---

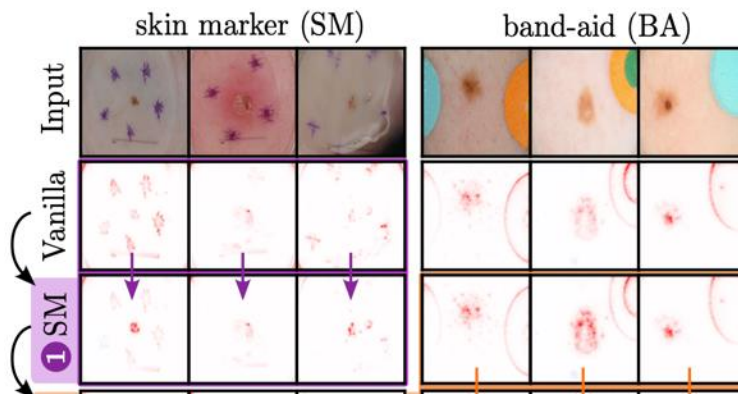
R2R iteration	corrected artifacts	↓ artifact relevance (%)			↑ F1 (%)			↑ accuracy (%)				
					<i>poisoned</i>	<i>original</i>		<i>poisoned</i>	<i>original</i>			
0	-	18.4	45.5	24.2	61.3	59.7	60.5	73.9	71.8	71.5	68.7	<b>80.1</b>
1	SM	13.1	35.0	21.3	61.6	61.0	60.7	73.8	72.2	72.6	68.4	80.0

[Pahde et al. 2023]

# Example

ISIC Dataset Artifacts:

- Skin Marker
- Band-Aid



--- Reveal Step ---

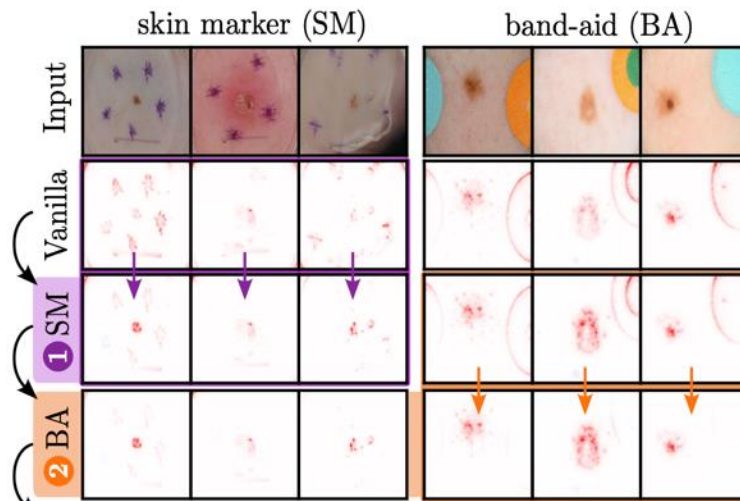
R2R iteration	corrected artifacts	↓ artifact relevance (%)			↑ F1 (%)			↑ accuracy (%)				
					<i>poisoned</i>	<i>original</i>		<i>poisoned</i>	<i>original</i>			
0	-	18.4	45.5	24.2	61.3	59.7	60.5	73.9	71.8	71.5	68.7	<b>80.1</b>
1	SM	13.1	35.0	21.3	61.6	61.0	60.7	73.8	72.2	72.6	68.4	80.0

[Pahde et al. 2023]

# Example

ISIC Dataset Artifacts:

- Skin Marker
- Band-Aid



--- Revise Step ---

R2R iteration	corrected artifacts	↓ artifact relevance (%)			↑ F1 (%)			↑ accuracy (%)				
					<i>poisoned</i>	<i>original</i>		<i>poisoned</i>	<i>original</i>			
0	-	18.4	45.5	24.2	61.3	59.7	60.5	73.9	71.8	71.5	68.7	<b>80.1</b>
1	SM	13.1	35.0	21.3	61.6	61.0	60.7	73.8	72.2	72.6	68.4	80.0
2	SM, BA	<b>12.8</b>	16.8	16.8	61.5	<b>63.6</b>	61.1	73.9	72.3	<b>74.6</b>	68.6	79.7

