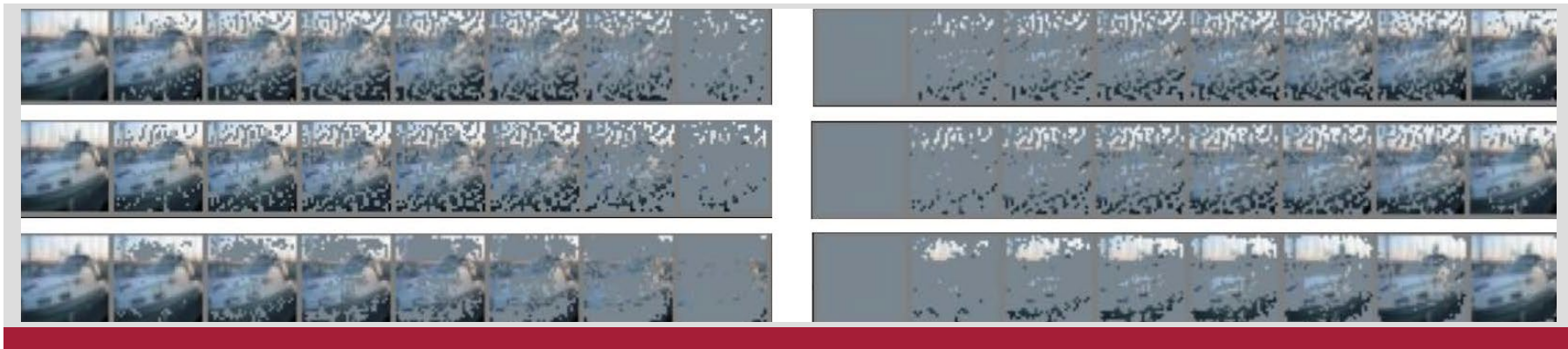




Mathematisch-Naturwissenschaftliche Fakultät, Human-Computer Interaction (HCI) and Data Science & Analytics Research (DSAR)



# A Consistent and Efficient Evaluation Strategy for Attribution Methods

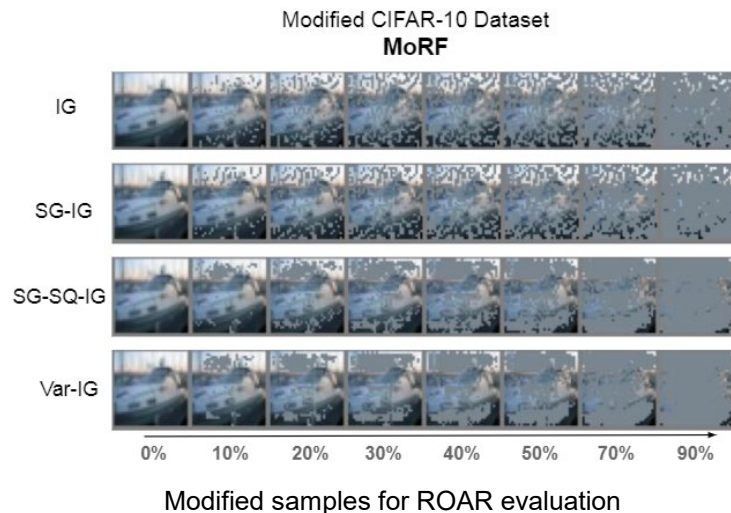
Yao Rong\*, Tobias Leemann\*, Vadim Borisov, Gjergji Kasneci, Enkelejda Kasneci

