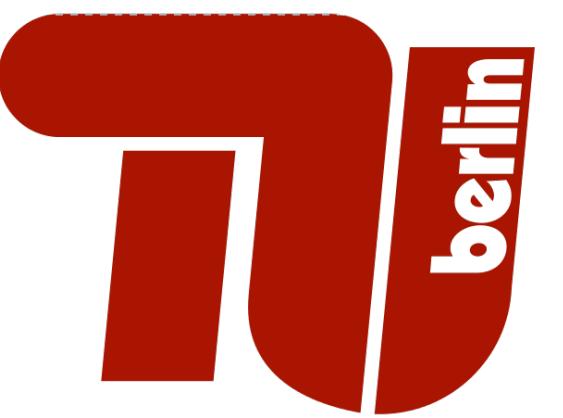


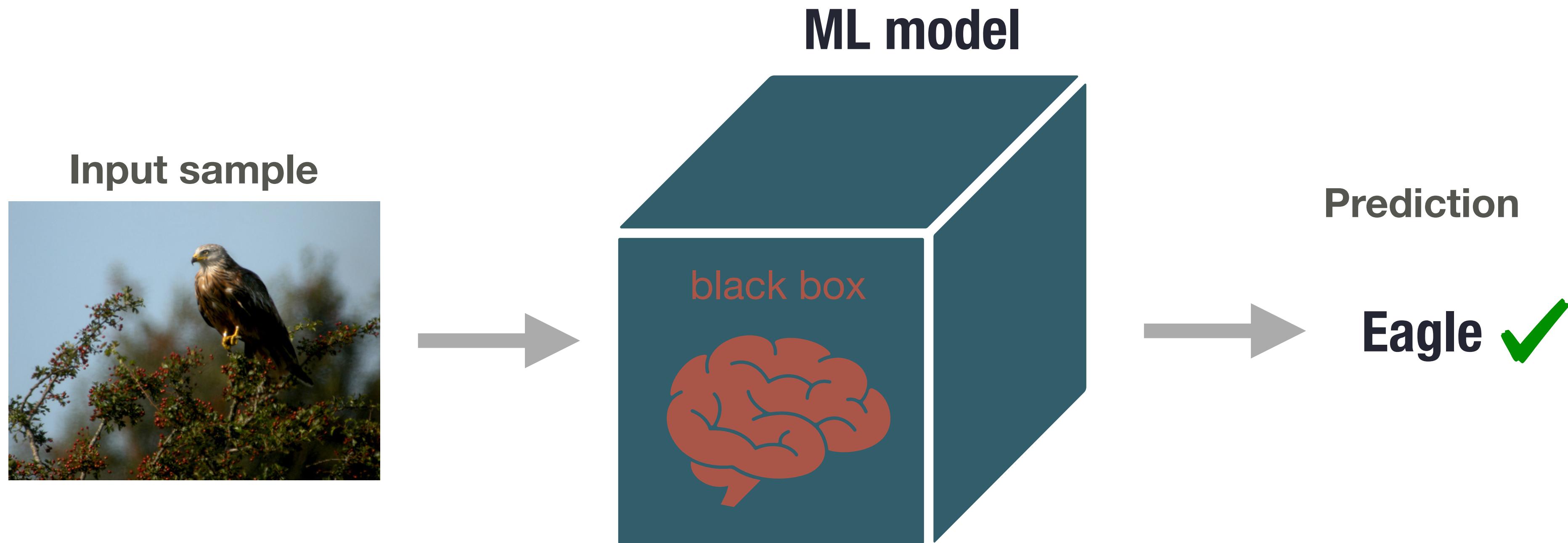
PredDiff: Explanation and Interactions from Conditional Expectations

Stefan Blücher (PhD @ BIFOLD @ Prof. Müller)
Joint work with Johanna Vielhaben and Nils Strodthoff





Warm-up – XAI review



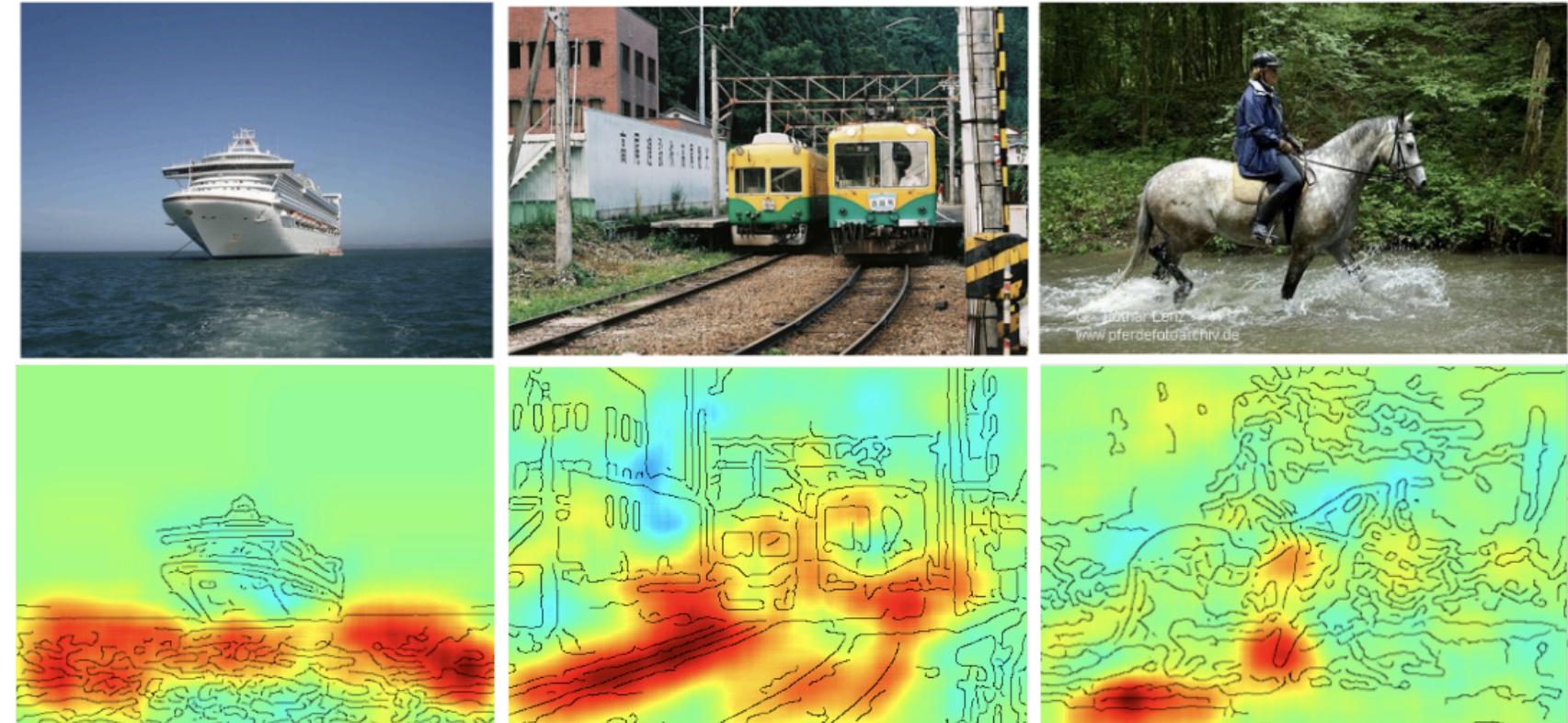
**How does the model arrive at this decision?
What can we learn from this?**

Explainable AI (XAI) – overview



Local → Global

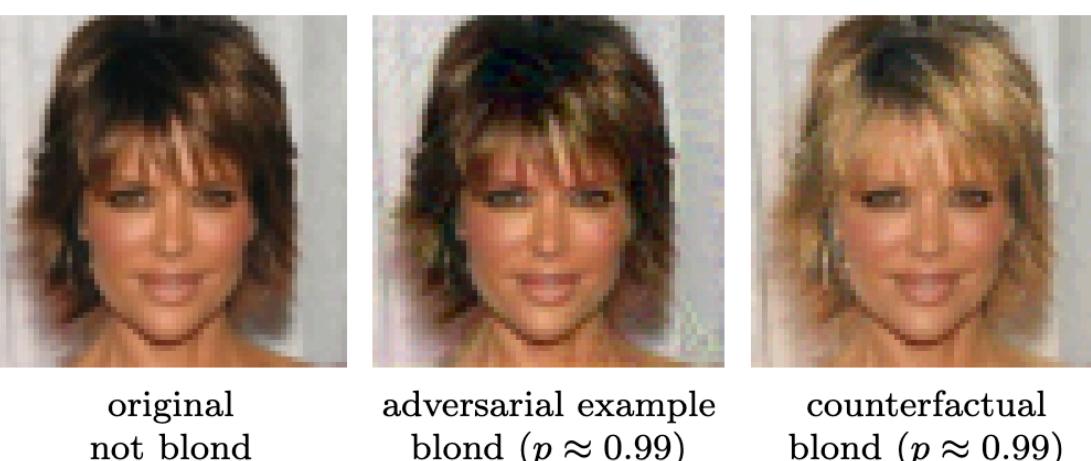
Heatmaps



Lapuschkin, et al., *Nature communications* 10.1 (2019)
"Unmasking Clever Hans predictors and assessing what machines really learn."

