**Training image** 



Robust Features: dog Non-Robust Features: dog

Adversarial example towards "cat"

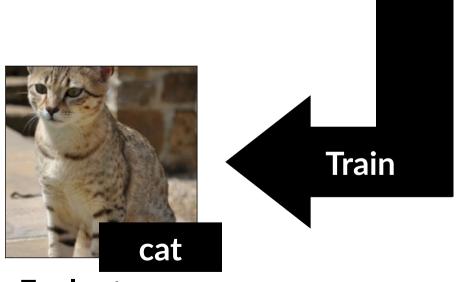
max P(cat)

Relabel as cat



Robust Features: dog Non-Robust Features: cat

good accuracy



Evaluate on original test set