

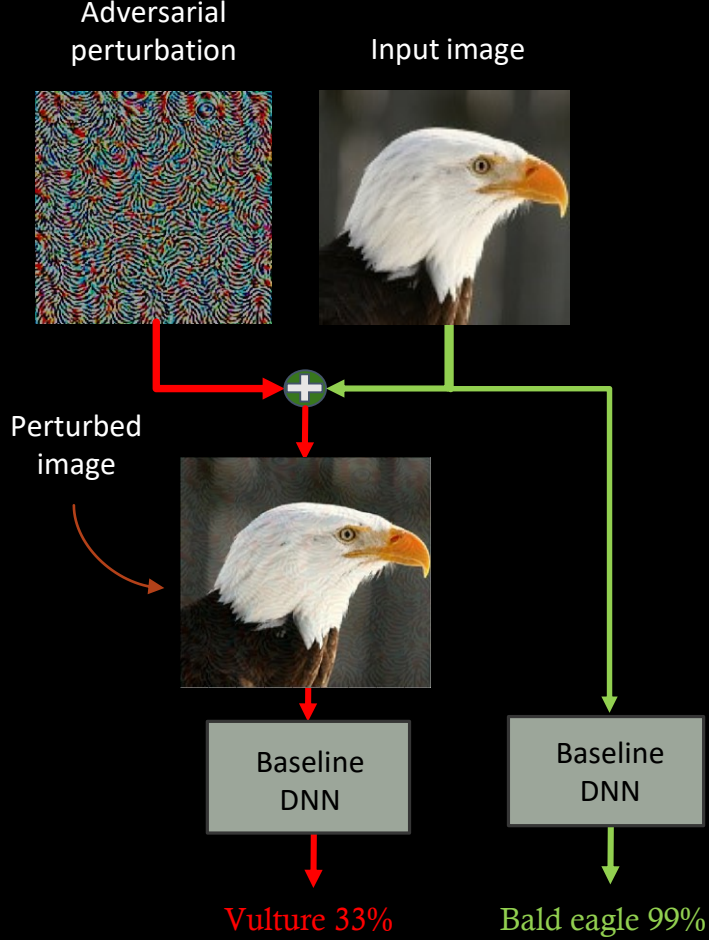
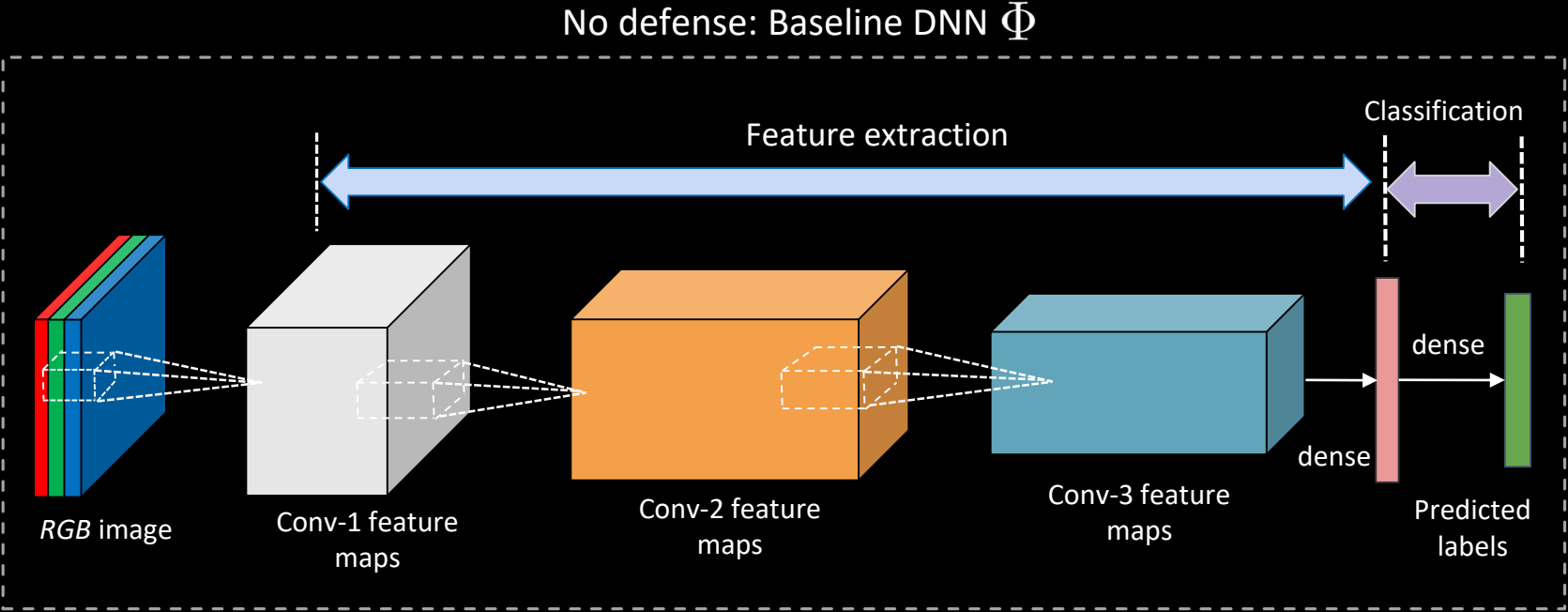
# Defending Against Universal Attacks Through Selective Feature Regeneration

