

A APPENDIX

A.1 NULL COMPONENTS

Algorithm 1 Learning $\tilde{\mathbf{v}}_\phi$

```

1 #  $f_\theta, f_\phi, f_\psi$ 
2 # epochs: Training epochs # 500
3 #  $\epsilon$ : Learning rate # 0.01
4 # steps: Epochs to reduce the learning rate # [150, 300, 400]
5 # lim: Range to randomly sample initial values for  $\tilde{\mathbf{v}}_\phi$  # [.01, .05, .1, .25, .5, 1, 1.5, 2]
6
7  $\tilde{\mathbf{v}}_\phi \sim U(-lim, lim)$ 
8 for epoch in range(epochs):
9     for  $\mathbf{x}$  in loader: # load a minibatch  $\mathbf{x}$  of images
10         with torch.no_grad():
11              $\mathbf{gt} = f_\psi(f_\phi(f_\theta(\mathbf{x})))$  # original output logits
12
13              $\mathbf{u} = f_\theta(\mathbf{x}) + [0; \tilde{\mathbf{v}}_\phi]$  # introducing noise
14              $\mathbf{out} = f_\psi(f_\phi(\mathbf{u}))$ 
15              $\text{loss} = (\mathbf{out} - \mathbf{gt})^2 \cdot \text{sum}(\text{dim}=-1)^{0.5} \cdot \text{mean}()$  # L2 computation
16
17              $\text{loss.backward}()$  # back-propagate
18              $\mathbf{grad} = \tilde{\mathbf{v}}_\phi \cdot \mathbf{grad}()$ 
19              $\tilde{\mathbf{v}}_\phi = \tilde{\mathbf{v}}_\phi - \epsilon * \mathbf{grad}$  # gradient step
20
21         if epoch in steps:
22              $\epsilon = \epsilon * 0.1$ 
23
24 return  $\tilde{\mathbf{v}}_\phi$ 

```

A.2 WATERMARKING IMAGES

As mentioned in the main paper, we use SQLP for minimising equation 7. We use the implementation provided by SciPY (Virtanen et al., 2020) and run it 5000 iterations. Experimentally, we observed that using $p = 1$ norm provides better watermarking behaviour than $p = 2$. All reported and displayed nullspace noise content is with $p = 1$ unless stated otherwise.

To quantitatively assess the robustness of model to the watermarking process, we will have to watermark every image in the dataset. This process requires considerable time and compute to execute. Instead, we perform the evaluation for randomly selected 20 images and compute the % match predictions and absolute difference in the predicted probabilities. We found that 85% of the watermarked images were classified as the source image category. For the mean absolute difference, we compute it between the predicted probabilities for source image category both for the source image and the watermarked image. We observed that the difference in confidence varied 11.63% on an average.

In figure 6 we show the original, $|\mathbf{v}_\theta|$ and the resulting watermark images. The watermark image is the same as one reported in the main paper. Using nullspace watermarking, we notice that shape details are more likely to be transferred than other information from the watermark image.

A.3 TARGETED NULLSPACE NOISE

Instead of directly minimising the unconstrained equation 7 with huge number of variables ($r \times r \times c$), we manually fix the values of 2 channels (green and blue) and only perform the minimisation to learn the red channel values for the transformed image. This reduces the number of parameters to one-third and also retains a lot of target image information without any loss. This is also the reason why we observe different colored tint for the transformed images. With respect to the implementation, the details are identical to that for watermarking.

A.4 FOOLING XAI METHODS

In figure 7 we show the saliency maps generated by various XAI methods. Even though the maps generated by methods other than LRP are poor (hard to interpret), we see that the source and transformed respond similarly to these methods.



Figure 6: Watermark superposition using the nullspace basis vectors.

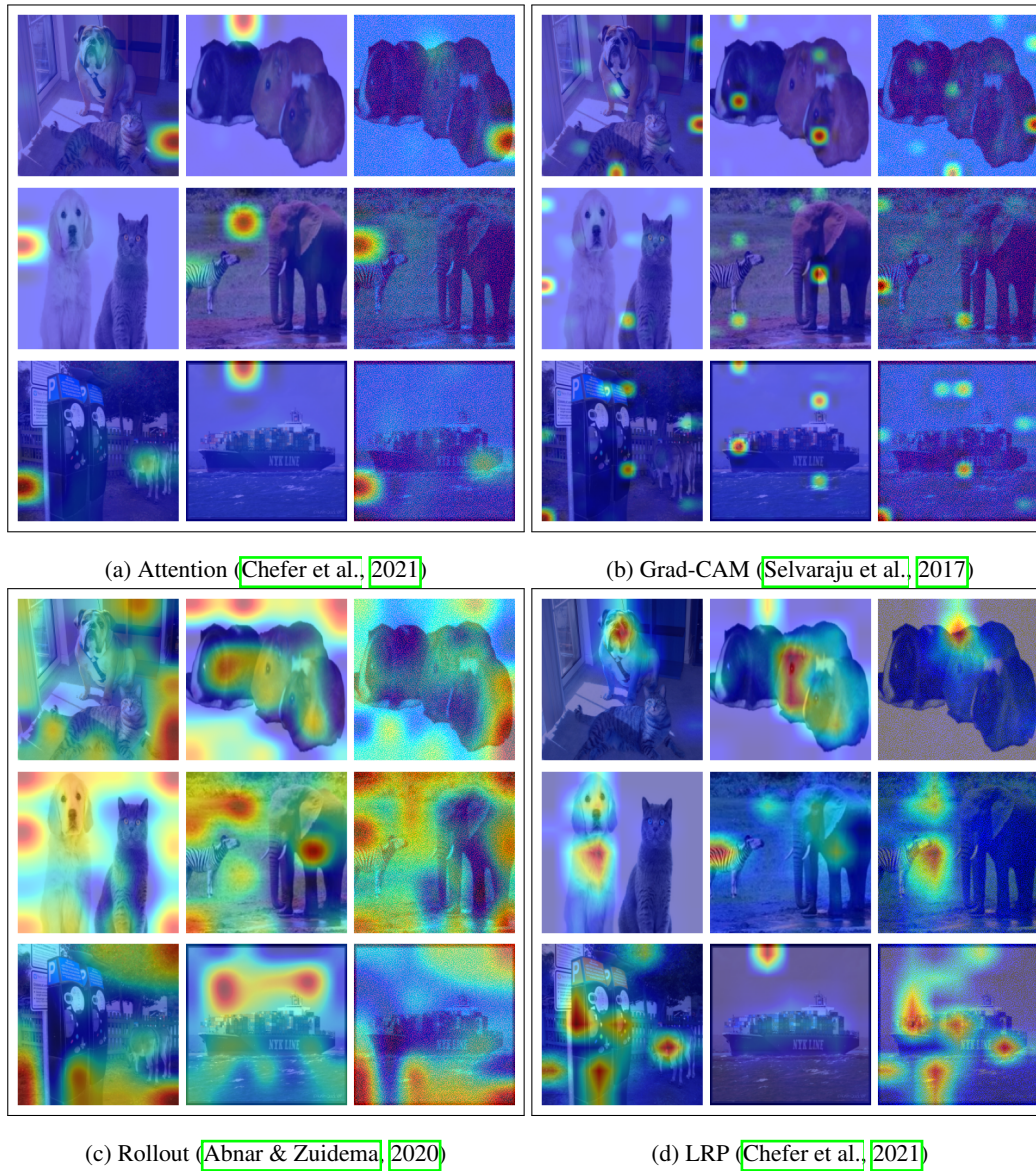


Figure 7: Interpretability maps generated via different methods for (source, target, transformed) images