

Finding adversarial examples

© 2019 Philipp Krähenbühl and Chao-Yuan Wu

Finding adversarial examples

- For input \mathbf{x}
- Find ϵ
- Such that
$$f(\mathbf{x} + \epsilon) \neq f(\mathbf{x})$$

Fast gradient sign

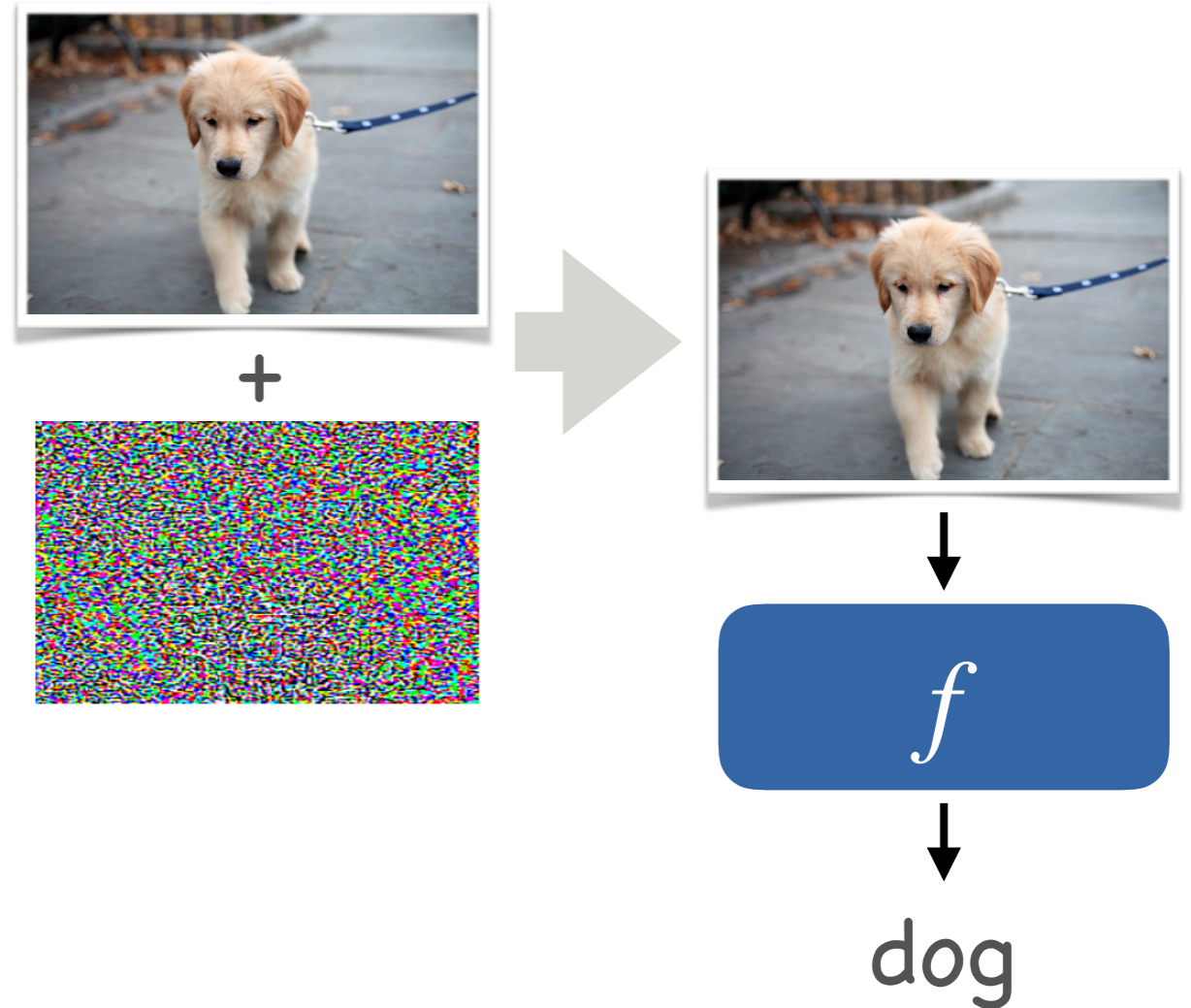
- Assume networks are locally linear
- Optimal attack with $\|\epsilon\|_\infty \leq c$
- $\epsilon = \text{sign} \left(\nabla_{\mathbf{x}} \ell (f(\mathbf{x}), y) \right)$



dog

Projected gradient descent

- Networks are not linear
- Optimize for the attack using gradient descent
- maximize $\ell (f(\mathbf{x} + \epsilon), y)$
- s.t. $\|\epsilon\|_{\infty} < c$



Global adversarial attacks

- Attacks all possible inputs at once
 - PGD on entire dataset
- Attack not input specific
- Attack transfers between architectures
 - Dataset specific?

