





## Gifsplanation via Latent Shift: A Simple Autoencoder Approach to Counterfactual Generation for Chest X-rays

Given:		
E(x)D(z) $f(x)$	<ul> <li>Latent Shift Method:</li> <li>Opposite of an adversarial attack.</li> <li>Perturb the input so the classifier reduces its prediction regularized by the decoder.</li> <li>Compute the gradient of the output of the classifier with respect to the latent space.</li> </ul>	





IoU analysis with expert annotations

IoU is generally low, little variation between methods. Seems inconsistent with how method is qualitatively better.

	Dataset	$\mathrm{Model} \rightarrow$	XRV-all		XRV-mimic_ch	
Task		2D Method	AUC	IoU	AUC	IoU
Mass	NIH	grad guided integrated latentshift-max	0.82	$\begin{array}{c} 0.16{\pm}0.14\\ \textbf{0.19{\pm}0.16}\\ 0.13{\pm}0.13\\ 0.14{\pm}0.17\end{array}$	Mo	odel does t predict
Lung Opacity	RSNA	grad guided integrated latentshift-max	0.84	0.21±0.11 0.21±0.12 0.17±0.10 0.20±0.13	0.75	0.13±0.09 0.09±0.07 0.08±0.07 0.15±0.14
Pneumothorax	SIIM-ACR	grad guided integrated latentshift-max	0.78	0.01±0.02 0.03±0.05 0.01±0.02 0.02±0.04	0.67	0.01±0.02 0.02±0.03 0.01±0.01 0.03±0.07

## <u>Reader study:</u> Two radiologists evaluated how confident they were in a models predictions.

240 Chest X-ray images Radiologists asked: 50% are false positives "How confident are you in the model's prediction? (1-5)"





True Positives:  $0.15\pm0.95$  confidence increase using Latent Shift (p=0.01). False Positives:  $0.04\pm1.06$  increase which is not significant (p=0.57)