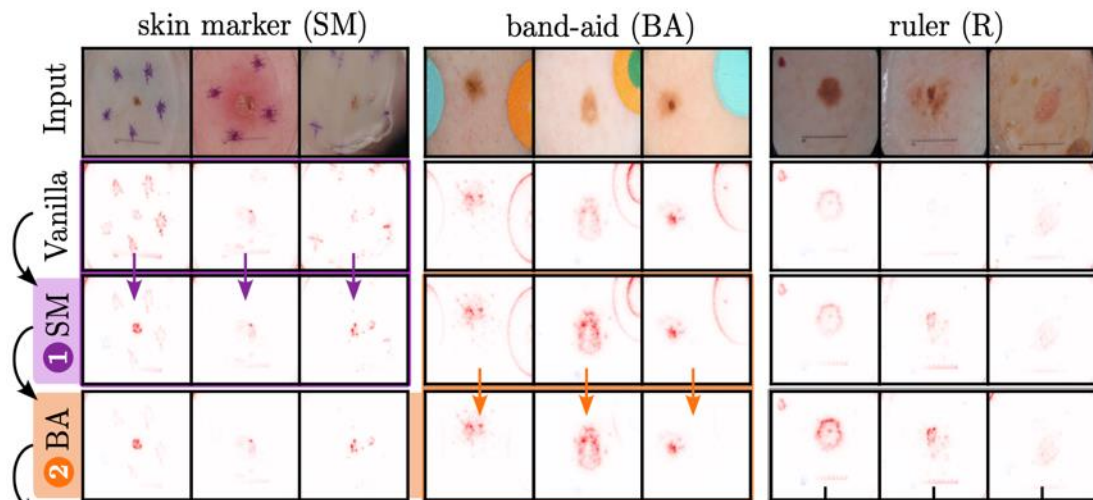
[Pahde et al. 2023]



# Example

ISIC Dataset Artifacts:

- Skin Marker
- Band-Aid
- Ruler



--- Reveal Step ---

R2R iteration	corrected artifacts	↓ artifact relevance (%)			↑ F1 (%)			↑ accuracy (%)				
					<i>poisoned</i>	<i>original</i>		<i>poisoned</i>	<i>original</i>			
0	-	18.4	45.5	24.2	61.3	59.7	60.5	73.9	71.8	71.5	68.7	<b>80.1</b>
1	SM	13.1	35.0	21.3	61.6	61.0	60.7	73.8	72.2	72.6	68.4	80.0
2	SM, BA	<b>12.8</b>	16.8	16.8	61.5	<b>63.6</b>	61.1	73.9	72.3	<b>74.6</b>	68.6	79.7

[Pahde et al. 2023]

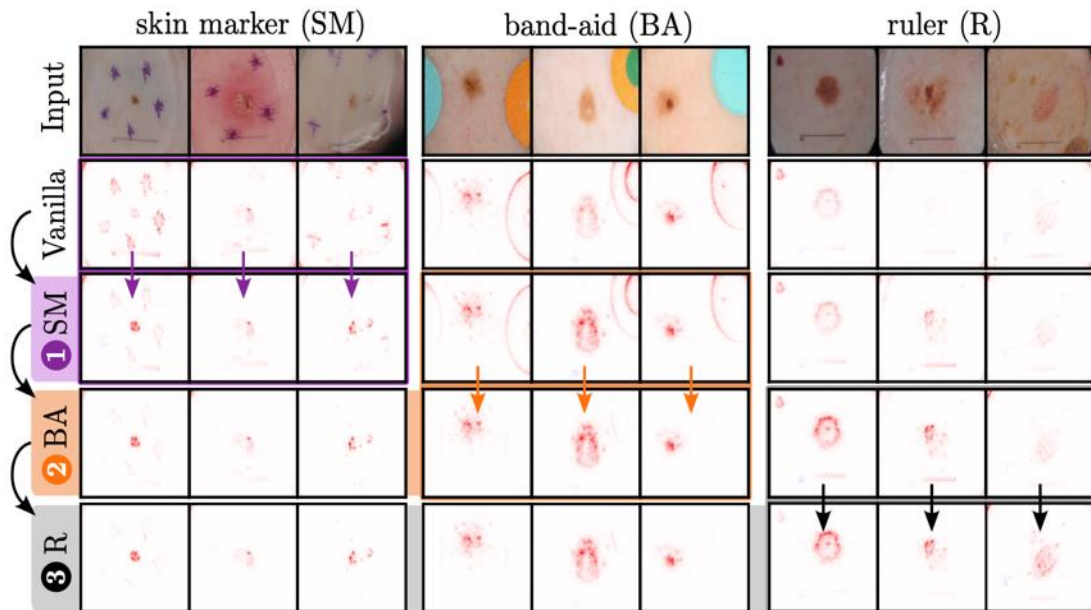


# Example

ISIC Dataset Artifacts:

