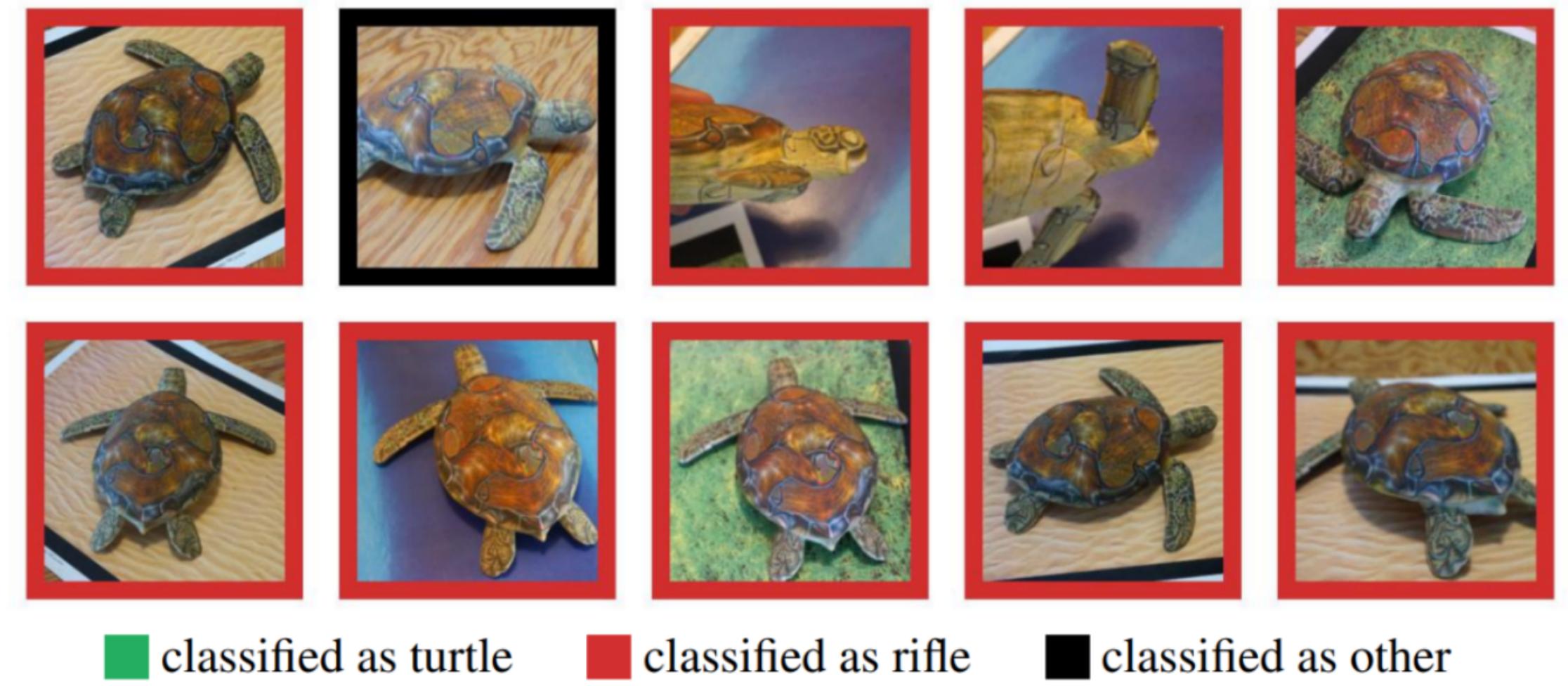


- **Adversarial пример** – вектор, подающийся на вход алгоритму, на котором алгоритм выдает некорректный выход.
- **Adversarial атака** – алгоритм действий, целью которого является получение Adversarial примера.



