# Bounded Adversarial Attack on Deep Content Features
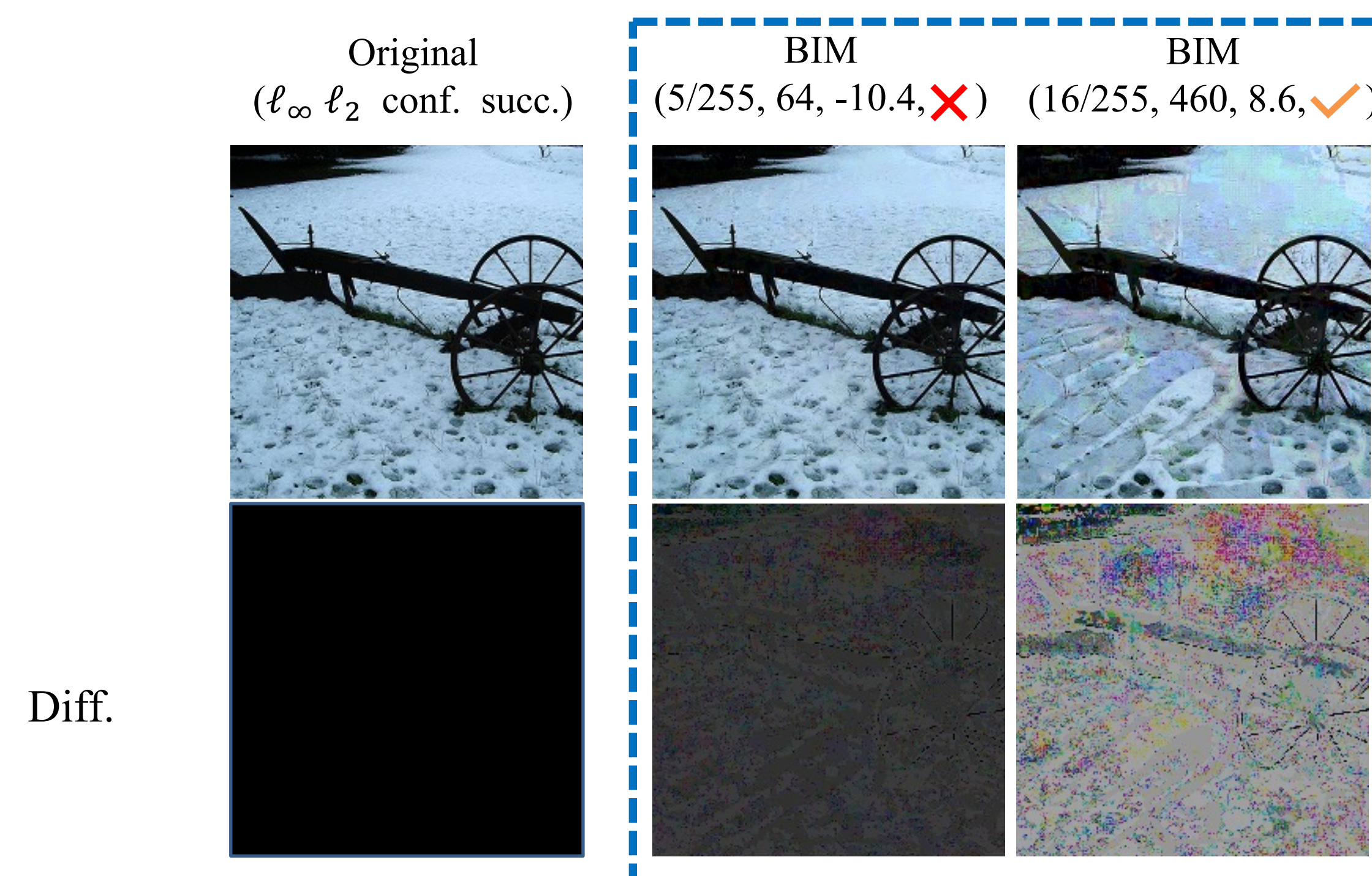
Qiuling Xu, Guanhong Tao and Xiangyu Zhang @ Purdue University