University of Tübingen, \*equal contribution



## Introduction

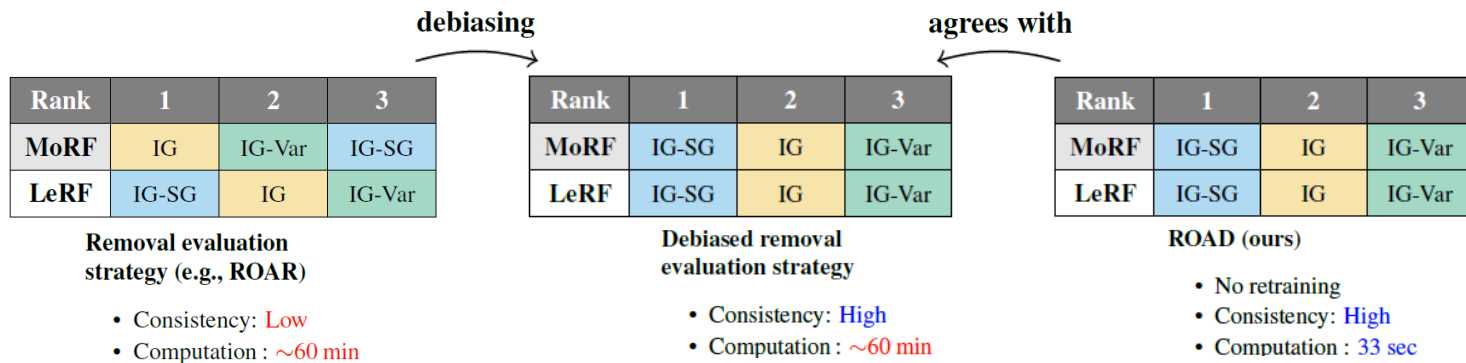
- *Feature attribution methods* are highly popular explanation techniques
- Need for quantitative evaluation strategies that assess faithfulness



[1] S. Hooker et al.: A Benchmark for Interpretability Methods in Deep Neural Networks. NeurIPS, 2019

## Problems of Existing Evaluation Strategies

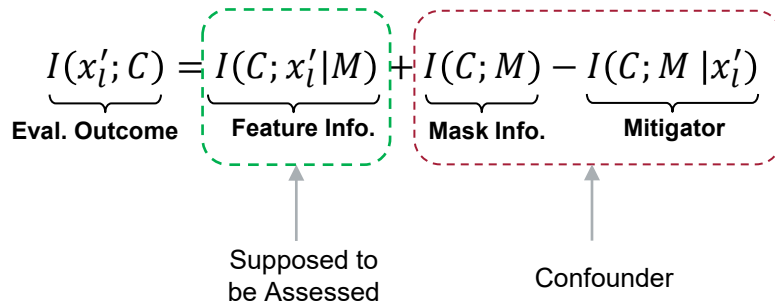
- Inefficiency (retraining step)
- Inconsistency



- We propose RemOve And Debias (ROAD)

# Analysis

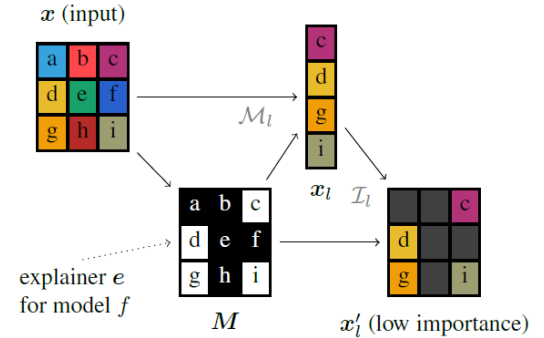
- Analysis from an information-theoretic perspective:



$C$ : Class information

$M$ : (binary) Mask from attribution methods

$x'_l$ : Imputed input with low importance features



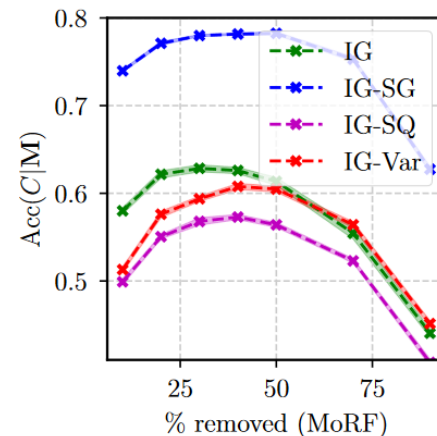
setup of our analysis

## Analysis

- Analysis from an information-theoretic perspective:

$$I(x'_l; C) = \underbrace{I(C; x'_l | M)}_{\text{Feature Info.}} + \underbrace{I(C; M) - I(C; M | x'_l)}_{\text{Mask Info.} \quad \text{Mitigator}}$$

Eval. Outcome      Supposed to be assessed      Confounder



**Class Information Leakage** is significant for real-world data. ←

High accuracy obtained using **only the binary masks** (no values) for class prediction.

