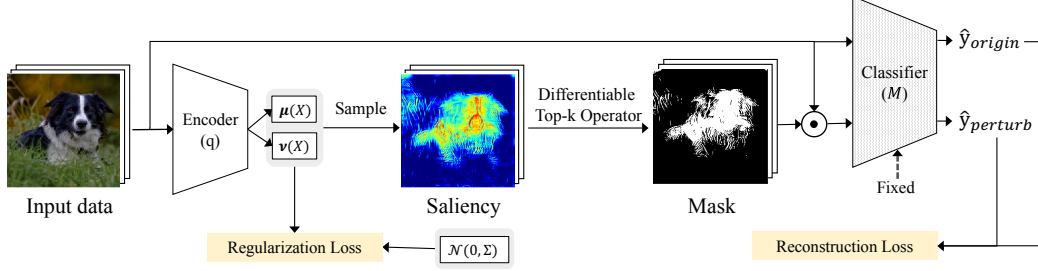APPENDIX

## A    SCHEMATIC DESCRIPTION



Figure 1: Schematic description.

An input image is fed into the encoder network that gives the mean and the variance of Guassian distribution (which is the approximate posterior). A saliency map is sampled from it, followed by passing a differentiable top-$k$ operator to provide a binary mask. With this mask, the input image is perturbed. The perturbed image is passes a classifier $M$ that gives categorical probability $\hat{\boldsymbol{y}}_{perturb}$. The loss function is composed of two terms: the reconstruction term between $\hat{\boldsymbol{y}}_{perturb}$ and the categorical probability obtained from the original image $\hat{\boldsymbol{y}}_{origin}$, and the regularization term between the approximate posterior $q(\boldsymbol{s}|\boldsymbol{x}) = \mathcal{N}(\boldsymbol{s}; \boldsymbol{\mu}(\boldsymbol{x}), \mathrm{diag}(\boldsymbol{\nu}(\boldsymbol{x})))$ and the prior distribution $\mathcal{N}(\boldsymbol{s}; \boldsymbol{0}, \boldsymbol{\Sigma})$. As the classifier $M$ is the model that we aim to interpret, it is fixed so as the parameter not to be updated while training the encoder network.
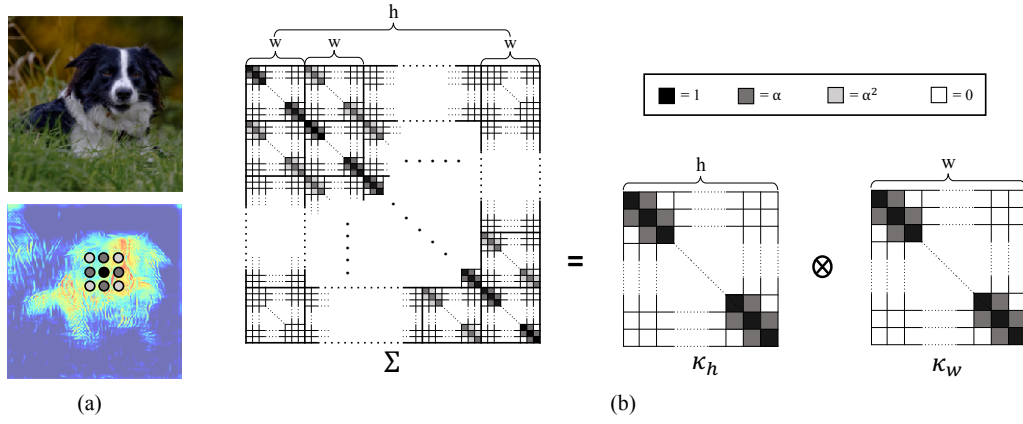
# B SOFT-TV GAUSSIAN PRIOR



Figure 2: Soft-TV Gaussian prior.

We consider the prior knowledge that adjacent pixels in a saliency map have positive correlation (Figure 2(a)). After modeling the prior as Gaussian distribution $\mathcal{N}(s; \mathbf{0}, \mathbf{\Sigma})$, the covariance matrix $\mathbf{\Sigma}$ is designed to infuse this prior knowledge into the prior distribution (Figure 2(b)). The covariance matrix is then decomposed by Kronecker product to better calculate the KL divergence of the regularization loss.