- Skin Marker
- Band-Aid
- Ruler

--- Revise Step ---

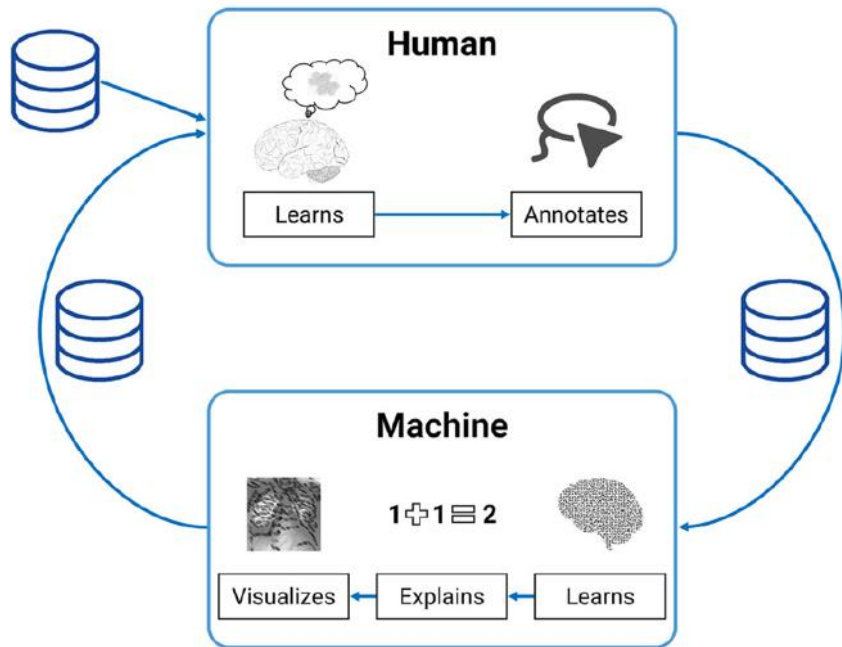


R2R iteration	corrected artifacts	↓ artifact relevance (%)			↑ F1 (%)			↑ accuracy (%)				
					<i>poisoned</i>	<i>original</i>		<i>poisoned</i>	<i>original</i>			
0	-	18.4	45.5	24.2	61.3	59.7	60.5	73.9	71.8	71.5	68.7	<b>80.1</b>
1	SM	13.1	35.0	21.3	61.6	61.0	60.7	73.8	72.2	72.6	68.4	80.0
2	SM, BA	<b>12.8</b>	16.8	16.8	61.5	<b>63.6</b>	61.1	73.9	72.3	<b>74.6</b>	68.6	79.7
3	SM, BA, R	14.6	<b>15.7</b>	<b>8.5</b>	<b>62.0</b>	63.4	<b>64.0</b>	<b>74.0</b>	<b>72.4</b>	74.5	<b>71.8</b>	79.9

[Pahde et al. 2023]