Relevant input regions?

Counterfactuals



Dombrowski, et al., *arXiv:2206.05075* (2022).
"Diffeomorphic Counterfactuals with Generative Models."

Which input changes prediction?

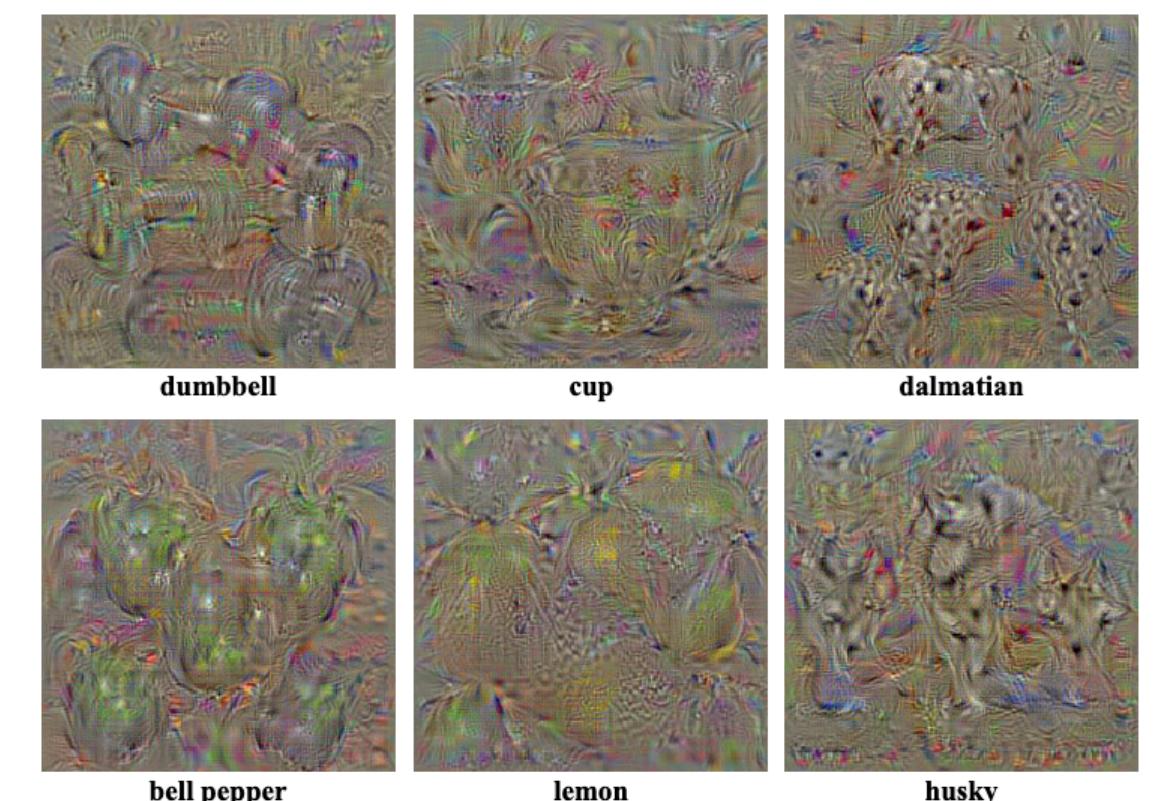
Discover concepts



Vielhaben, et al., *arXiv:2203.06043* (2022).
"Sparse Subspace Clustering for Concept Discovery (SSCCD)."

Investigate hidden model representations

Activation maximization



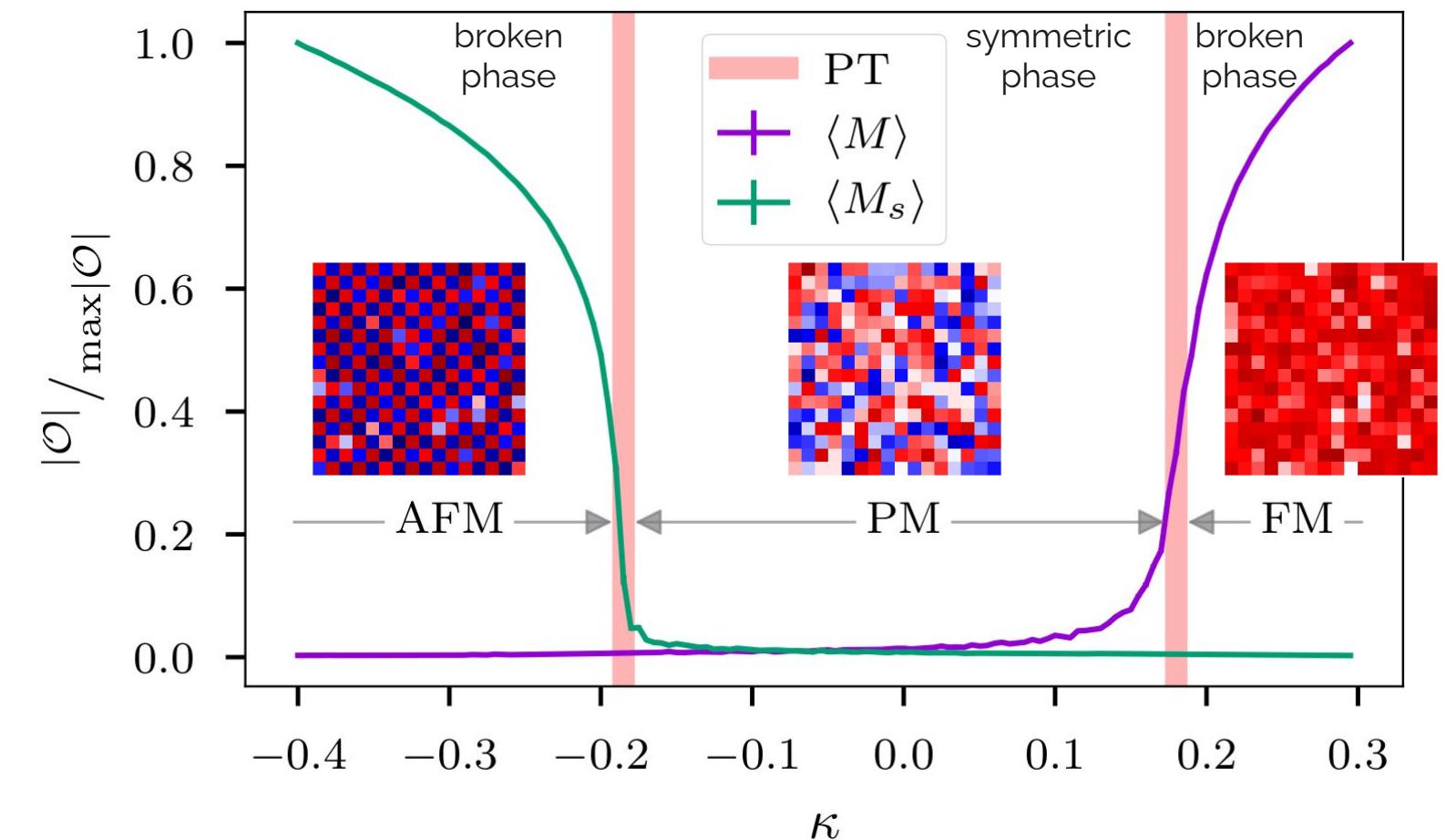
Simonyan, et al., *arXiv:1312.6034* (2013).
"Deep inside convolutional networks: Visualising image classification models and saliency maps."