- Fast Sign Gradient Method

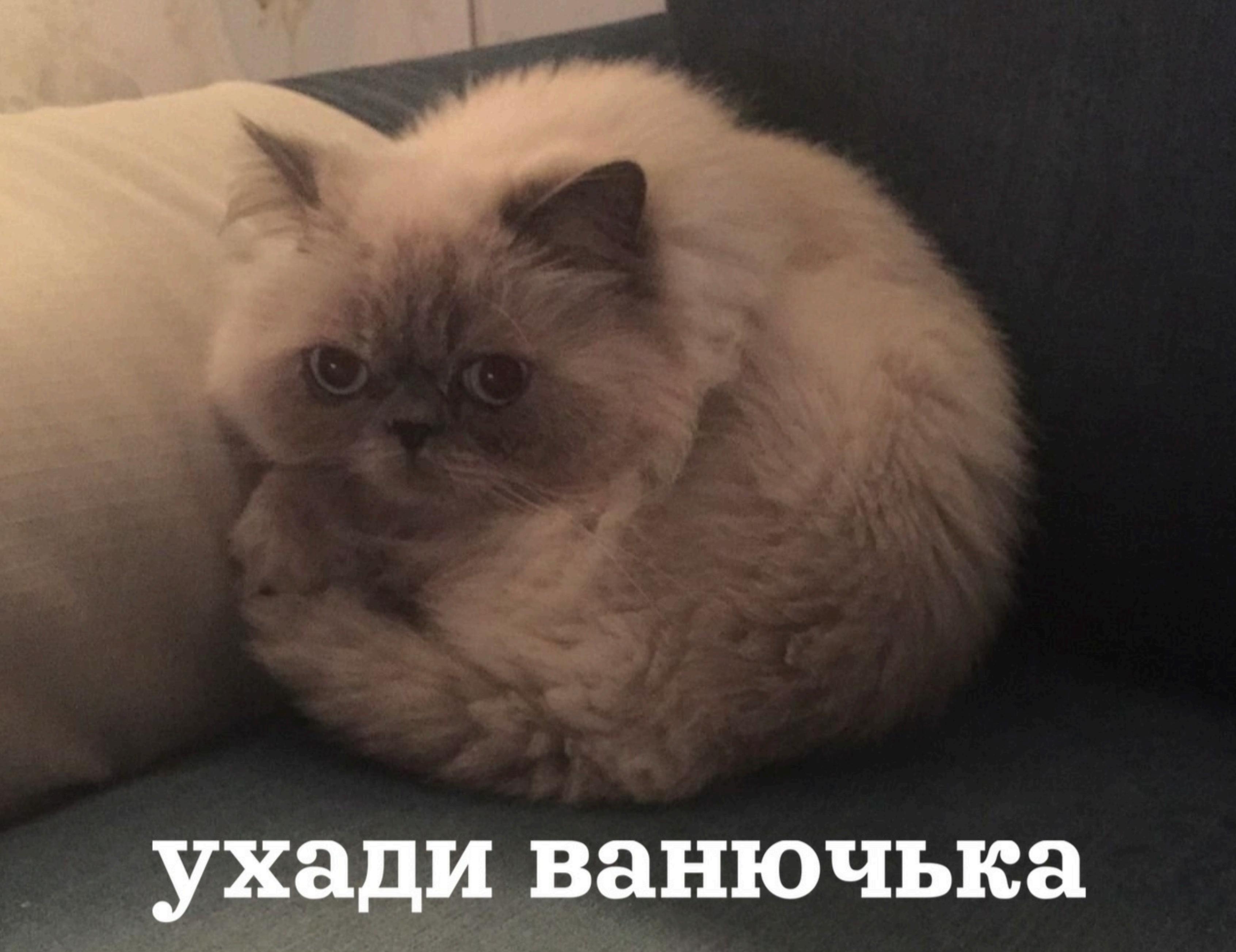
$$X' = X + \epsilon * sign(\nabla_x J(X, y_{true}))$$

- DeepFool атака (про минимальную шумовую карту)
- GAN-based методы (встраивание BlackBox модели в архитектуру генеративно-состязательной сети)

Adversarial Training

- Добавлять adversarial примеры в обучающую выборку
- Добавлять немного зашумленные объекты обучающей выборки
- Label Smoothing
- Построение ансамблей
- Feature squeezing

Немного про тестирование



ухади ванючъка

**Adversarial Examples Are Not
Bugs, They Are Features**

Setup

- binary classification $(x, y) \in \mathcal{X} \times \{\pm 1\}$ $C : \mathcal{X} \rightarrow \{\pm 1\}$
- $\mathcal{F} = \{f : \mathcal{X} \rightarrow \mathbb{R}\}$
- $\mathbb{E}_{(x,y) \sim \mathcal{D}}[f(x)] = 0$
- $\mathbb{E}_{(x,y) \sim \mathcal{D}}[f(x)^2] = 1$

