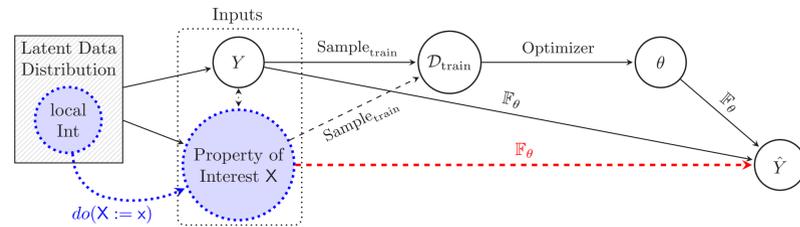


Locally Explaining Prediction Behavior via Gradual Interventions and Measuring Property Gradients

1. Causal Perspective on Explanations

- Data generating systems can be studied as **Structural Causal Models**⁴
- Neural networks are data generating systems and we are interested in **changes in the model outputs \hat{Y}**
- The Causal Hierarchy Theorem⁵ states that in the general case the causal hierarchy does not collapse

Interventional data is necessary for interventional insights



Associational Explanations⁶

Testing conditional dependence
 $H_0 : X \perp\!\!\!\perp \hat{Y} | Y$

Interventional Explanations⁸

Sampling from
 $P(\hat{Y} | do(X := x), Y)$

3. Measuring Changes in Model Behavior

$$\mathbb{E}_x[|\nabla_x \mathbb{F}_\theta(I_x)|] = \int |\nabla_x \mathbb{F}_\theta(I_x)| \cdot p(x) dx$$

$$\stackrel{(*)}{=} \frac{1}{|\mathcal{X}|} \sum_{x \in \mathcal{X}} |\nabla_x \mathbb{F}_\theta(I_x)|$$

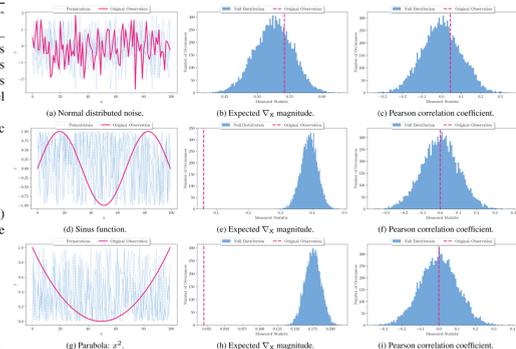
X property of interest
 \mathcal{X} set of equidistant samples
 \mathbb{F}_θ trained model
 I_x input with specific property realization

- Gradient magnitude w.r.t. the property of interest to measure change
- Generalization of the **Causal Concept Effect**⁸ for gradual interventions

Algorithm 1 Hypothesis test for changes in prediction behavior for variations in a property X .

Require: ordered list of predictions $\mathbb{F}_\theta(I_x) \triangleright N$ elements
Require: test statistic $S \triangleright$ for the outputs
Require: integer $K > 0 \triangleright$ Number of Permutations
Require: $\delta \in (0, 1) \triangleright$ Significance Level

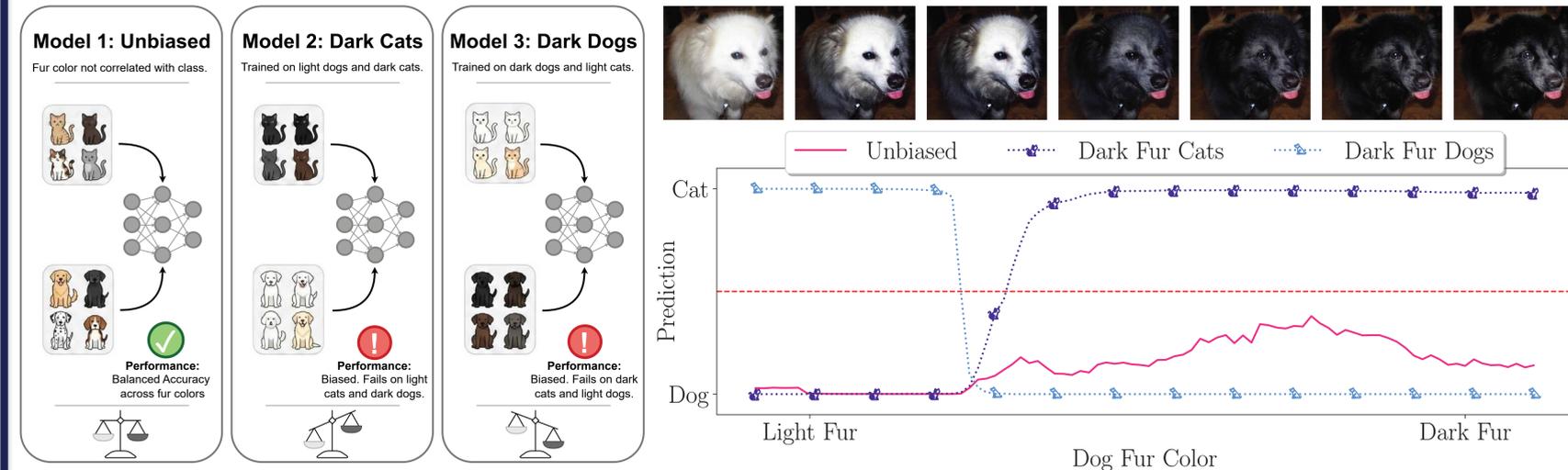
$p \leftarrow 0.0$
 $p_{orig} \leftarrow S(\mathbb{F}_\theta(I_x)) \triangleright$ Estimate the original statistic
for $i \in \{1, \dots, K\}$ **do**
 $\mathbb{F}_\theta(I_x^{perm_i}) \leftarrow \text{permute}(\mathbb{F}_\theta(I_x))$
 $\sigma_{perm_i} \leftarrow S(\mathbb{F}_\theta(I_x^{perm_i}))$
 if compare $(\sigma_{orig}, \sigma_{perm_i})$ **then**
 $p \leftarrow p + 1/K \triangleright$ Increment the p -value
 end if
end for
if $p < \delta$ **then**
 return significant.
else
 return not significant.
end if