## Method

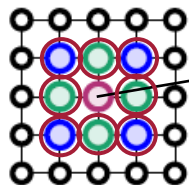
- Minimally Revealing Imputation:  $I(x'_i; M | C) = 0, I(x'_i; M) \approx 0$

$$I(x'_i; C) = \underbrace{I(C; x'_i | M)}_{\text{Feature Info.}} + \underbrace{I(C; M)}_{\text{Mask Info.}} - \underbrace{I(C; M | x'_i)}_{\text{Mitigator}}$$

Eval. Outcome

Stop the class Information Leakage:  $I(C; M) \approx I(C; M | x'_i)$

- Noisy Linear Imputation



$$x_{i,j} = w_d(x_{i,j+1} + x_{i,j-1} + x_{i+1,j} + x_{i-1,j})$$

$$+ w_i(x_{i+1,j+1} + x_{i-1,j+1} + x_{i-1,j-1} + x_{i+1,j-1})$$



Minimally Revealing Imputation:  
**Noisy Linear Imputation**

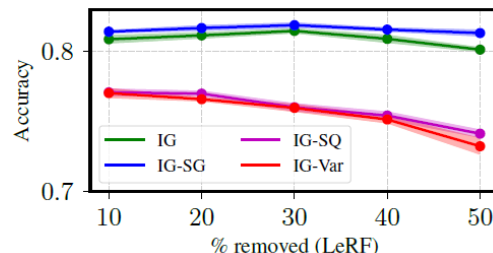
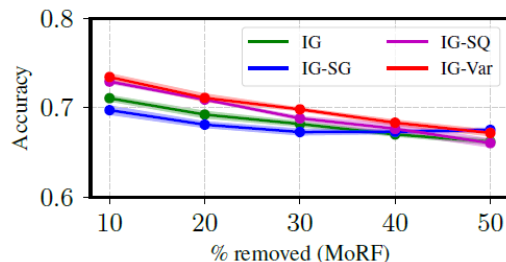
# Experiments

- Consistency under Removal Orders

- CIFAR-10, eight attribution methods.
- Quantitative results: Spearman rank correlation between removal orders:

Retrain		No-Retrain	
MoRF vs. LeRF		MoRF vs. LeRF	
fixed	lin	fixed	lin
-0.01±0.01	<b>0.61±0.01</b>	0.01±0.00	<b>0.58±0.01</b>

- Qualitative results: Our Noisy Linear Imputation in “Retrain” strategy:





# Experiments

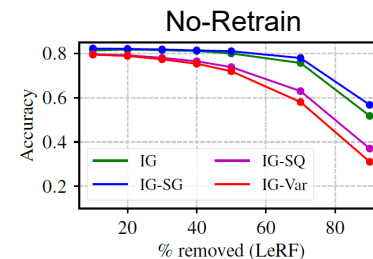
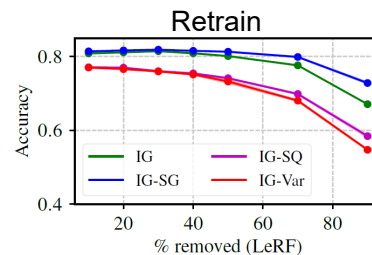
- Efficiency
  - CIFAR-10, eight attribution methods.
  - Consistency evaluation:

**Spearman rank correlation** between evaluation strategies

MoRF		LeRF	
Retain vs. No-Retr.		Retain vs. No-Retr.	
fixed	lin	fixed	lin
$0.15 \pm 0.01$	<b><math>0.84 \pm 0.01</math></b>	$0.09 \pm 0.01$	<b><math>0.94 \pm 0.01</math></b>

- Runtime evaluation:

Strategy	Retrain		No-Retrain	
	fixed <sup>†</sup>	lin	fixed	lin*
Time	$3903 \pm 117$ s	$4686 \pm 2$ s	$18.0 \pm 0.1$ s	$33.3 \pm 0.1$ s
Relative	100 %	120 %	0.5 %	0.9 %







## Thank you!

- More results using GAN imputation and the Food-101 dataset are in the paper!



Fixed value



GAN

- ROAD is available at: [https://github.com/tleemann/road\\_evaluation](https://github.com/tleemann/road_evaluation)  
& Quantus: <https://github.com/understandable-machine-intelligence-lab/Quantus>