Explanatory Interactive ML

# Explanatory Interactive ML (XIL)



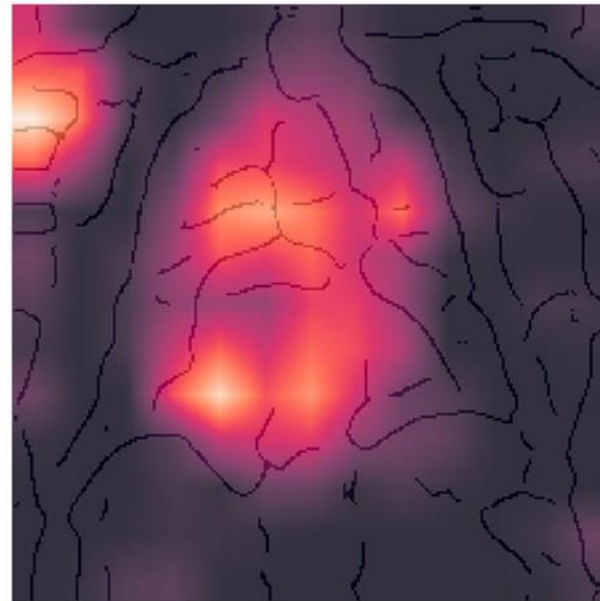
*“In XIL, a learner can interactively query the user (or some other information source) to obtain the desired outputs of the data points. The interaction takes the following form. At each step, the learner considers a data point (labeled or unlabeled), predicts a label, and provides explanations of its prediction. The user responds by correcting the learner if necessary, providing a slightly improved – but not necessarily optimal – feedback to the learner.”*

[Pfeuffer et al. 2023]

# Example

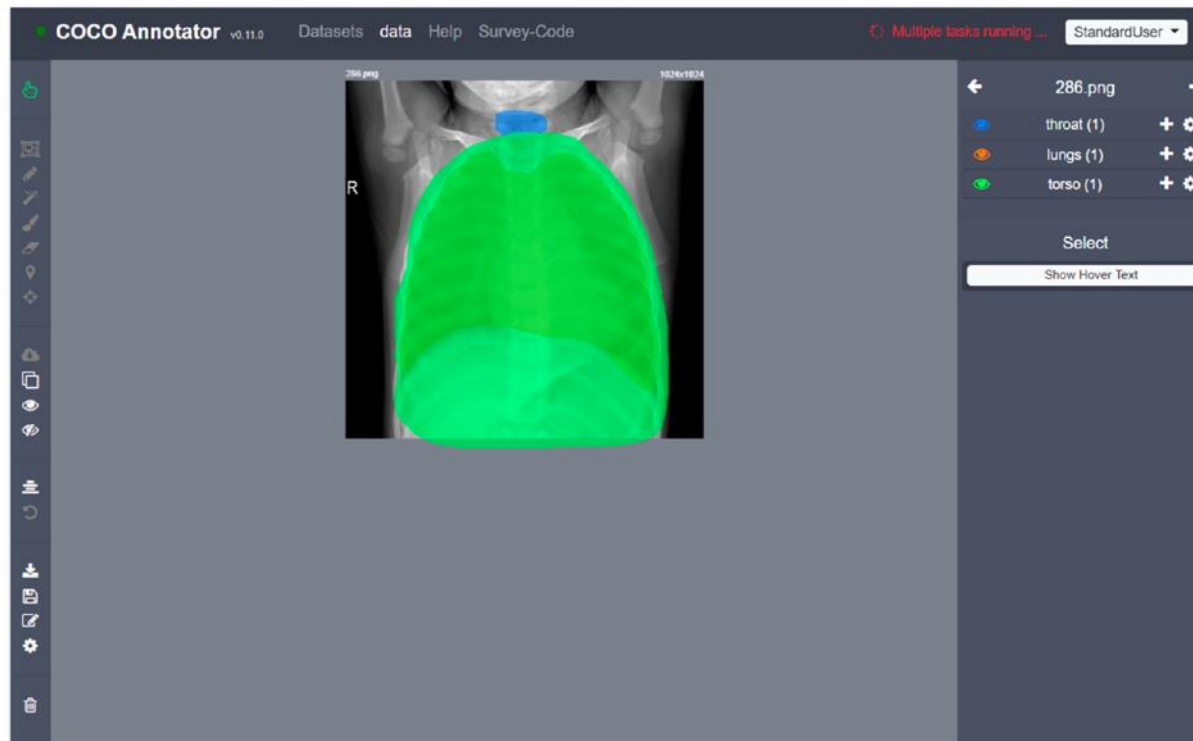
## 1. Step: Reveal

**Fig. 6** Loop 1 XAI. Original image on the left, XAI showing confounders with Grad-CAM method (Selvaraju et al. 2017) on an overlay of the edge-filtered image on the right (color figure online)