## C  IMPLEMENTATION DETAILS

**Encoder architecture**  We use 17 convolution layers for the encoder network. To make the spatial size of the encoder's input and output to be same, we do not use a pooling layer. Every convolution layer is comprised of a convolution with kernel size $3 \times 3$, stride 1, and padding 1, followed by batch normalization and a rectified linear unit. The number of output channels for each convolution is as follows: $[64, 64, 64, 64, 32, 32, 32, 32, 32, 16, 16, 16, 16, 16, 16, 2]$. The encoder network provides $\boldsymbol{\mu} \in \mathbb{R}^{h \times w}$ and $\eta \in \mathbb{R}^{h \times w}$ for each channel of the output where $\boldsymbol{\mu}$ is the mean of the Gaussian distribution and $\eta = \log \nu$ with $\nu \in \mathbb{R}^{h \times w}$ the variance of the Gaussian distribution.

**hyper-parameters**  We use Adam (Kingma & Ba, 2014) optimizer with learning rate to be $0.0001$, weight decay to be $0.0005$, and betas to be $(0.9, 0.99)$. We use batch size of $128$ while training the encoder network. We run 10 epochs for the Imagenet dataset, and save the network that has the lowest loss.

# D PROOF OF REGULARIZATION LOSS EQUATION

Let us first define the notation. $\otimes$ is the Kronecker product, and $\odot$ is the element-wise multiplication. $\text{sum}(\cdot)$ is the summation of elements. For a vector $\boldsymbol{b} \in \mathbb{R}^{hw}$, we denote $\text{rsh}(\boldsymbol{b}) \in \mathbb{R}^{w \times h}$ as reshaping the vector $\boldsymbol{b}$ to the matrix where the $(i,j)$-th entry is $\boldsymbol{b}[i + wj]$, and $\text{diag}(\boldsymbol{b})$ as the diagonal matrix where the diagonal is $\boldsymbol{b}$. Also, for a square matrix $\boldsymbol{B}$, $\text{diag}(\boldsymbol{B})$ is the vector where the $i$-th entry is $B[i,i]$. For a matrix $\boldsymbol{B}$, $\text{vec}(\boldsymbol{B})$ denotes the vectorization by stacking the columns of the matrix $\boldsymbol{B}$ to a single column vector.

Recall that the approximate posterior is $q(\boldsymbol{s}|\boldsymbol{x}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}_0)$ where $\boldsymbol{\Sigma}_0 = \text{diag}(\boldsymbol{\nu})$ with $\boldsymbol{\mu}, \boldsymbol{\nu} \in \mathbb{R}^{hw}$, and the prior distribution is $p(\boldsymbol{s}|\boldsymbol{x}) = \mathcal{N}(\boldsymbol{0}, \boldsymbol{\Sigma}_1)$ where $\boldsymbol{\Sigma}_1 = \boldsymbol{\kappa}_h \otimes \boldsymbol{\kappa}_w$ with $\boldsymbol{\kappa}_h \in \mathbb{R}^{h \times h}$ and $\boldsymbol{\kappa}_w \in \mathbb{R}^{w \times w}$. KL divergence between two Gaussian distribution is:

$$
\begin{aligned}
D_{\text{KL}}[\,q(\boldsymbol{s}|\boldsymbol{x}) \,\|\, p(\boldsymbol{s}|\boldsymbol{x})\,] &= D_{\text{KL}}[\,\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}_0) \,\|\, \mathcal{N}(\boldsymbol{0}, \boldsymbol{\Sigma}_1)\,] \\
&= \frac{1}{2}\left( \text{tr}\left(\boldsymbol{\Sigma}_1^{-1}\boldsymbol{\Sigma}_0\right) + \boldsymbol{\mu}^T\boldsymbol{\Sigma}_1^{-1}\boldsymbol{\mu} - hw + \log\frac{|\boldsymbol{\Sigma}_1|}{|\boldsymbol{\Sigma}_0|}\right) .
\end{aligned}
\tag{1}
$$