What Happens Under Gradual Interventions? Using Generative Models to Trace Causal Drivers of Local Predictions.



Project Page

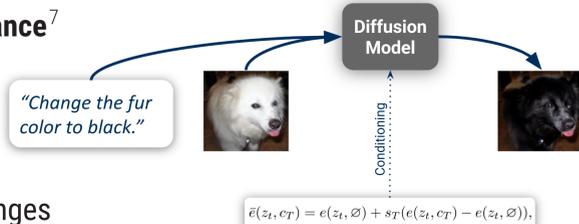


2. Generating Interventional Data

Three Sources



- Classifier-Free Guidance**⁷ for conditioning
- Sampling nonlinear inbetween steps
- Gradual interventions enables study of changes



4. Baseline Comparisons

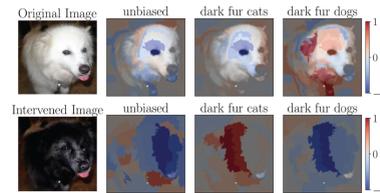
Qualitative¹⁰

Our Approach to Quantify Impact

$\mathbb{E}[|\nabla_x|]$ of the fur color and background property for our three CvD models. Additionally, we report significance ($p < 0.01$) abbreviated as "S" and prediction flips denoted as "F".

Model	Fur Color			Background		
	$\mathbb{E}[\nabla_x]$	S	F	$\mathbb{E}[\nabla_x]$	S	F
Unbiased	.0099	✓	✗	.00060	✓	✗
Dark Cats	.0109	✓	✗	.00006	✓	✗
Dark Dogs	.0110	✓	✓	.00013	✓	✗

LIME as a Local Attribution Baseline



Quantitative⁹

Mean accuracy (↑) and standard deviation in percent (%) of local XAI methods when predicting locally biased model behavior for an intervention. The first column denotes the dataset: Cats vs. Dogs (CvD) and the ISIC archive (ISIC). For CvD, we evaluate the three ConvMixer models trained on separate training datasets. We investigate fur color and synthetic colorful patch interventions, respectively.

