

Input

Feature Extractor

Feature Maps

Prediction

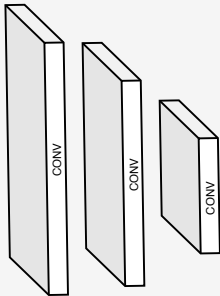
Channel-Wise Gradient
Average

Masks

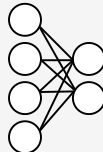
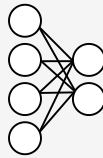
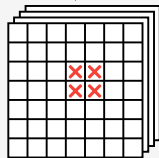
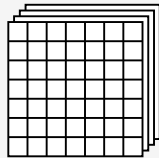


...

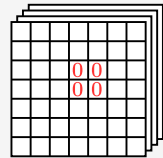
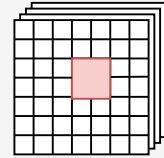
e.g. PACS



e.g. ResNet



$$\frac{\partial y_c}{\partial \mathbf{z}}$$



$$y_c - \tilde{y}_c$$

Apply mask only for
Top-b percentile where it
decreases confidence the
most (e.g. only elephant)

Backpropagate Loss

$$\nabla_{\theta} \mathcal{L}(\cdot, \cdot)$$