Categorize features

- **ρ -useful features** if it is correlated with the true label in expectation
$$\mathbb{E}_{(x,y) \sim \mathcal{D}} [y \cdot f(x)] \geq \rho.$$
- **γ -robustly useful features** if ρ -useful feature remains γ -useful under adversarial perturbation
$$\mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\inf_{\delta \in \Delta(x)} y \cdot f(x + \delta) \right] \geq \gamma.$$
- **Useful, non-robust features** is a feature which is ρ -useful for some ρ bounded away from zero, but is not a γ -robust feature for any $\gamma \geq 0$

Classification

$$C(x) = \operatorname{sgn} \left(b + \sum_{f \in F} w_f \cdot f(x) \right)$$

Standard Training

- performed by minimizing a loss function (via *empirical risk minimization* (ERM)) that decreases with the correlation between the weighted combination of the features and the label

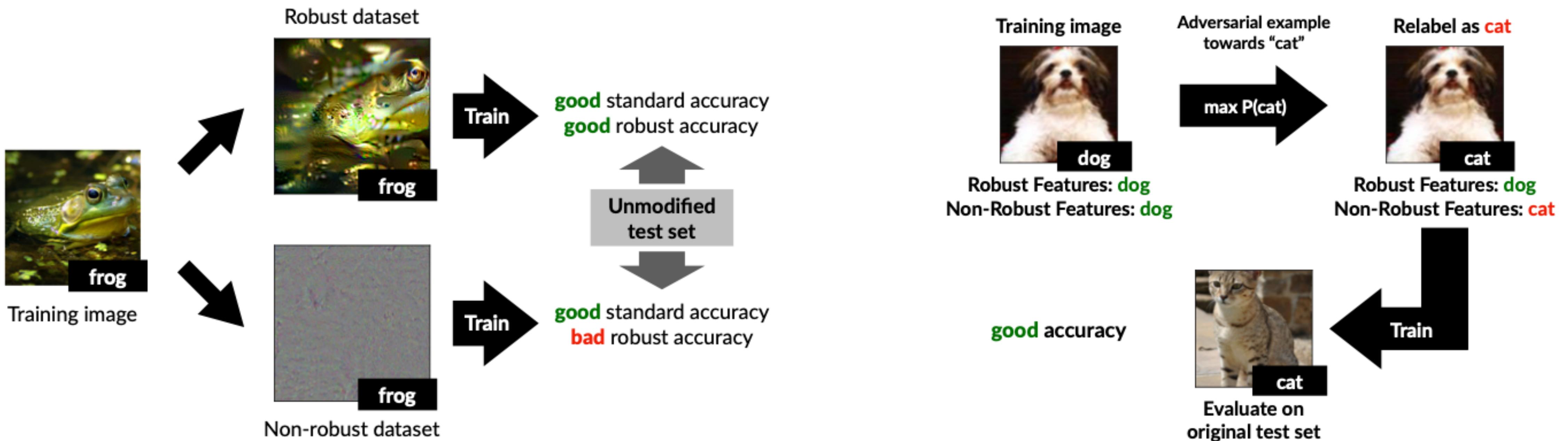
$$\mathbb{E}_{(x,y) \sim \mathcal{D}} [\mathcal{L}_\theta(x, y)] = -\mathbb{E}_{(x,y) \sim \mathcal{D}} \left[y \cdot \left(b + \sum_{f \in F} w_f \cdot f(x) \right) \right]$$

Robust training

- use an *adversarial* loss function that can discern between robust and non-robust features:

$$\mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\max_{\delta \in \Delta(x)} \mathcal{L}_\theta(x + \delta, y) \right]$$

