

Benchmarking the Attribution Quality of Vision Models

Robin Hesse¹

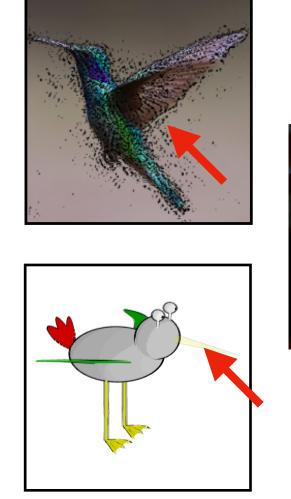
Simone Schaub-Meyer^{1,2}

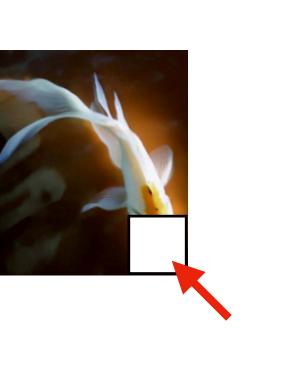
The chicken and egg problem ...of evaluating XAI methods We need XAI methods to know how deep models are working To To To

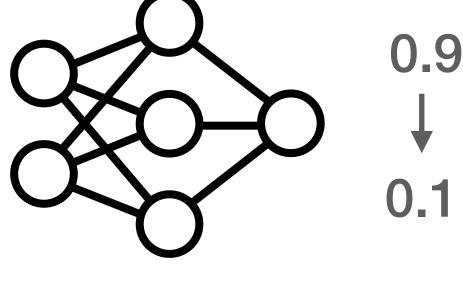
Related work

...and its limitations

Common protocols to evaluate attribution quality







Goldfish

→ Delete patches, pixels, or concepts to measure the effect on the output confidence or accuracy

Protocol	Domain alignment	No inf. leakage	Inter-model comp.	Natural data
Incremental deletion	X	N/A	X	\checkmark
Single deletion	X	N/A	\checkmark	\checkmark
ROAR	\checkmark	X	X	\checkmark
FunnyBirds	\checkmark	\checkmark	\checkmark	X
IDSDS (ours)	\checkmark	\checkmark	\checkmark	\checkmark

Stefan Roth^{1,2}