Visualize class prototypes

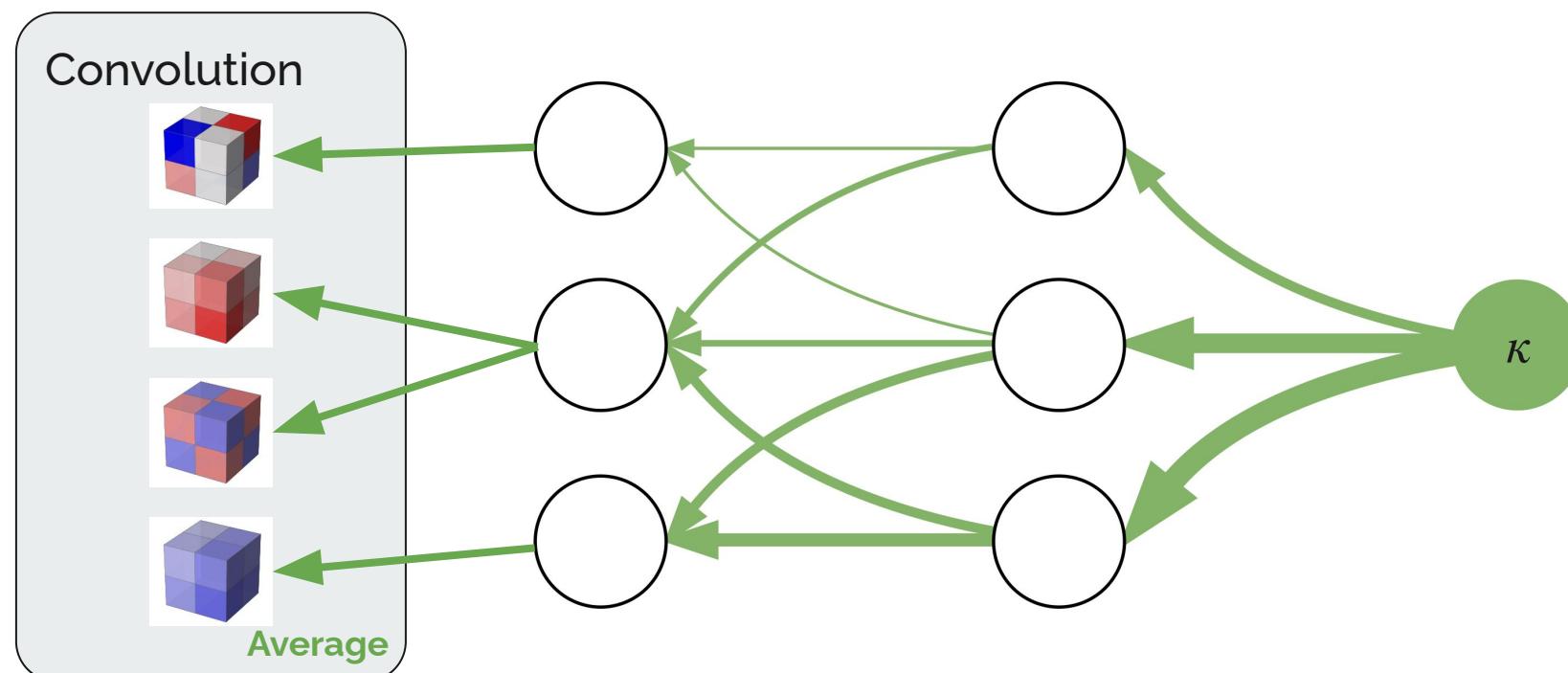
Use-case: scientific insights with XAI



Yukawa phase diagram - fixed fermion coupling

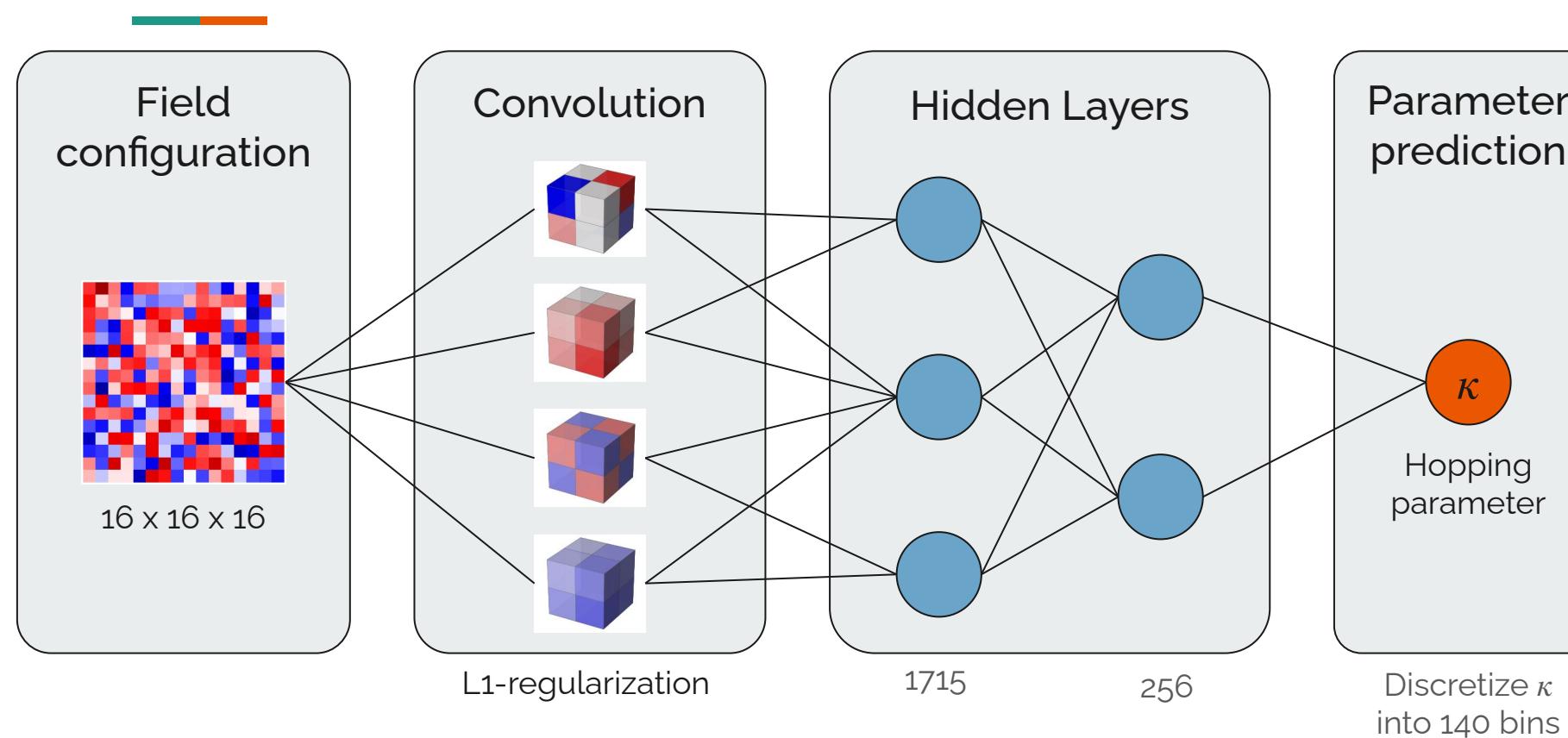


Filter relevance - averaging conv layer



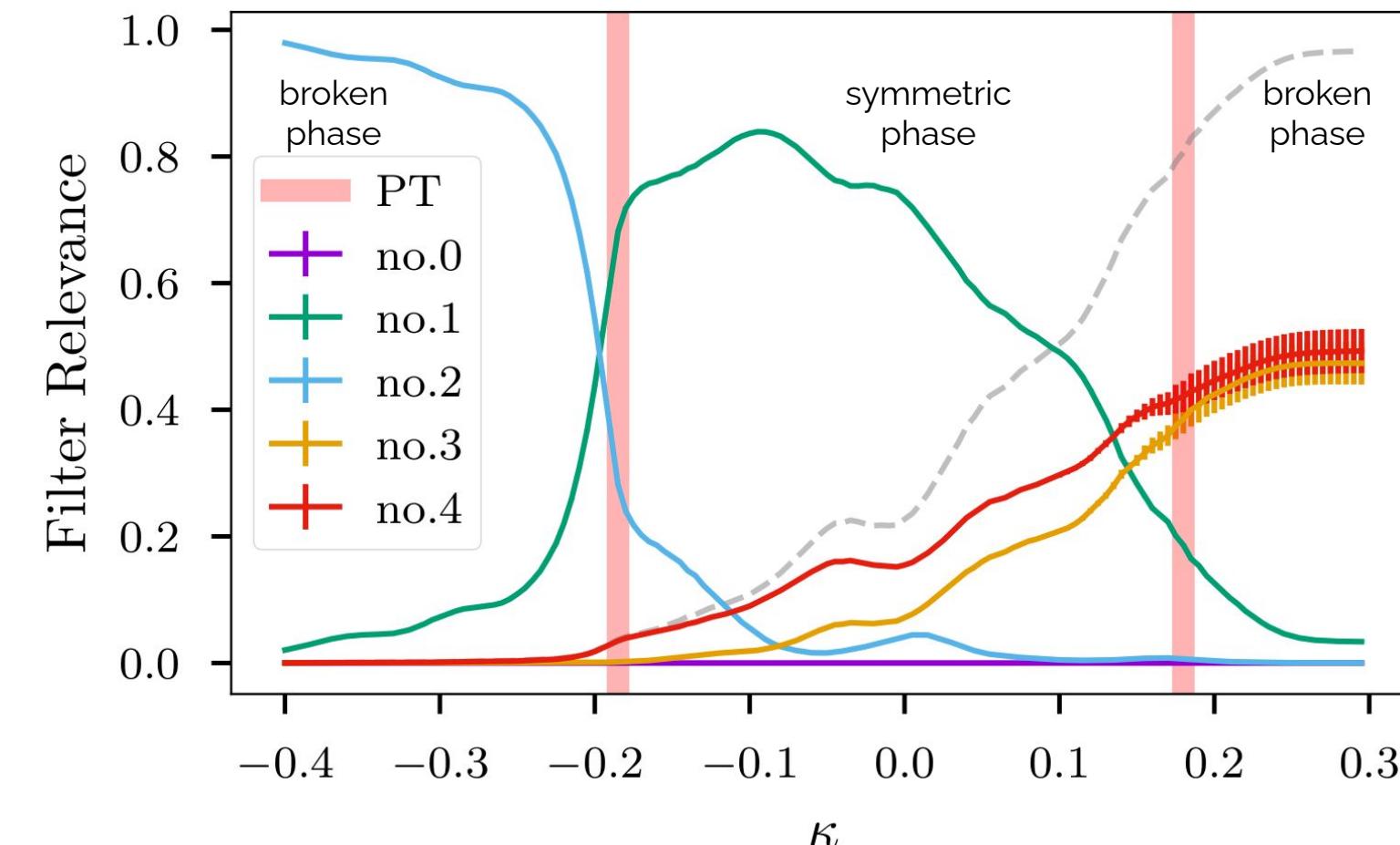
Graphic adapted from heatmapping.org
G. Montavano, W. Samek, K.-R. Müller Digital Signal Processing 73 (2018) 1–15

Network architecture - shallow convolution