Finding Robust (and Non-Robust) Features



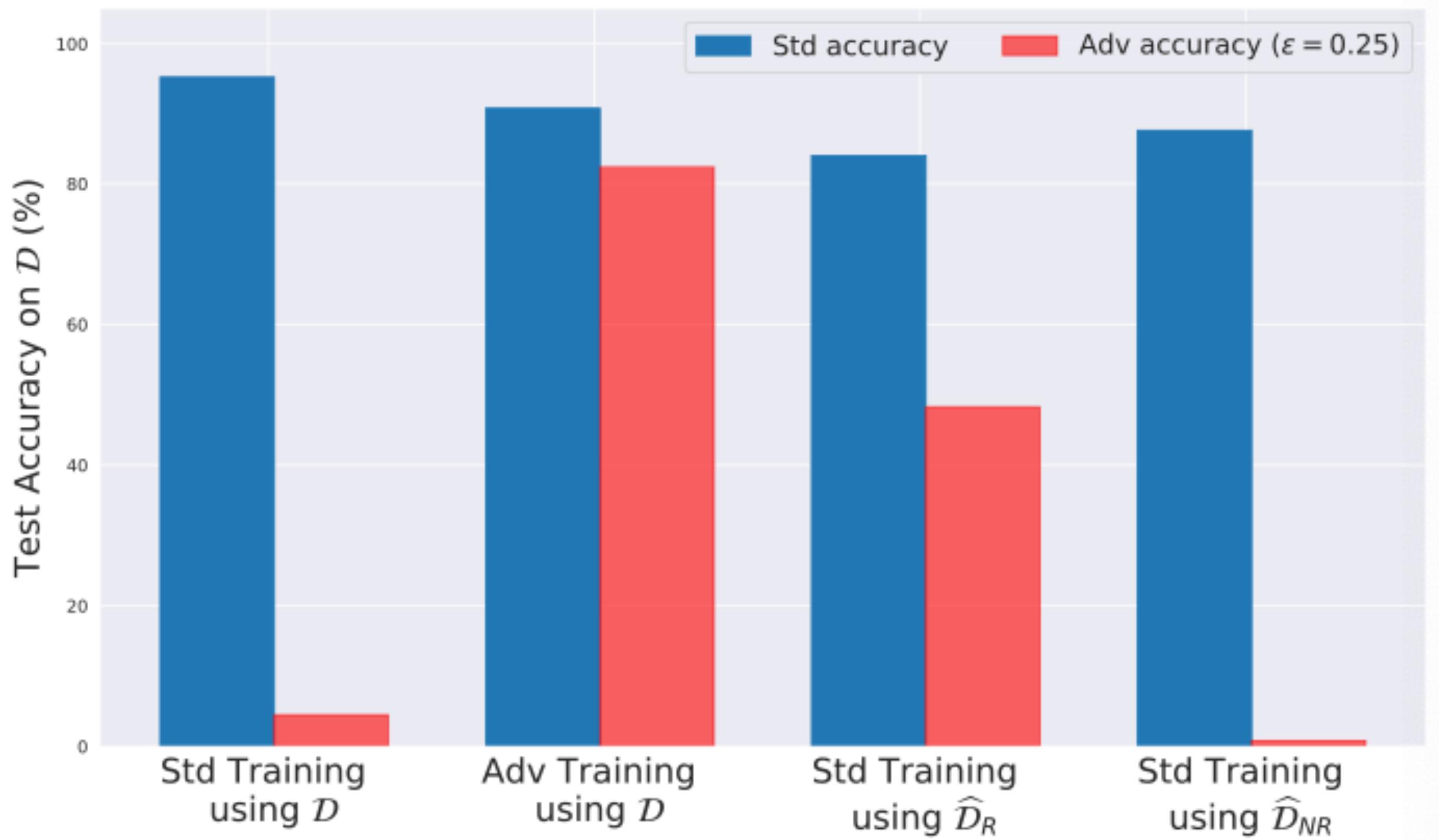
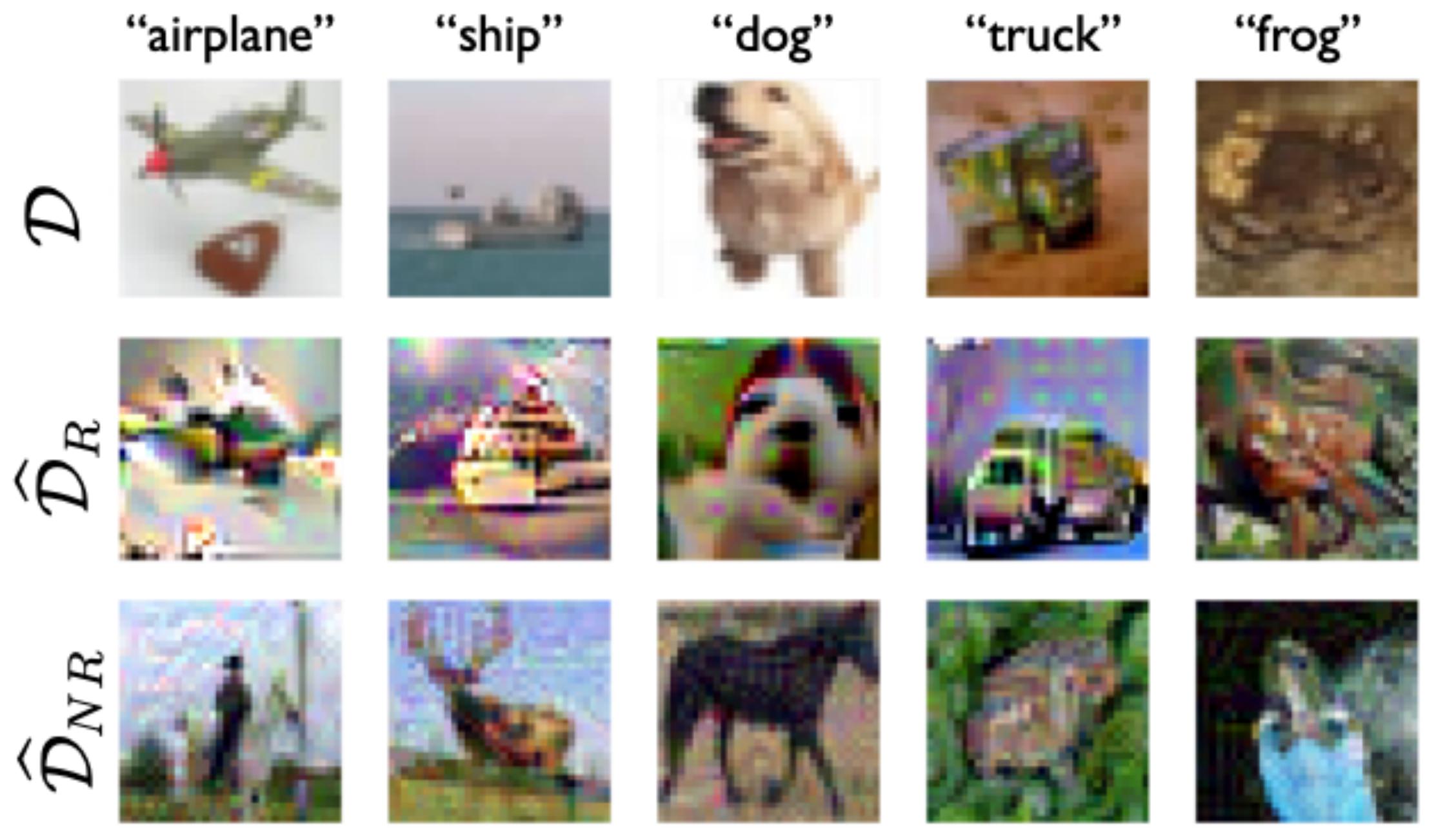
Disentangling robust and non-robust features