Tejas Borkar<sup>1</sup>, Felix Heide<sup>2,3</sup>, Lina Karam<sup>1,4</sup>

<sup>1</sup>Arizona State University <sup>2</sup>Princeton University <sup>3</sup>Algolux <sup>4</sup>Lebanese American University

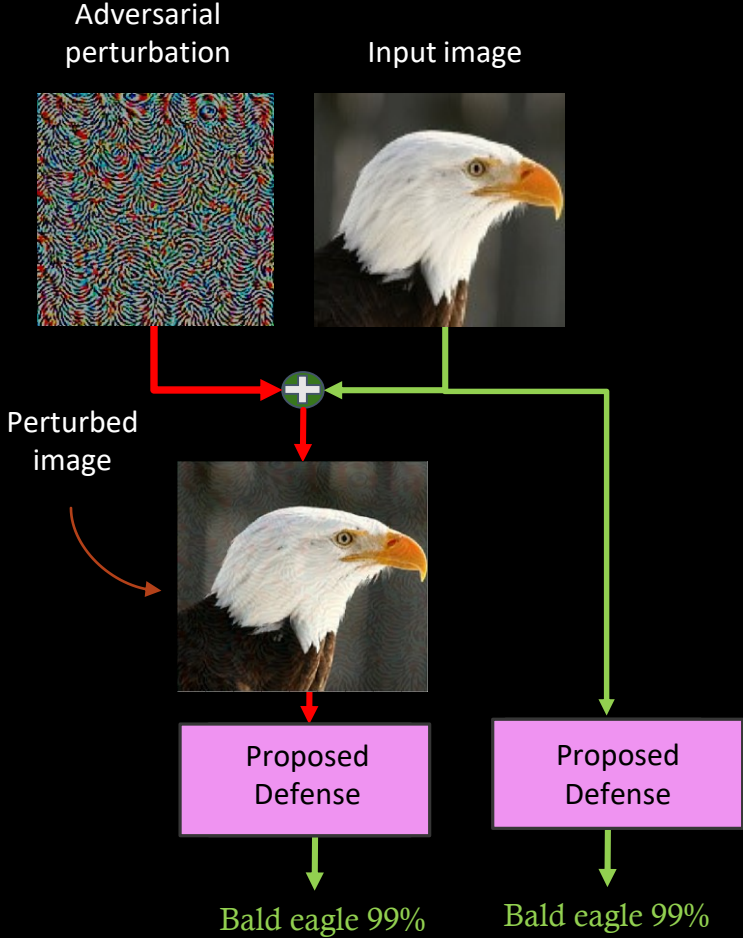
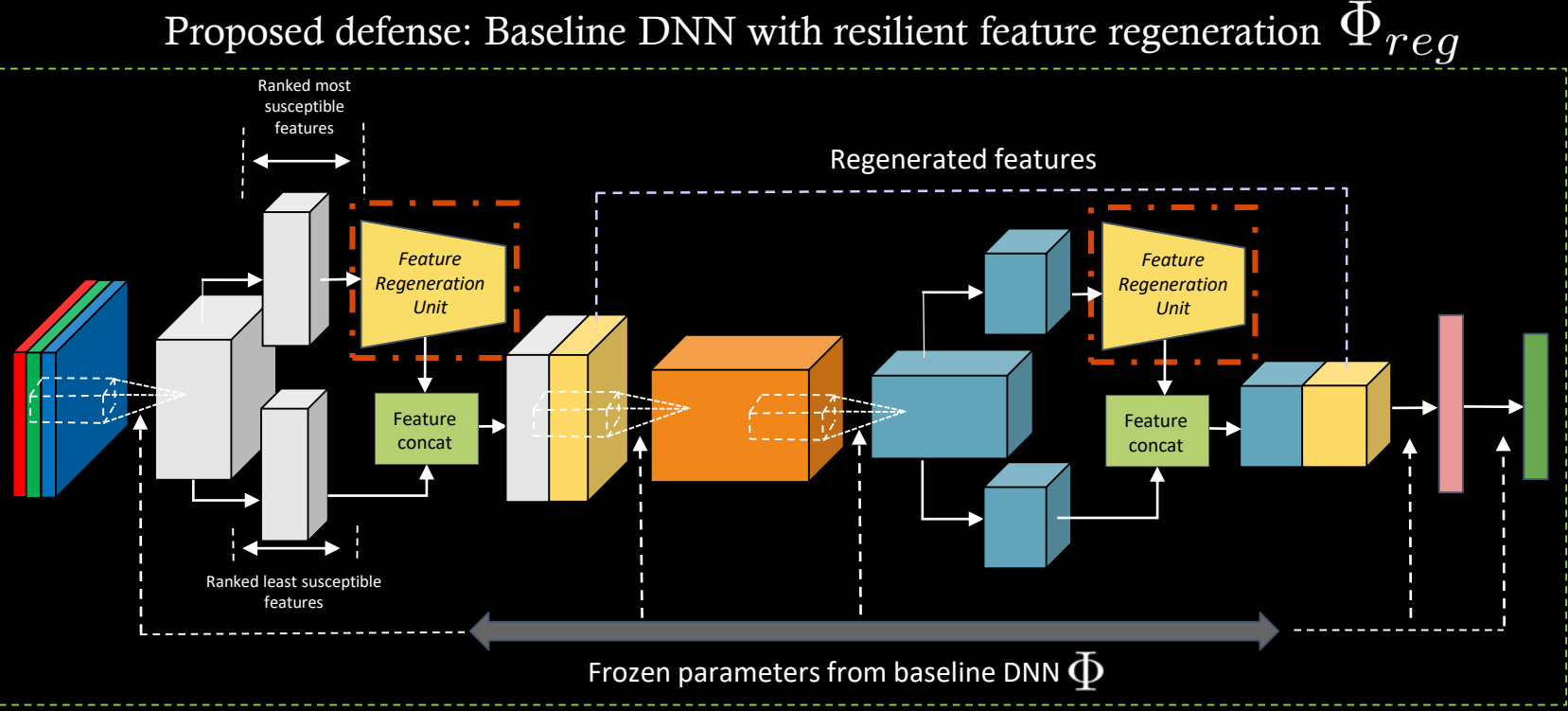
# Universal Adversarial Attacks