Analysis
using XAI

Filter relevances - discovering structures



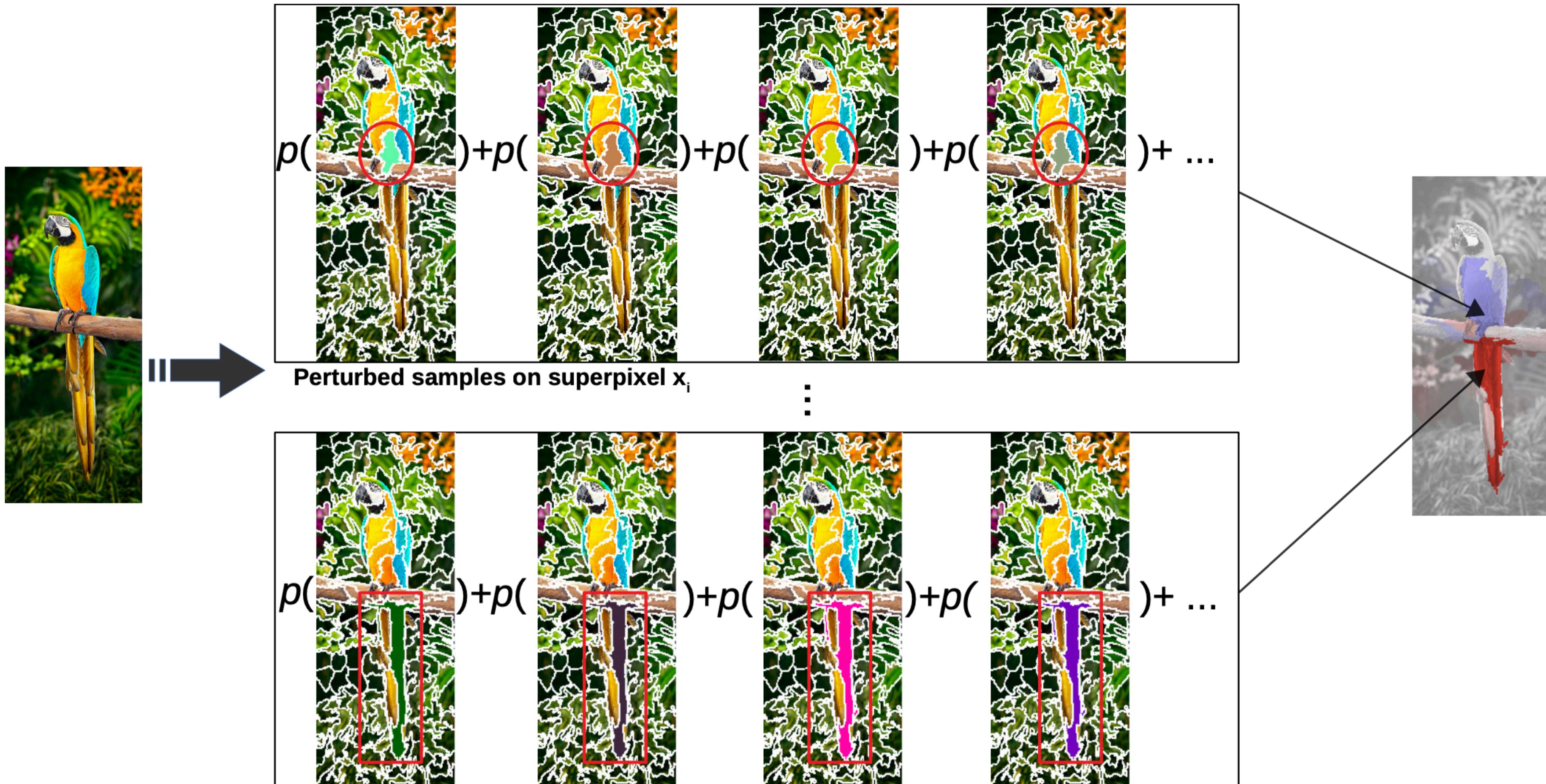


Local model-agnostic attributions

Explaining by removing

Prediction Difference Analysis (*PredDiff*)

Original – Perturbations = Attribution



Prediction Difference Analysis (*PredDiff*)



- Simple attribution framework
- Firmly rooted in probability theory
- Only occlude target feature Y
 - > numerically cheap
 - > on-manifold