CVPR JUNE 19-24 2022 NEW ORLEANS • LOUISIANA

## Existing Adversarial Attack Bound is Hard to Scale

➤ Pixel Space

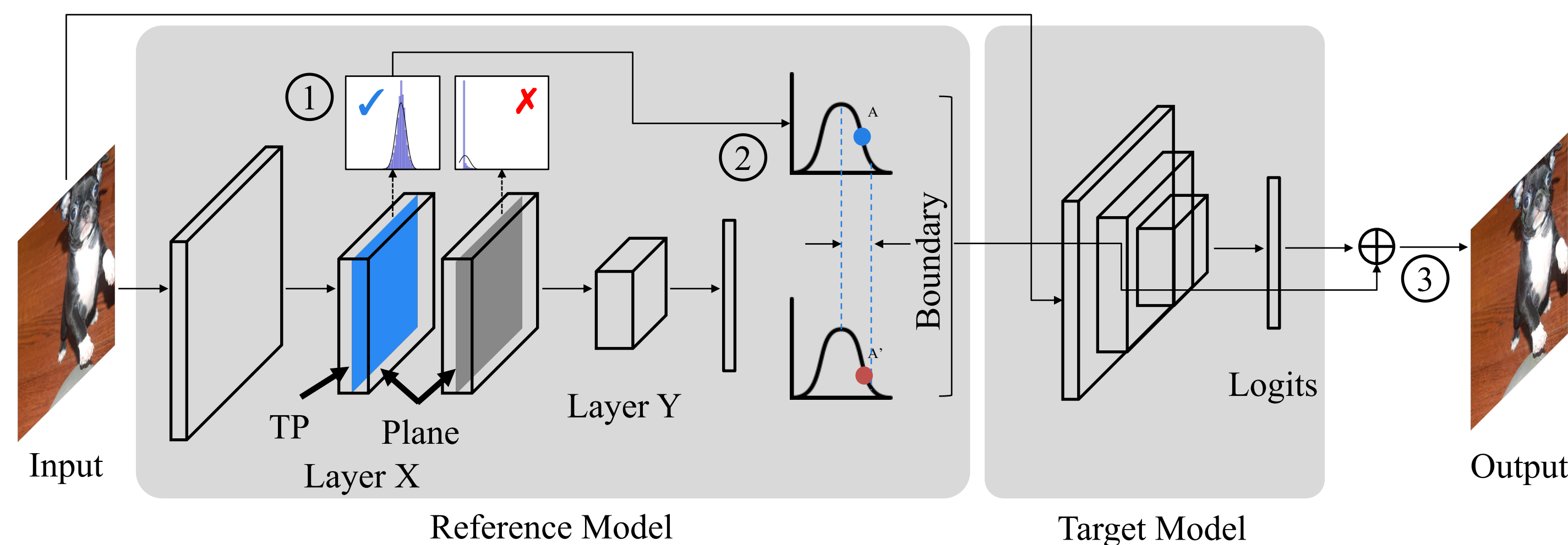Attack fails given a small bound; Attack is detectable on a larger bound.



Original
($\ell_\infty$ $\ell_2$ conf. succ.)

BIM
(5/255, 64, -10.4, ✗)

BIM
(16/255, 460, 8.6, ✓)

Diff.

➤ Feature Space

Attack is detectable on samples with high confidence.



Ours (D²B)
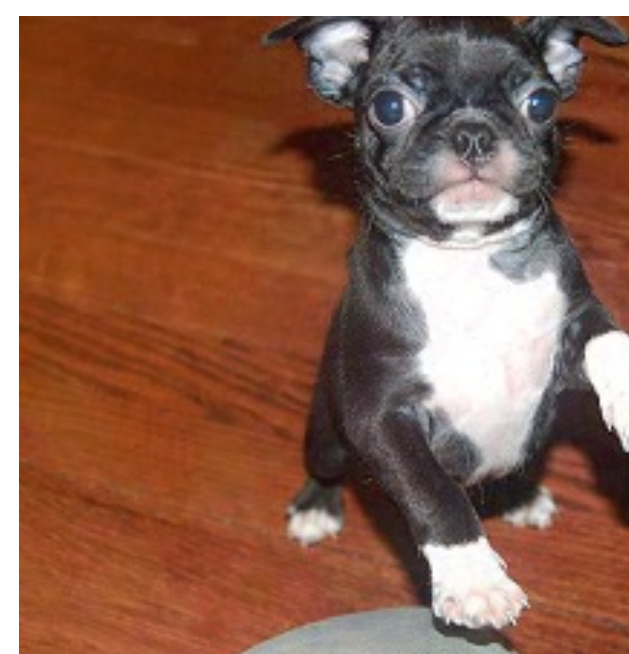(0.22, 586, **34.8**, ✓)
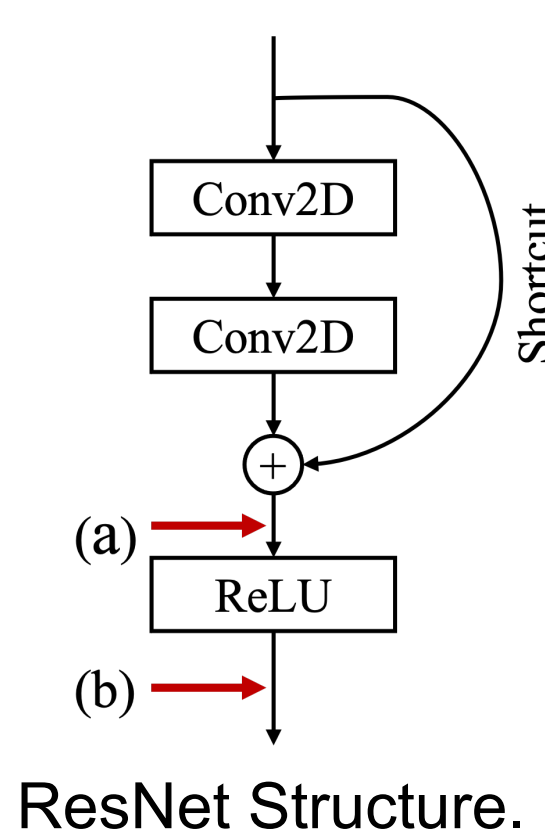
FS
(1.0, 3.8k, 15.4, ✓)