- Image agnostic and transferable across networks



# Defending against Universal Adversarial Attacks

- Selective feature regeneration effectively restores robustness

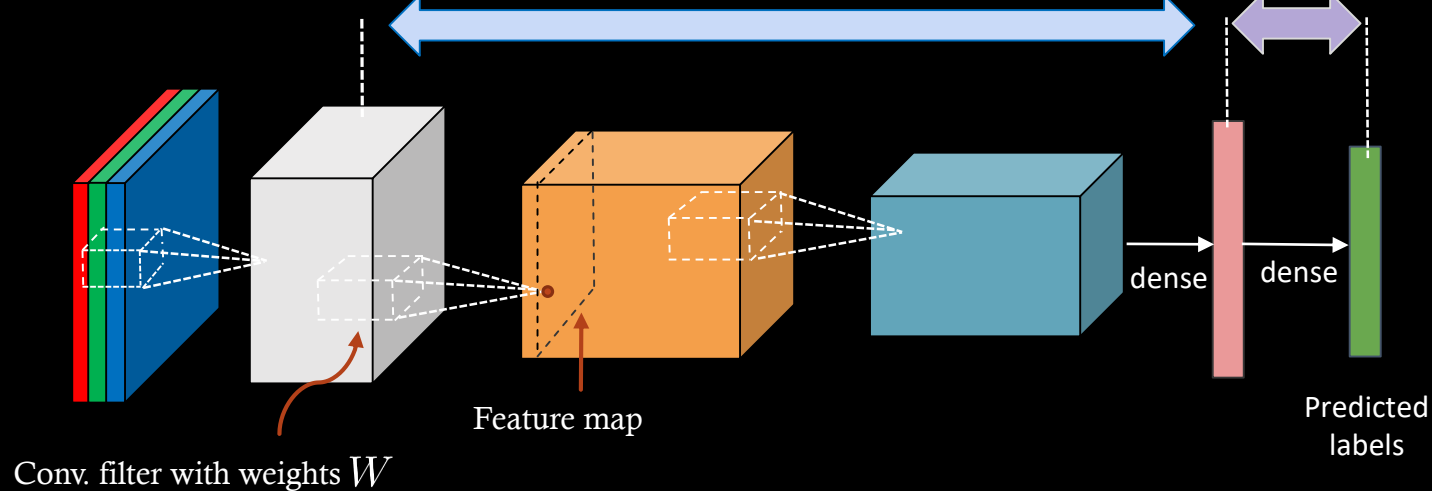


# Ranking CNN Filters Based on Noise Susceptibility

Sample Baseline DNN  $\Phi$

Feature extraction

Classification



We show:

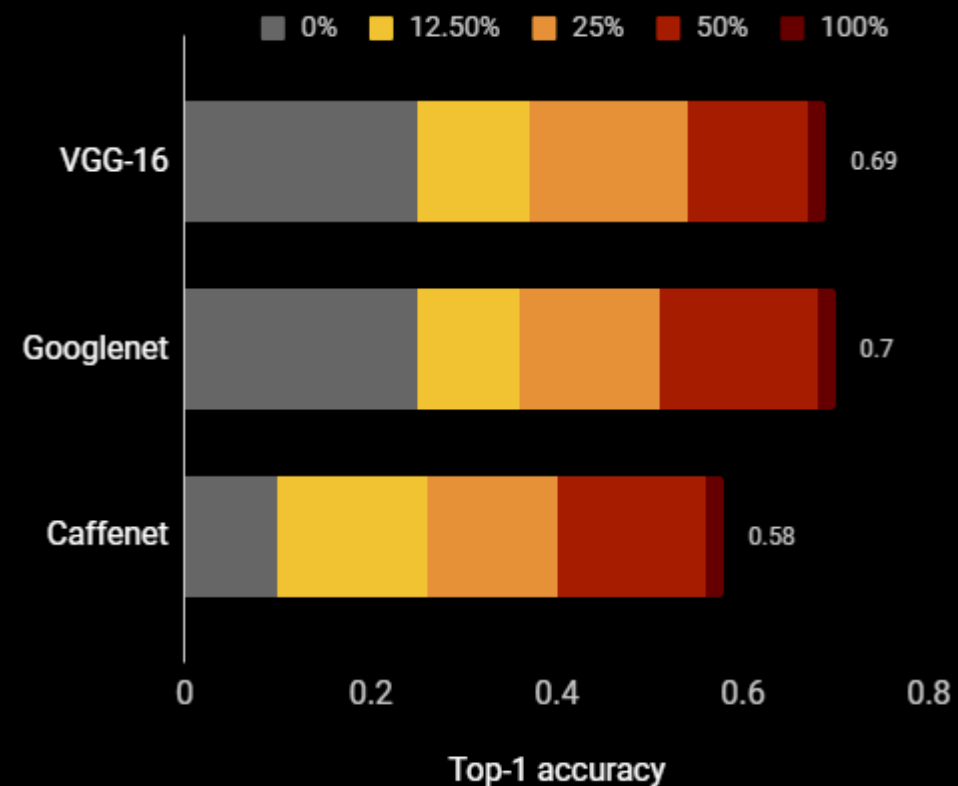
- Max perturbation level induced in feature map

$\propto$

$l_1$ -norm of the filter weight ( $\|W\|_1$ )