Occluded prediction $\rightarrow m_{Y|x_i}^f = \int dY f(x, Y) p(Y|x)$

Model \downarrow

Imputer \downarrow

Regression

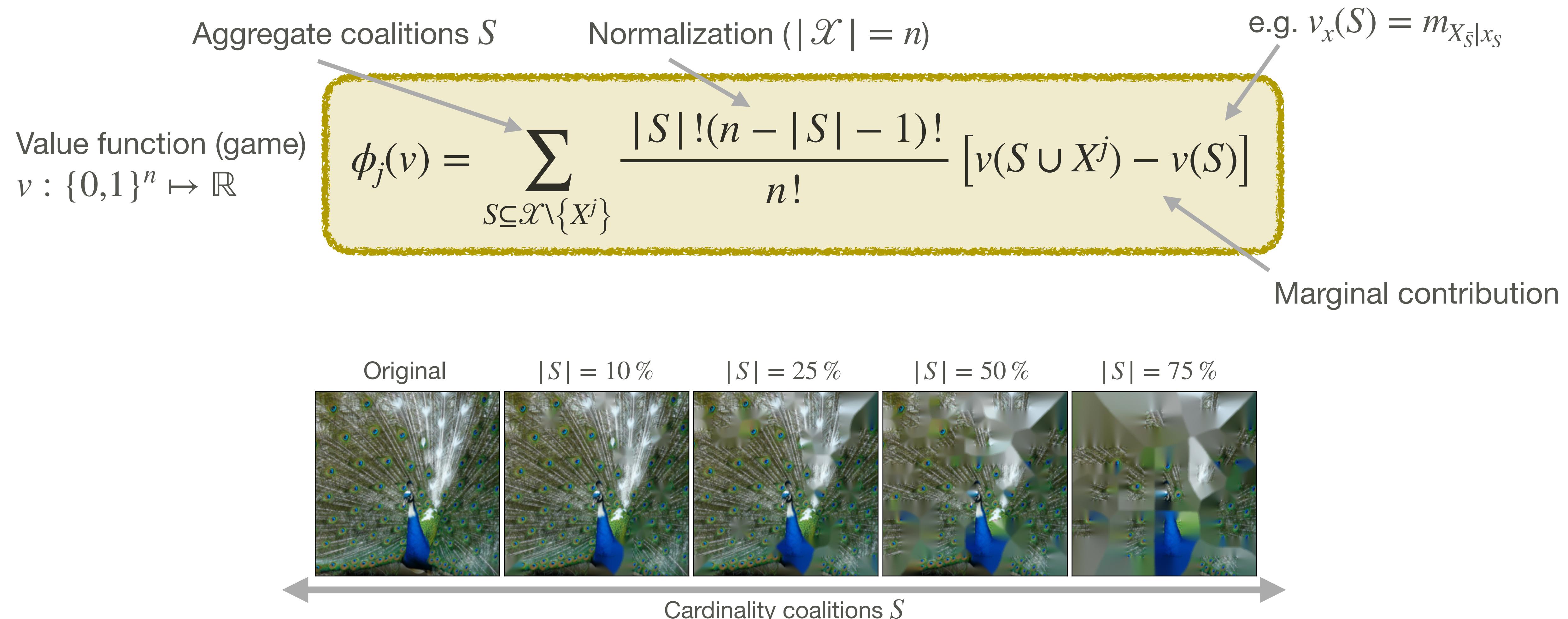
$$\bar{m}_{Y|x_i}^f := f(x_i, y_i) - m_{Y|x_i}^f$$

Classification

$$\bar{m}_{Y|x_i}^{f_c} := \log_2 f_c(x_i, y_i) - \log_2 m_{Y|x_i}^{f_c}$$

Shapley values – a „*n*-player game“

Marginalize spectator features



Shapley values – desirable axioms

Efficiency

$$\sum_i \phi_i = v(\mathcal{X}) - v(\emptyset)$$

Symmetry

$$\forall S \subseteq \mathcal{X} \setminus \{i, j\} : v(S \cup i) = v(S \cup j)$$

$$\Rightarrow \phi_i = \phi_j$$

Dummy player

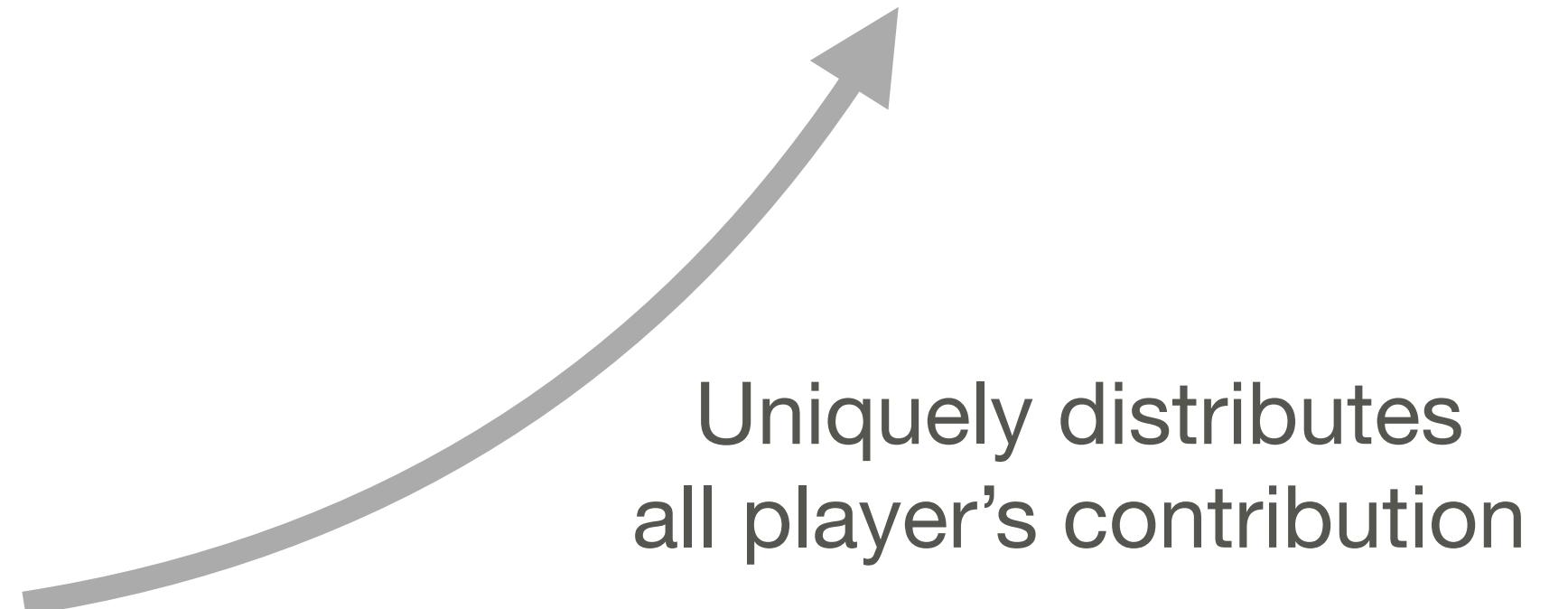
$$\forall S \subseteq \mathcal{X} \setminus i : v(S \cup i) = v(S)$$