We compute each term in RHS of equation 1 to make it computationally efficient.

$$
\begin{aligned}
\text{tr}\left(\boldsymbol{\Sigma}_1^{-1}\boldsymbol{\Sigma}_0\right) &= \text{tr}\left((\boldsymbol{\kappa}_h \otimes \boldsymbol{\kappa}_w)^{-1} \cdot \text{diag}(\boldsymbol{\nu})\right) \\
&= \text{tr}\left((\boldsymbol{\kappa}_h^{-1} \otimes \boldsymbol{\kappa}_w^{-1}) \cdot \text{diag}(\boldsymbol{\nu})\right) \\
&= \text{diag}\left(\boldsymbol{\kappa}_h^{-1} \otimes \boldsymbol{\kappa}_w^{-1}\right) \odot \text{diag}(\boldsymbol{\nu}) \\
&= \text{diag}\left(\boldsymbol{\kappa}_w^{-1}\right)^T \cdot \text{rsh}(\boldsymbol{\nu}) \cdot \text{diag}\left(\boldsymbol{\kappa}_h^{-1}\right) .
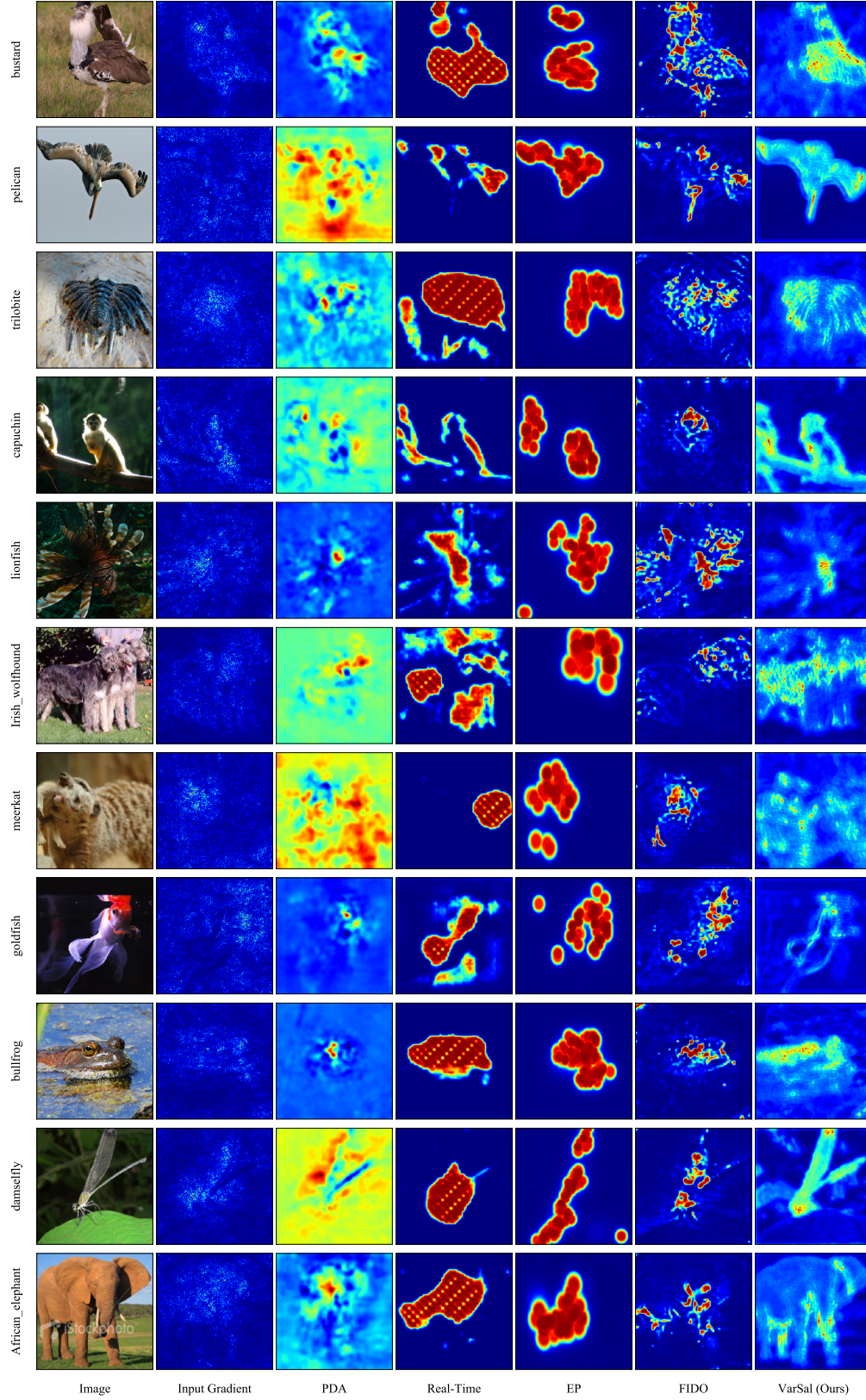\end{aligned}
\tag{2}
$$