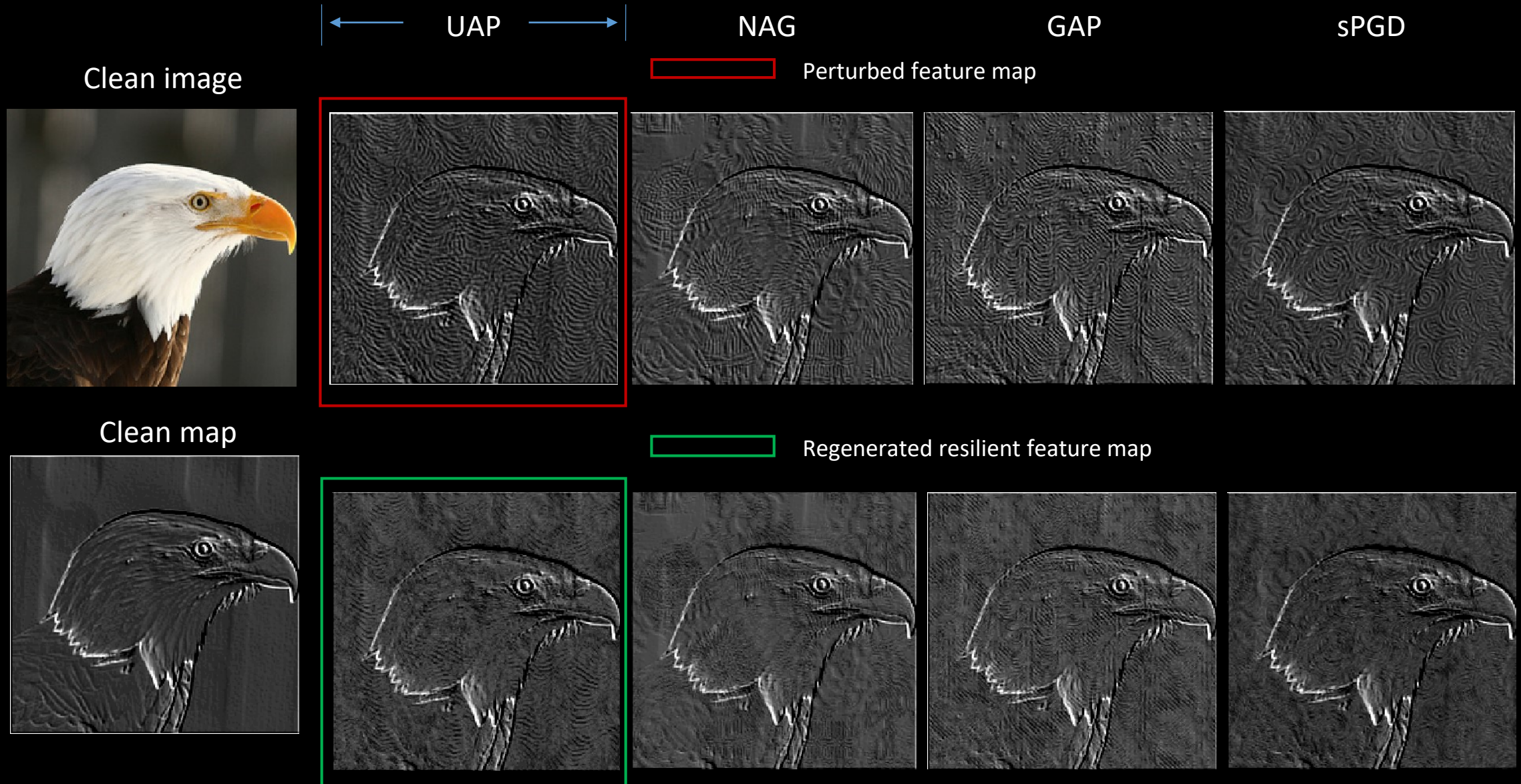
## Suppressing perturbations in ranked filters' output maps