# Example

## 2. Step: Data curation

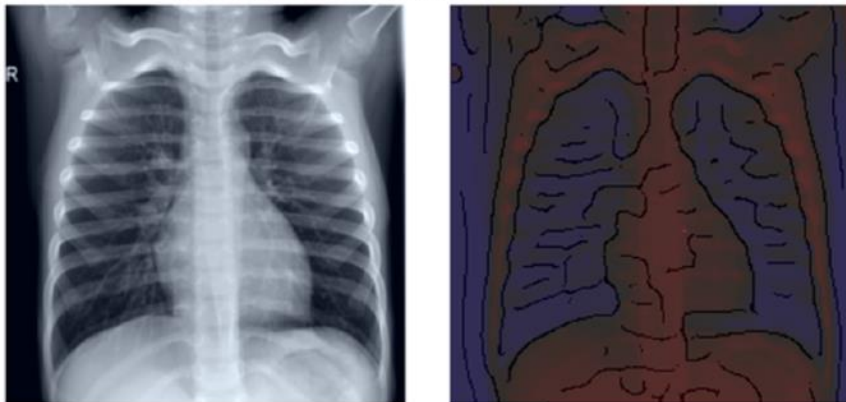


**Fig. 7** The user interface of the customized COCO Annotator used by the crowdworkers to annotate the images. Using the tools on the left on each X-ray, the areas for the throat, lungs, and torso were highlighted (color figure online)

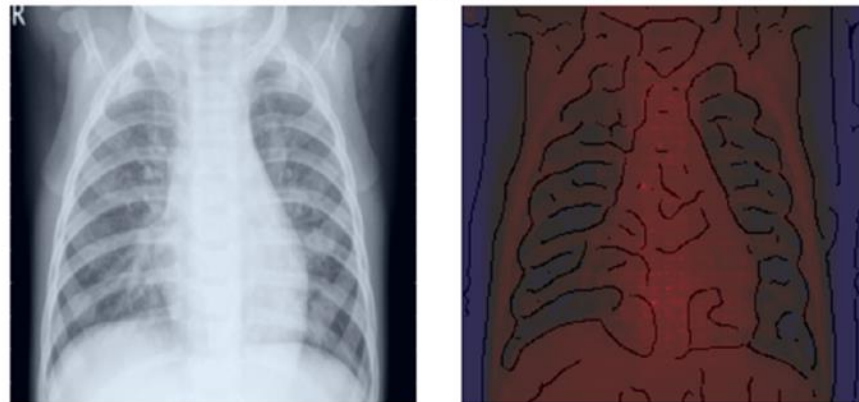
# Example

## 3. Step: Refine data and retrain model

True Class: Normal; Pred Class: Normal



True Class: Viral Pneumonia; Pred Class: Viral Pneumonia



**Fig. 9** XAI showing Explanations for Classification in loop 3a (using our customized red-blue filter solution based on VarGrad (Adebayo et al. 2018)) (color figure online)

# 10th Take Home Message

*Integrating the human into XAI / XIL is crucial.*

*Concept-Level explanations can be a useful interface here (especially for non-image data).*

Future of XAI



# Next-Generation Explanations

contextual and semantically rich:

--> explain whole inference (where, what, how)

adequate for receiver:

--> presented in comprehensible form (also for time series etc.)

actionable:

--> useful beyond visualization, e.g., for debugging, improving or auditing model.

of various types:

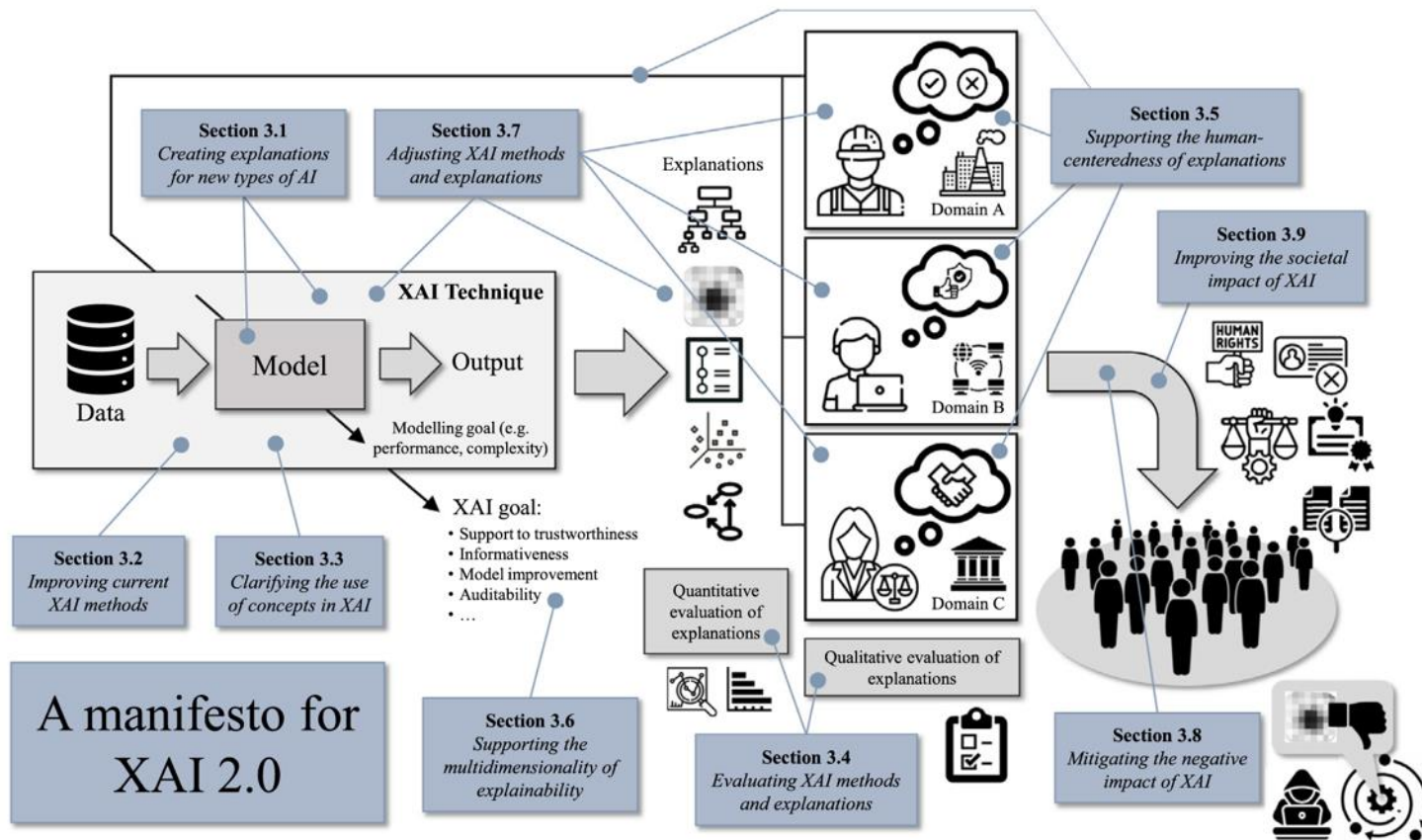
--> broader definition of explainability / transparency

# Next-Generation Explanations

## Explainable Artificial Intelligence (XAI) 2.0: A manifesto of open challenges and interdisciplinary research directions

Luca Longo <sup>1,2,\*</sup>, Mario Brcic <sup>3</sup>, Federico Cabitza <sup>4,5</sup>, Jaesik Choi <sup>6,7</sup>, Roberto Confalonieri <sup>8</sup>,  
Javier Del Ser <sup>9,10,11</sup>, Riccardo Guidotti <sup>12</sup>, Yoichi Hayashi <sup>13</sup>, Francisco Herrera <sup>11</sup>,  
Andreas Holzinger <sup>14</sup>, Richard Jiang <sup>15</sup>, Hassan Khosravi <sup>16</sup>, Freddy Lecue <sup>17</sup>,  
Gianclaudio Malgieri <sup>18</sup>, Andrés Páez <sup>19,20</sup>, Wojciech Samek <sup>21,22,23</sup>, Johannes Schneider <sup>24</sup>,  
Timo Speith <sup>25,26</sup>, Simone Stumpf <sup>27</sup>

# Next-Generation Explanations



# Thank you for your attention

<http://www.heatmapping.org>

- ▶ Tutorials
- ▶ Software
- ▶ Online Demos