$$\Rightarrow \phi_i = 0$$

Linearity

$$\phi_i(v + \alpha w) = \phi_i(v) + \alpha \phi(w)$$

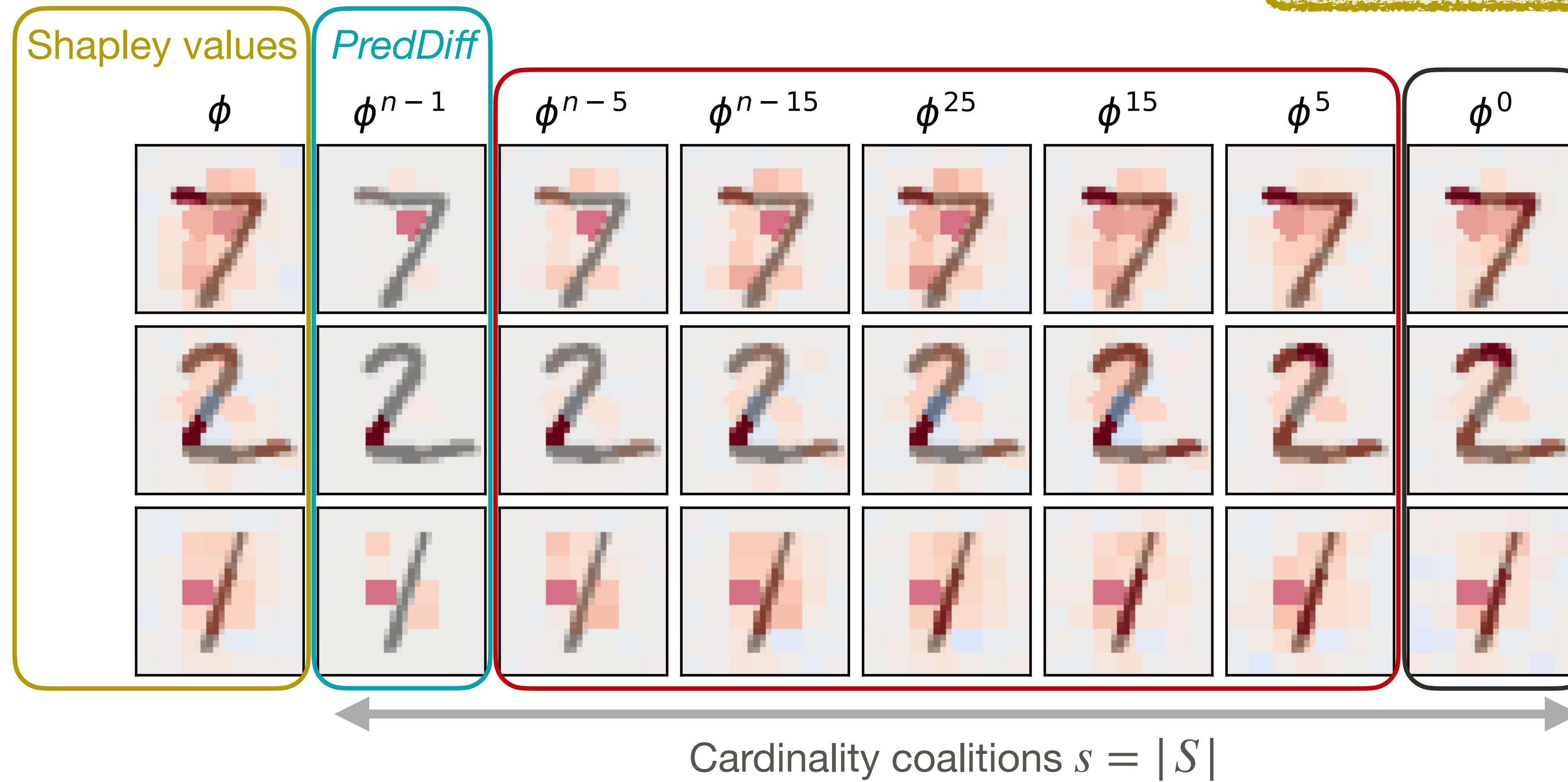
$$\phi_j(v) = \sum_{S \subseteq \mathcal{X} \setminus \{X^j\}} \frac{|S|!(n - |S| - 1)!}{n!} [v(S \cup X^j) - v(S)]$$



Shapley values – marginal contributions

$$\phi_j(v) = \sum_{S \subseteq \mathcal{X} \setminus \{X^j\}} \frac{|S|!(n - |S| - 1)!}{n!} [v(S \cup X^j) - v(S)]$$

ϕ^s



„no-interaction“ property

Regression

$$f(X, Y) = Y^2 + 3Z$$

Classification

Conditional informative interaction

$$p(y, z | c, x) = p(y | c, x) p(z | c, x)$$

Factorizing imputer distribution

$$\bar{m}_{YZ|x}^{f^{YZ}} = 0$$

Interaction measure

Decompose model prediction

$$f^V(V) = \sum_{W \subseteq V} (-1)^{|V|-|W|} P_{\mathcal{X} \setminus W} f(X, Y, Z)$$

for example $P_Y f(X, Y, Z) = f(X, y, Z)$

Two-point interaction

$$\bar{m}_{YZ|x}^{f^{YZ}} = \bar{m}_{YZ|x}^f - \bar{m}_{Y|x}^{f^Y} - \bar{m}_{Z|x}^{f^Z}$$