SM
(0.95, 15k, 31.5, ✓)

Diff.

## Method - Quantile Bound on Gaussian Representation



Input | TP | Plane | Layer Y | Boundary | Logits | Output
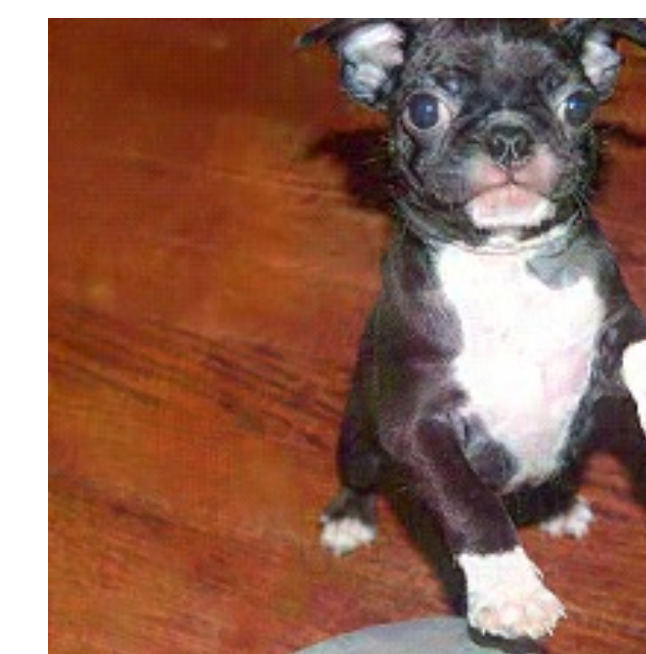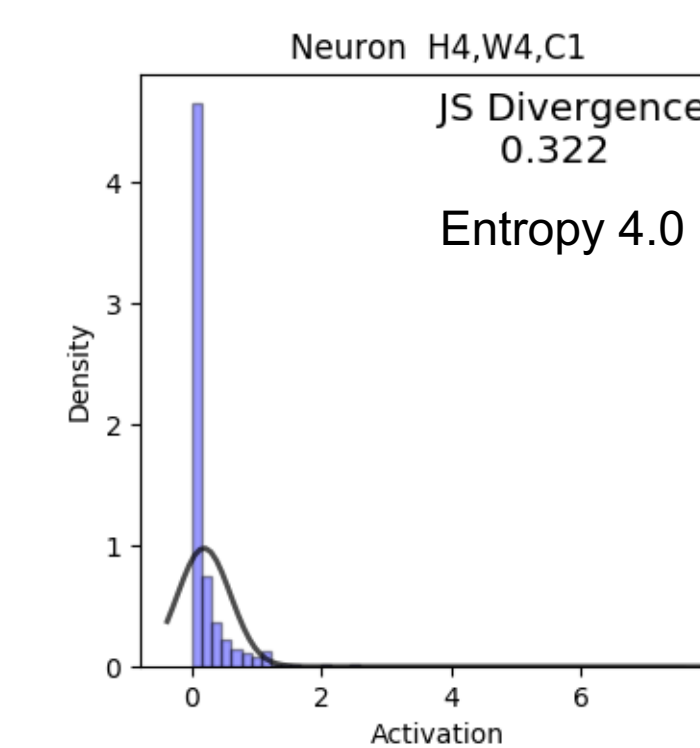Layer X
Reference Model | Target Model

