

Reducing ReLU Count for Privacy-Preserving CNN Speedup

Inbar Helbitz

Shai Avidan

Tel-Aviv University

ReLU Sharing:

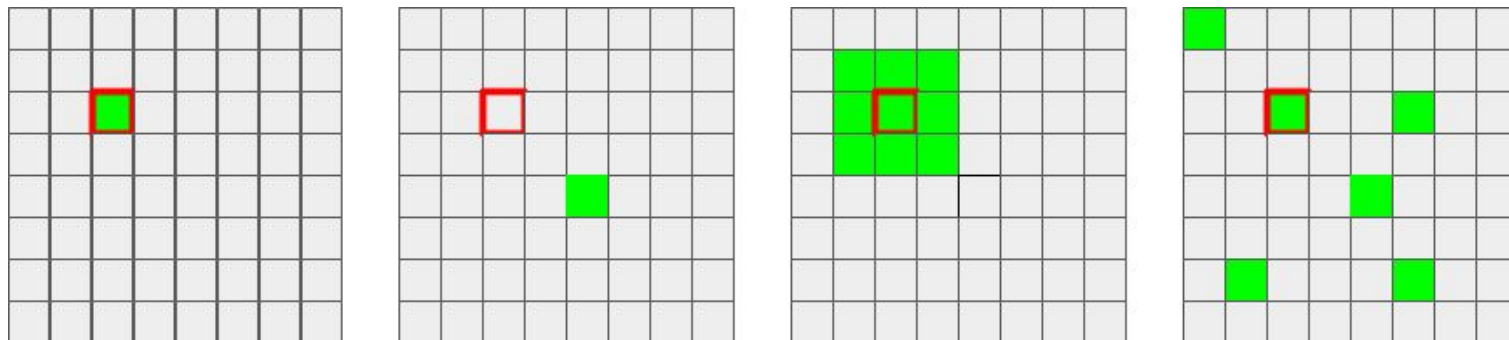
$$\text{gReLU}(s(\mathbf{p}), \mathbf{v}_i^T \mathbf{s}) = \begin{cases} s(\mathbf{p}) & \mathbf{v}_i^T \mathbf{s} \geq 0 \wedge \mathbf{p} \in i \\ 0 & \text{otherwise} \end{cases}$$

\mathbf{s} denote the convolution result

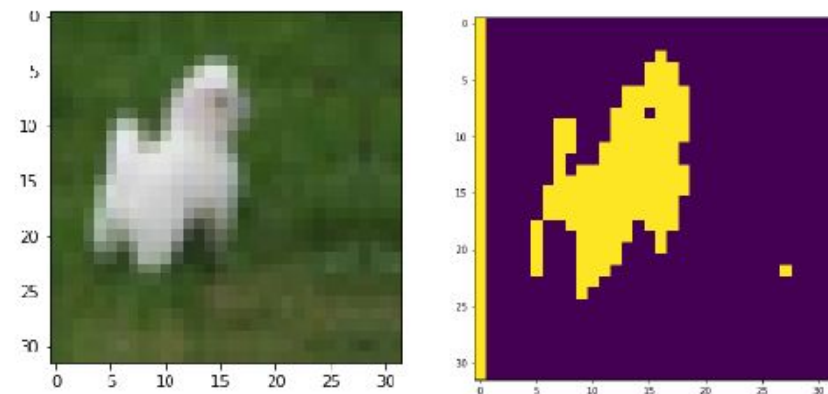
\mathbf{p} is the (2D) activation location

\mathbf{v} is a weight vector, can be set during training or can be learned

$\mathbf{p} \in i$ denotes activation \mathbf{p} belongs to patch i



ReLU of neighbor pixels are highly correlated



We can cut the number of ReLU operations by up to three orders of magnitude and, as a result, cut the communication bandwidth by 60% and the runtime by 55%.