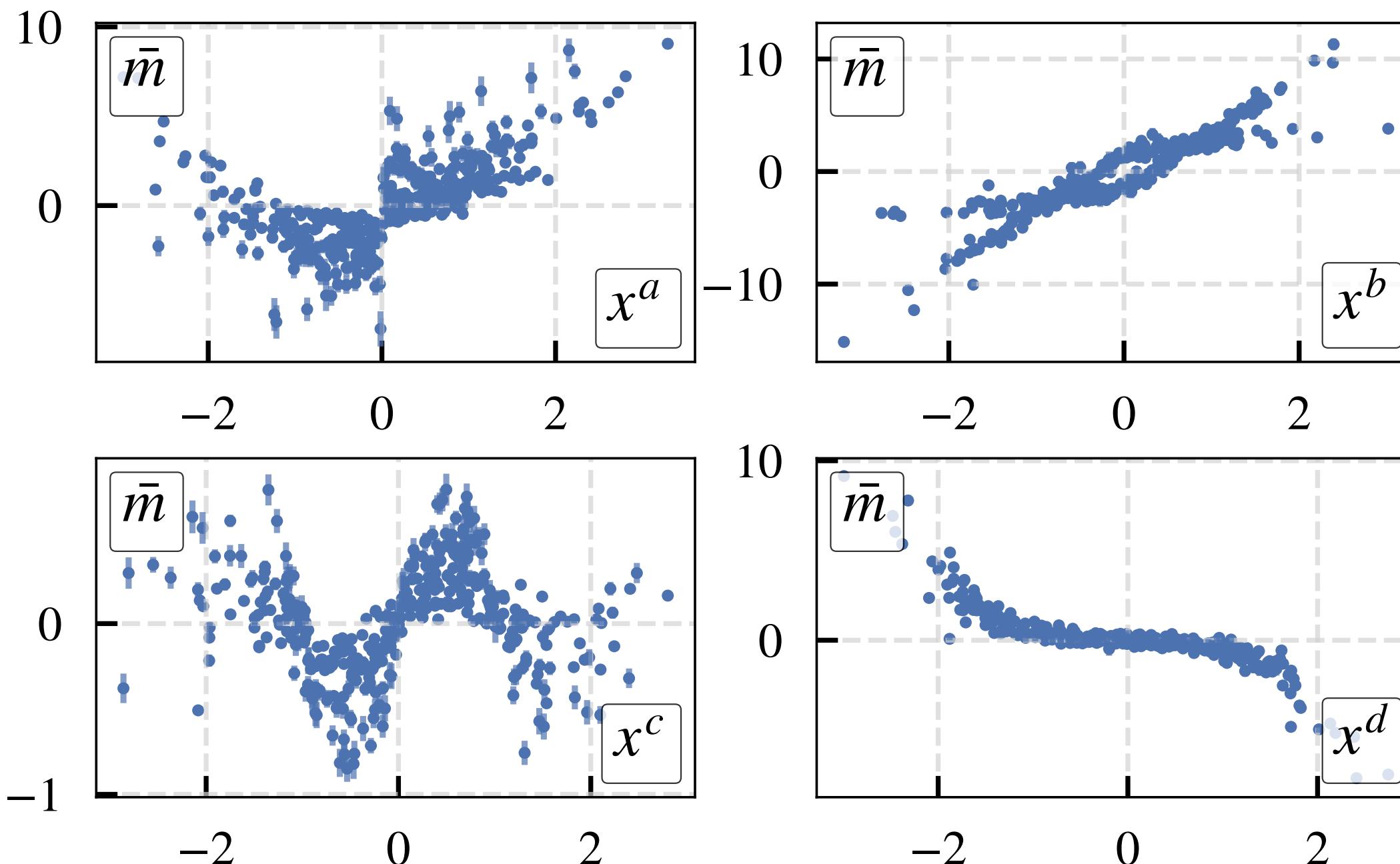


Experiments

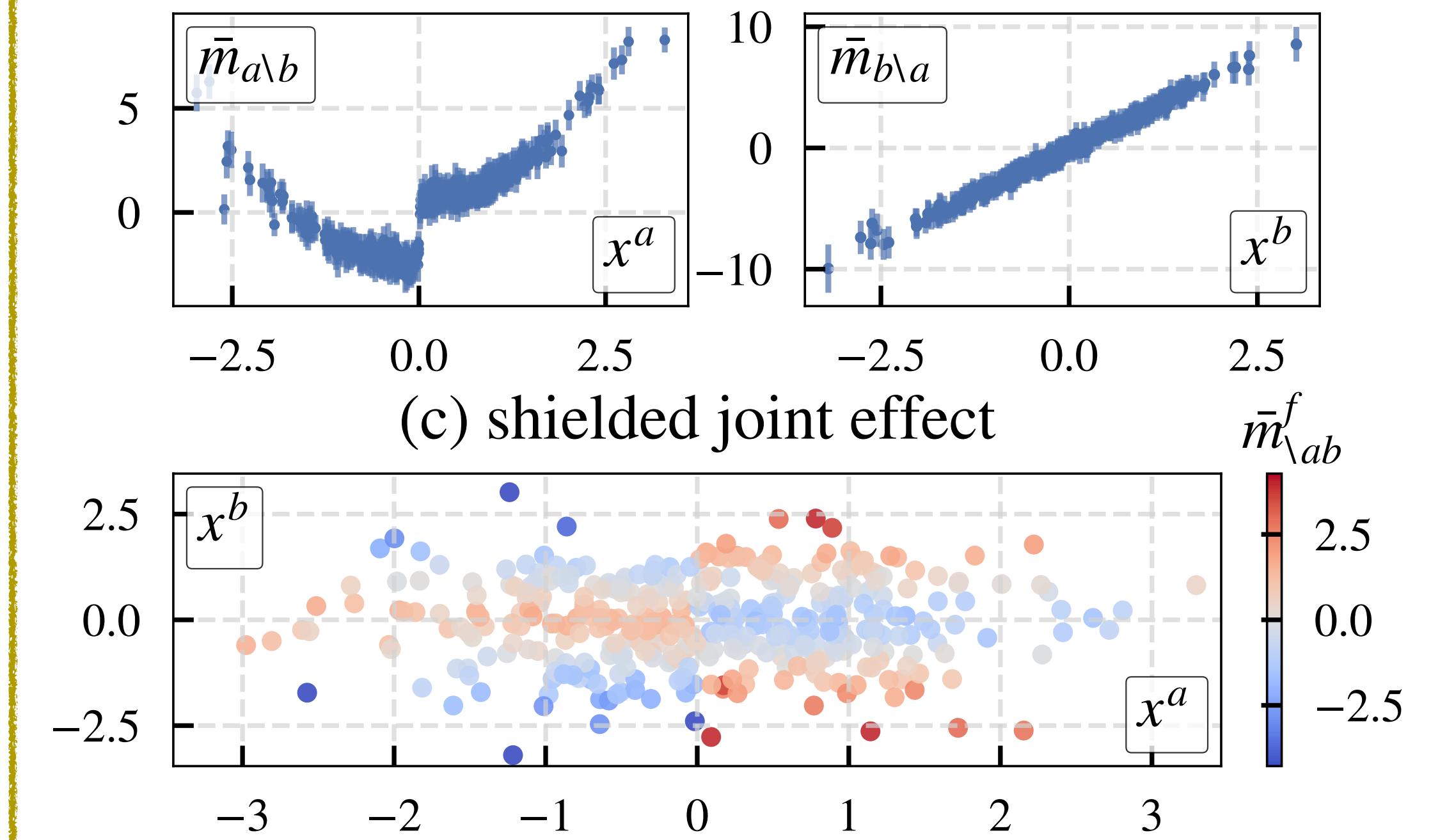
Recover additive regression components

$$f(x) = x_a^2 + 3x_b + \sin(\pi x_c) - \frac{x_d^3}{2} + 2 \operatorname{sgn}(x_a) \operatorname{abs}(x_b)$$