$$
\begin{aligned}
\boldsymbol{\mu}^T\boldsymbol{\Sigma}_1^{-1}\boldsymbol{\mu} &= \boldsymbol{\mu}^T \cdot \left(\boldsymbol{\kappa}_h^{-1} \otimes \boldsymbol{\kappa}_w^{-1}\right) \cdot \boldsymbol{\mu} \\
&= \boldsymbol{\mu}^T \cdot \left(\boldsymbol{\kappa}_h^{-1} \otimes \boldsymbol{\kappa}_w^{-1}\right) \cdot \text{vec}\left(\text{rsh}(\boldsymbol{\mu})\right) \\
&= \boldsymbol{\mu}^T \cdot \text{vec}\left(\boldsymbol{\kappa}_h^{-1} \cdot \text{rsh}(\boldsymbol{\mu}) \cdot \boldsymbol{\kappa}_w^{-1T}\right) \\
&= \text{sum}\left(\text{rsh}(\boldsymbol{\mu}) \odot \left(\boldsymbol{\kappa}_w^{-1} \cdot \text{rsh}(\boldsymbol{\mu}) \cdot \left(\boldsymbol{\kappa}_h^{-1}\right)^T\right)\right) .
\end{aligned}
\tag{3}
$$

$$
\begin{aligned}
\log\frac{|\boldsymbol{\Sigma}_1|}{|\boldsymbol{\Sigma}_0|} &= \log|\boldsymbol{\kappa}_h \otimes \boldsymbol{\kappa}_w| - \log|\text{diag}(\boldsymbol{\nu})| \\
&= \log|\boldsymbol{\kappa}_h|^w|\boldsymbol{\kappa}_w|^h - \log\prod_{i=1}^{hw}\boldsymbol{\nu}_i \\
&= w \cdot \log|\boldsymbol{\kappa}_h| + h \cdot \log|\boldsymbol{\kappa}_w| - \log\prod_{i=1}^{hw}\boldsymbol{\nu}_i \\
&= -\text{sum}(\log\boldsymbol{\nu}_i) + \text{const} .
\end{aligned}
\tag{4}
$$

Therefore, the equation 1 is derived as:

$$
\begin{aligned}
D_{\text{KL}}[\,q(\boldsymbol{s}|\boldsymbol{x}) \,\|\, p(\boldsymbol{s}|\boldsymbol{x})\,] &= \text{diag}\left(\boldsymbol{\kappa}_w^{-1}\right)^T \cdot \text{rsh}(\boldsymbol{\nu}) \cdot \text{diag}\left(\boldsymbol{\kappa}_h^{-1}\right) \\
&\quad + \text{sum}\left(\text{rsh}(\boldsymbol{\mu}) \odot \left(\boldsymbol{\kappa}_w^{-1} \cdot \text{rsh}(\boldsymbol{\mu}) \cdot \left(\boldsymbol{\kappa}_h^{-1}\right)^T\right)\right) \\
&\quad - \text{sum}(\log\boldsymbol{\nu}_i) + \text{const} .
\end{aligned}
\tag{5}
$$

# E  QUALITATIVE RESULTS



Image — Input Gradient — PDA — Real-Time — EP — FIDO — VarSal (Ours)

REFERENCES

Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.