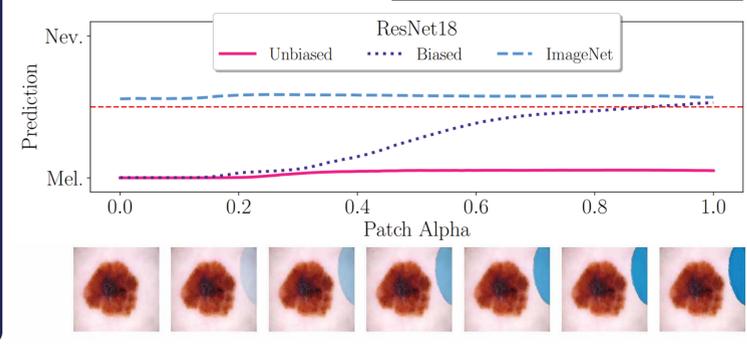
	Model	Ours	G-CAM	Int. Grad.	Occlusion	LIME	K-SHAP	DeepLift
CvD ²	Unbiased	86.13 ± 2.0	84.70 ± 2.1	84.77 ± 2.1	84.73 ± 2.0	84.70 ± 2.2	84.67 ± 2.1	84.83 ± 2.1
	Dark Cats Bias	84.27 ± 1.7	61.43 ± 3.3	64.77 ± 1.8	67.70 ± 1.9	58.63 ± 1.9	60.87 ± 1.0	67.13 ± 2.2
	Dark Dogs Bias	82.20 ± 0.8	64.20 ± 2.7	64.10 ± 2.2	67.10 ± 1.2	56.50 ± 2.1	58.50 ± 2.2	64.77 ± 2.1
ISIC ¹	ResNet18	95.50 ± 1.2	80.00 ± 4.0	78.63 ± 2.7	78.75 ± 4.0	76.12 ± 5.1	75.88 ± 4.6	76.88 ± 6.5
	EfficientNet-B0	94.25 ± 2.4	79.00 ± 4.9	73.75 ± 5.6	75.37 ± 6.3	73.88 ± 5.1	72.75 ± 5.5	76.38 ± 6.5
	ConvNeXt-S	92.88 ± 1.6	78.00 ± 4.7	75.50 ± 3.2	74.88 ± 4.7	75.00 ± 3.2	74.75 ± 2.7	74.88 ± 5.8
	ViT-B/16	95.12 ± 1.6	80.25 ± 3.4	75.87 ± 2.5	77.25 ± 3.5	76.00 ± 3.4	75.88 ± 2.9	77.50 ± 3.9

Melanoma Classification

- Application: Analysis of **melanoma**¹ classifiers
- Synthetic colorful patch intervention, a known bias
- Various backbones

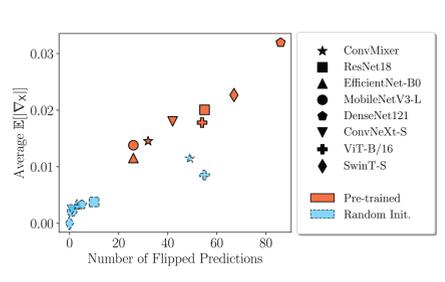
$\mathbb{E}[|\nabla_x|]$ for colorful patch interventions in skin lesion classifiers. We evaluate different models and training data.

Model	Training Data		
	Unbiased	Biased	ImageNet
ResNet18	.00061	.00531	.00062
EfficientNet-B0	.00018	.00495	.00066
ConvNeXt-S	.00001	.00519	.00081
ViT-B/16	.00016	.00208	.00129



Network Training Analysis

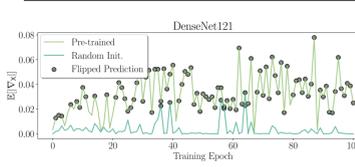
- Training analysis of CelebA³ young vs. old classifiers
- Split between random initialized and pretrained weights
- Intervention on the haircolor



Property gradient magnitudes correlate with prediction flips

Final accuracies in percent (%) achieved by various models trained to differentiate young versus old in CelebA. We split between ImageNet pre-training ("PT") and random initialization ("RI") and calculate the performance delta (Δ).

Model	PT	RI	Perf. Δ
ConvMixer	86.08	81.49	-4.59
ResNet18	85.37	84.38	-0.99
EfficientNet-B0	86.61	84.14	-2.46
MobileNetV3-L	85.94	83.05	-2.88
DenseNet121	85.75	84.20	-1.55
ConvNeXt-S	85.98	83.63	-2.35
ViT-B/16	85.51	72.49	-13.02
Swin-T-S	85.72	50.25	-35.47



Properties are also forgotten in later training stages

References: 1. International skin imaging collaboration, ISIC Archive. https://www.isic-archive.com/. 2. Will Cullen et al., Dogs vs. cats. https://kaggle.com/competitions/dogs-vs-cats. 2013. 3. Zhiwei Liu, Ping Luo, Xiangyong Wang, and Xiaoou Tang. Deep learning face attributes in the wild. ICCV 2015. 4. Judea Pearl. Causality. Cambridge University Press, 2009. 5. Elias Bareinboim, Juan David Correa, Dulger Bejing, and Thomas F. Icard. On pearl's hierarchy and the foundations of causal inference. Probabilistic and Causal Inference, 2022. 6. Christian Reimers, Jakob Rungt, and Joachim Denzler. Determining the relevance of features for deep neural networks. ECDV 2020. 7. Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In NeurIPS 2022 Workshop on Deep Generative Models and Downstream Applications, 2022. 8. Yash Goyal, Amir Feder, Uri Shalit, and Beom Kim. Explaining classifiers with causal concept effect (CACE). 2019. 9. Vitali Petsuk, Amir Das and Kate Saenko. RISE: Randomized Input Sampling for Explanation of Black-box Models. BMVC, 2018. 10. Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "Why should I trust you?" Explaining the predictions of any classifier. ACM SIGKDD, 2016.