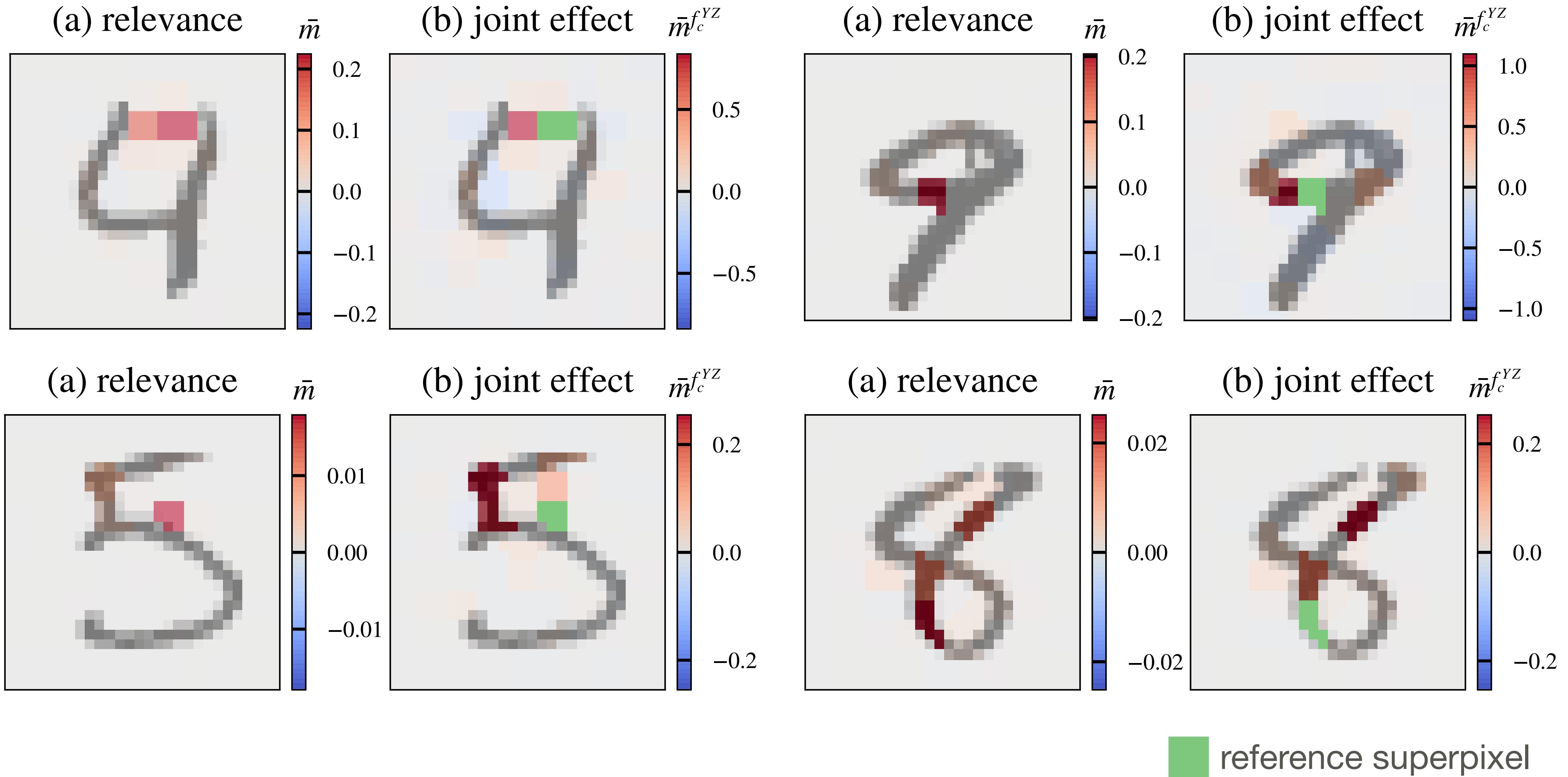
Main effects



Interactions



Meaningful interactions for MNIST digits



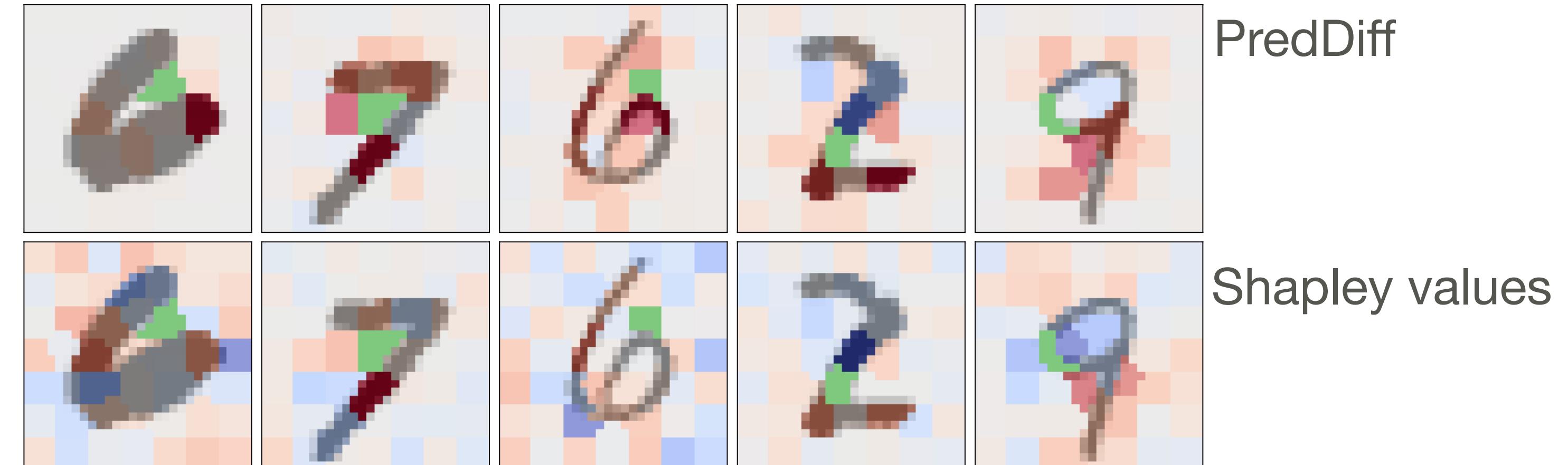
Similar attributions for *PredDiff* & Shapley



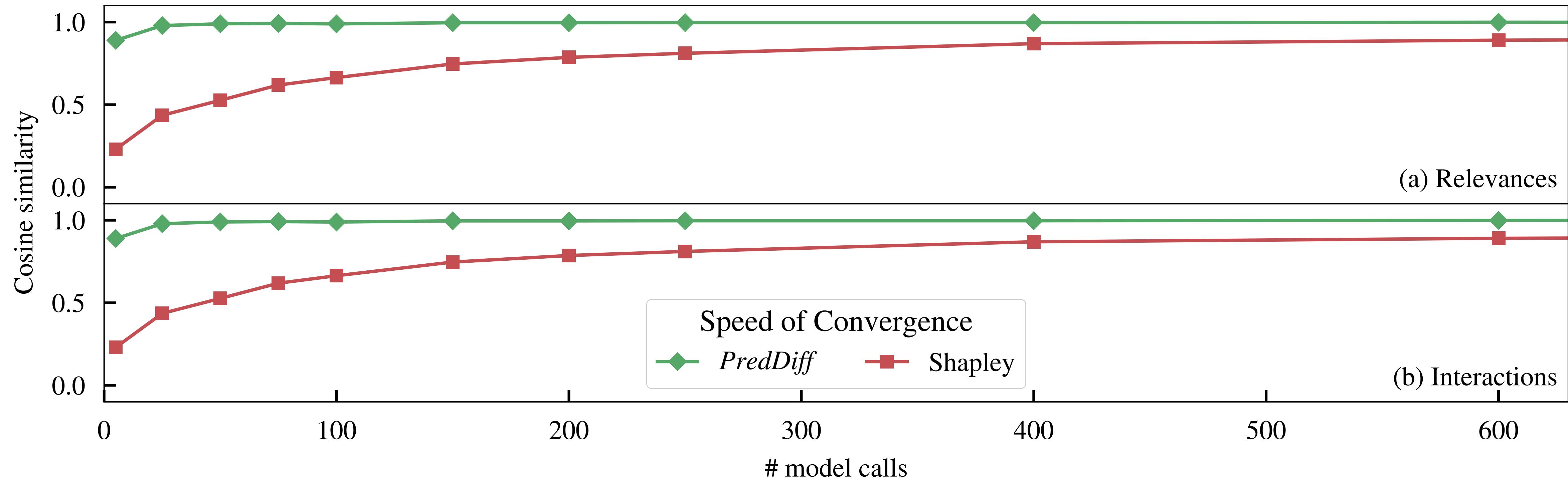
Relevances



Interactions



PredDiff: rapid convergence & scalability

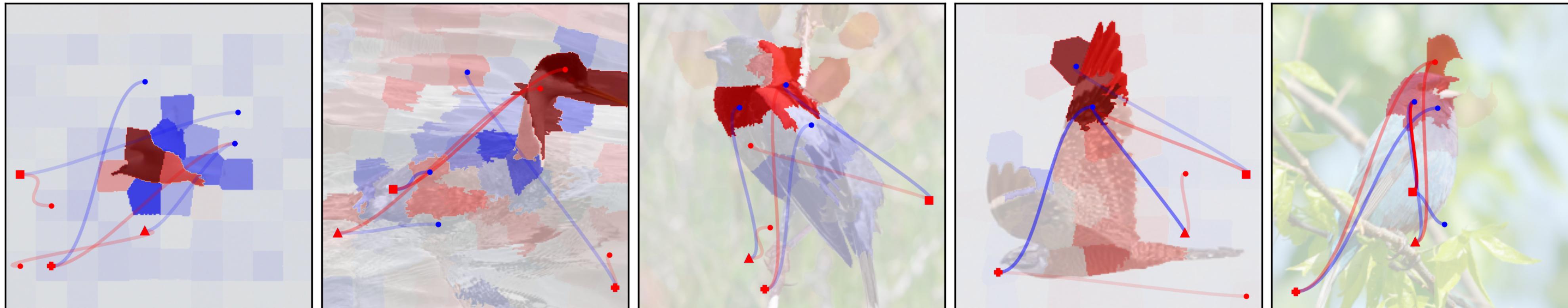


Insight: relevant features show interactions

Highest/Lowest relevance reference super pixels



Random reference super pixels



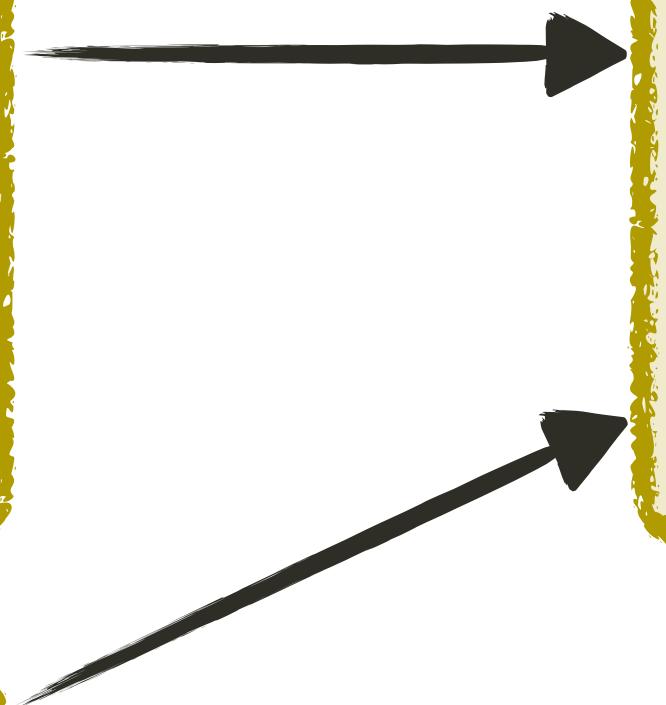


Take-home messages

Take-home messages

PredDiff

- Conceptual simplicity
- On-manifold imputations
- only target feature is occluded
- Linear scaling
- Provably *no-interaction* property



Higher-order attributions

NEW

- More detailed insights through interaction effects
- Allows to understand complex vision and text classifiers
- Software package to appear

Shapley values

- Very popular in the literature
- Axiomatic interpretation
- Many different variants proposed

Useful resources

- heatmapping.org
- <https://github.com/chr5tphr/zennit>
- <https://captum.ai/>
- <https://github.com/slundberg/shap>