

# **CLEVR-XAI:** A Benchmark Dataset for the **Ground Truth Evaluation of Neural Network** Explanations

Leila Arras, Ahmed Osman, Wojciech Samek

### Goal of eXplainable Artificial Intelligence (XAI) MOTIVATION



Blackbox







Post-hoc XAI methods:



Ground Truth  $\sum$  $R_{p_k}$  and  $R_{total} = \sum R_{p_k}$  $R_{within}$ Mass Accuracy: where  $R_{\scriptscriptstyle within} =$ R<sub>total</sub>  $k=1, s.t. p_k \in Ground Truth$  $|P_{top K} \cap Ground Truth|$ where  $P_{top K} = \{p_1, p_2, ..., p_K | R_{p_1} > R_{p_2} > ... > R_{p_K}\}$ Rank Accuracy: Ground Truth

|Ground Truth| = K : size of the ground truth mask N : size of the image  $R_{p_k}$ : relevance of pixel  $p_k$  after pooling along channels (we use **10 pooling techniques**)

#### Advantages:

selective recognition task (multiple target objects per image, modulated by the question) controlled setup (synthetic images, unbiased and uninformative image background) realistic images (complex questions, objects occluded, shadow, various object attributes) 2 sizes x 3 shapes x 2 materials x 8 colors, random locations, up to 10 objects per image) Note: In our benchmarking we use a **simple** Relation Network to first **validate** XAI methods, but future work could use any other computer vision model that can solve the VQA task.

RESULTS		Comp	oarison	of XAI	meth	ods
<u>Results using</u> CLEVR-XAI simple:	Best mean perfo LRP and Integra (+ LRP lowest v	ormance: ted Gradie ariance)!	Increased XAI accuracy with model confidence		Increased XAI accuracy with target object size	
Relevance Mass Accuracy	GT Single Object all		GT Single Object s.t. proba > 0.99999		GT Single Object s.t. nb pixels > 1000	
	Mean (std)	Median	Mean (std)	Median	Mean (std)	Median
LRP [20]	0.85 (0.17)	0.91	0.90 (0.09)	0.93	0.91 (0.09)	0.94
Excitation Backprop [37]	0.80 (0.20)	0.87	0.85 (0.14)	0.90	0.88 (0.13)	0.92
IG [19]	0.71 (0.27)	0.81	0.75 (0.25)	0.85	0.80 (0.24)	0.90
Guided Backprop [22]	0.58 (0.20)	0.62	0.63 (0.16)	0.66	0.76 (0.13)	0.78
Guided Grad-CAM [39]	0.58 (0.24)	0.63	0.64 (0.21)	0.68	0.83 (0.13)	0.86
SmoothGrad [46]	0.60 (0.33)	0.69	0.61 (0.33)	0.72	0.80 (0.25)	0.91
VarGrad [34]	0.58 (0.34)	0.68	0.60 (0.34)	0.71	0.80 (0.25)	0.91
Gradient [16]	0.49 (0.35)	0.49	0.51 (0.34)	0.53	0.82 (0.25)	0.93
Gradient $\times$ Input [18]	0.43 (0.34)	0.37	0.44 (0.34)	0.39	0.77 (0.27)	0.89
Deconvnet [42]	0.18 (0.16)	0.13	0.18 (0.16)	0.13	0.42 (0.21)	0.41
Grad-CAM [39]	0.09 (0.10)	0.05	0.10 (0.10)	0.07	0.29 (0.10)	0.28
orest performance: econynet and Grad-CAM! Highest variance: all gradient-based XAI methods!					More results and discussion in our paper	

## SUMMARY

## Outlook

- First objective ground truth and VQA based evaluation framework for XAI in computer vision
- Systematic comparison of 10+ popular XAI methods using novel quantitative metrics
- Surprising findings regarding the strengths and limitations of current XAI methods
- Dataset, code and framework are publicly available for future XAI research

• Open questions: Can an XAI method reach the 1.0 relevance accuracy upper bound? Do we need further model constraints to achieve this?

Leila Arras, Ahmed Osman, Wojciech Samek

**Fraunhofer** BIFOLD Heinrich-Hertz-Institut

CLEVR-XAI: A benchmark dataset for the ground truth evaluation of neural network explanations In: Information Fusion, Vol. 81, Pages 14-40, 2022

Dataset & Code: github.com/ahmedmagdiosman/clevr-xai

Further information: www.heatmapping.org