Percentage of suppressed maps in conv-1



# Robustness to Unseen Universal Adversarial Attacks

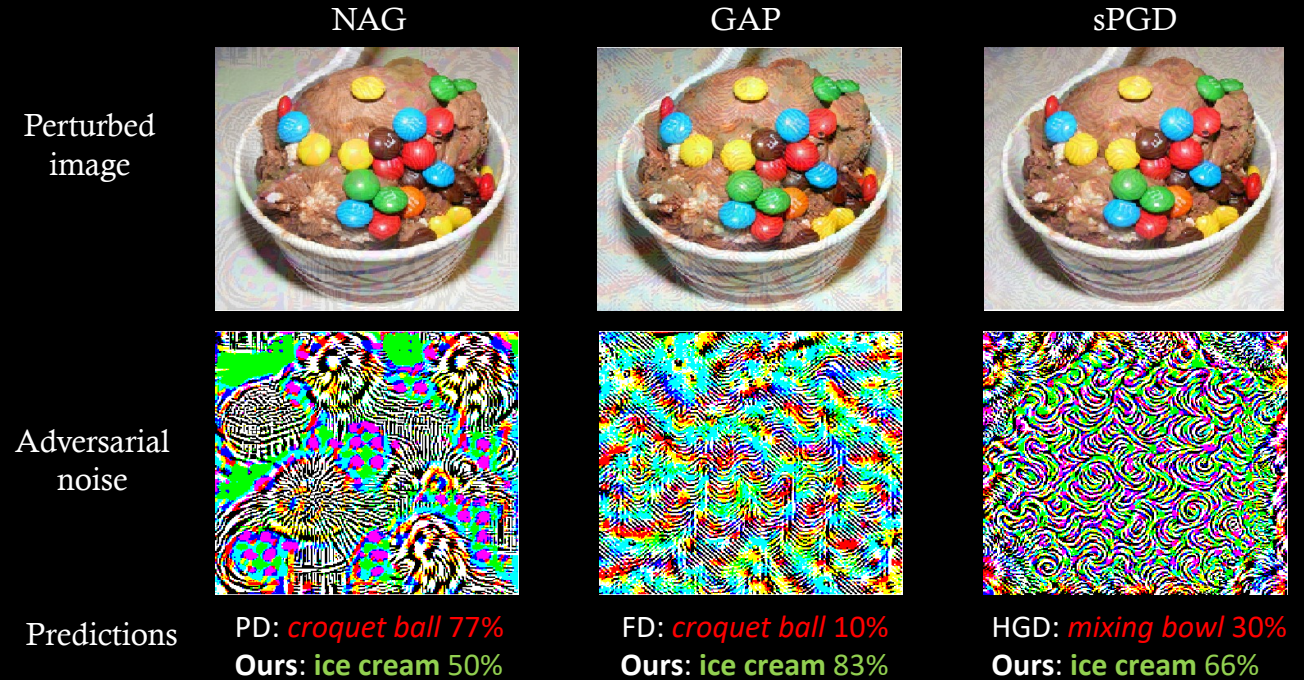
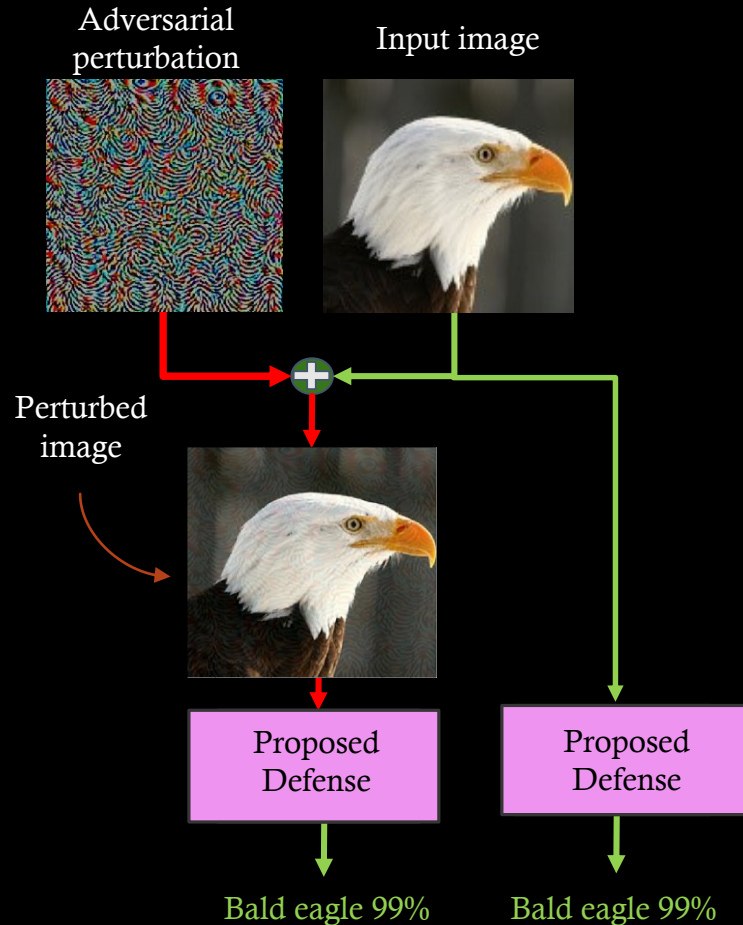
- Defense trained on only UAP noise samples



# Defending Against Universal Attacks Through Selective Feature Regeneration

Robustness to image-agnostic noise:

Robustness to unseen universal attacks:



Summary:

- Novel  $\ell_1$ -norm measure identifies and ranks adversarially susceptible feature maps
- Selective regeneration of only the most vulnerable feature maps restores robustness

Code: <https://github.com/tsborkar/Selective-feature-regeneration>