➤ ① Throttle plane (TP) selection
➤ ② Internal distribution boundary constraint
➤ ③ Adversarial sample generation with combined losses
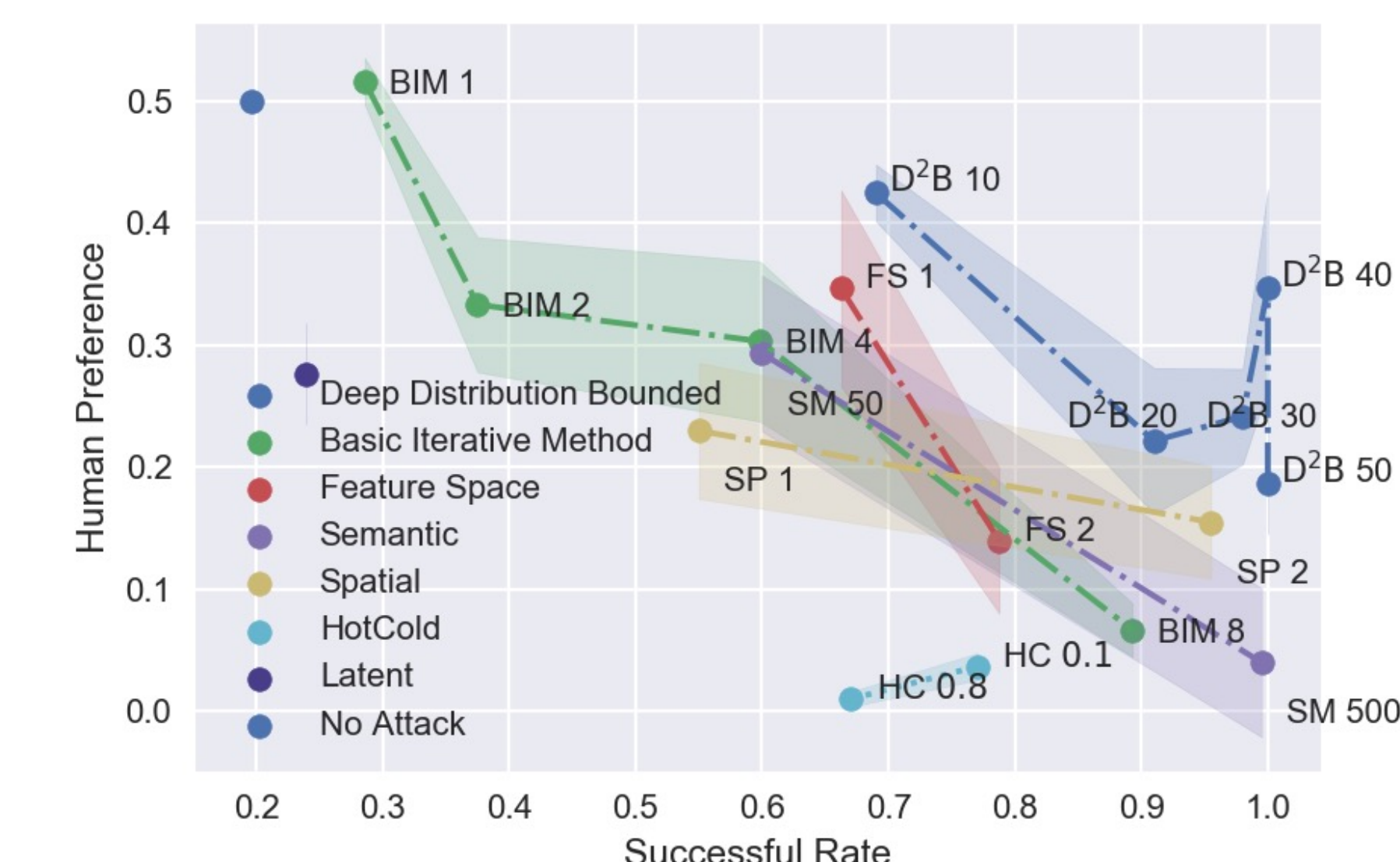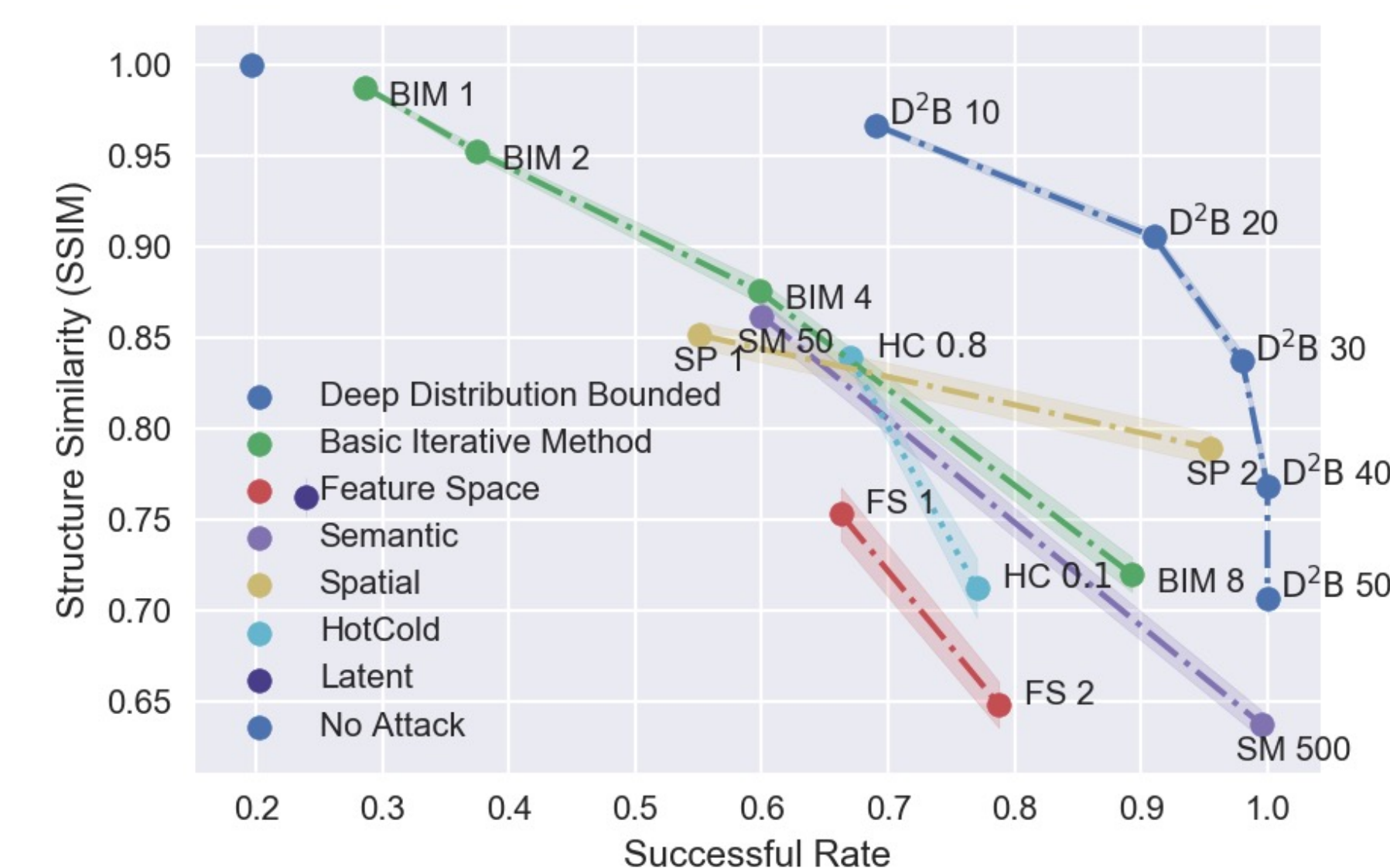
## Distribution of the Representation Matters



ResNet Structure.

<span style="color:red">Which representation should we use for quantile bound? Look for Gaussian.</span>

Neuron H4,W4,C1
JS Divergence 0.177
Entropy 5.1

Neuron H4,W4,C1
JS Divergence 0.322
Entropy 4.0

(a) Before ReLU

(b) After ReLU

## STOA Tradeoff on Imperceptibility versus Success



D2B has higher success rate at the same level of human preference. And with the same success rate, our adversarial examples are consistently more favored by the testers.

## Summary

➤ Identify the scalability problem of existing adversarial attacks.
➤ Proposed the quantile bound on deep content features.
➤ Proposed an efficient way for optimizing the adversarial samples
➤ Show the state-of-the-art trade-off between imperceptibility and attack success using the quantile bound.

## Contact

➤ Email: xu1230@purdue.edu

Home Page:

SCAN ME