- Хочу: гарантировать, что полезны только надежные функции, тогда бы обучение привело к созданию надежного классификатора
- Могу: использовать надежную модель и модифицировать набор данных, чтобы он содержал только те функции, которые имеют отношение к этой модели

$$\mathbb{E}_{(x,y) \sim \hat{\mathcal{D}}_R} [f(x) \cdot y] = \begin{cases} \mathbb{E}_{(x,y) \sim \mathcal{D}} [f(x) \cdot y] & \text{if } f \in F_C \\ 0 & \text{otherwise,} \end{cases}$$

GETROBUSTDATASET(D)

1. $C_R \leftarrow \text{ADVERSARIALTRAINING}(D)$
 $g_R \leftarrow \text{mapping learned by } C_R \text{ from the input to the representation layer}$
2. $D_R \leftarrow \{\}$
3. For $(x, y) \in D$
 $x' \sim D$
 $x_R \leftarrow \arg \min_{z \in [0,1]^d} \|g_R(z) - g_R(x)\|_2$ # Solved using ℓ_2 -PGD starting from x'
 $D_R \leftarrow D_R \cup \{(x_R, y)\}$
4. Return D_R



Non-robust features suffice for standard classification

GETNONROBUSTDATASET(D, ε)

1. $D_{NR} \leftarrow \{\}$

2. $C \leftarrow \text{STANDARDTRAINING}(D)$

3. For $(x, y) \in D$

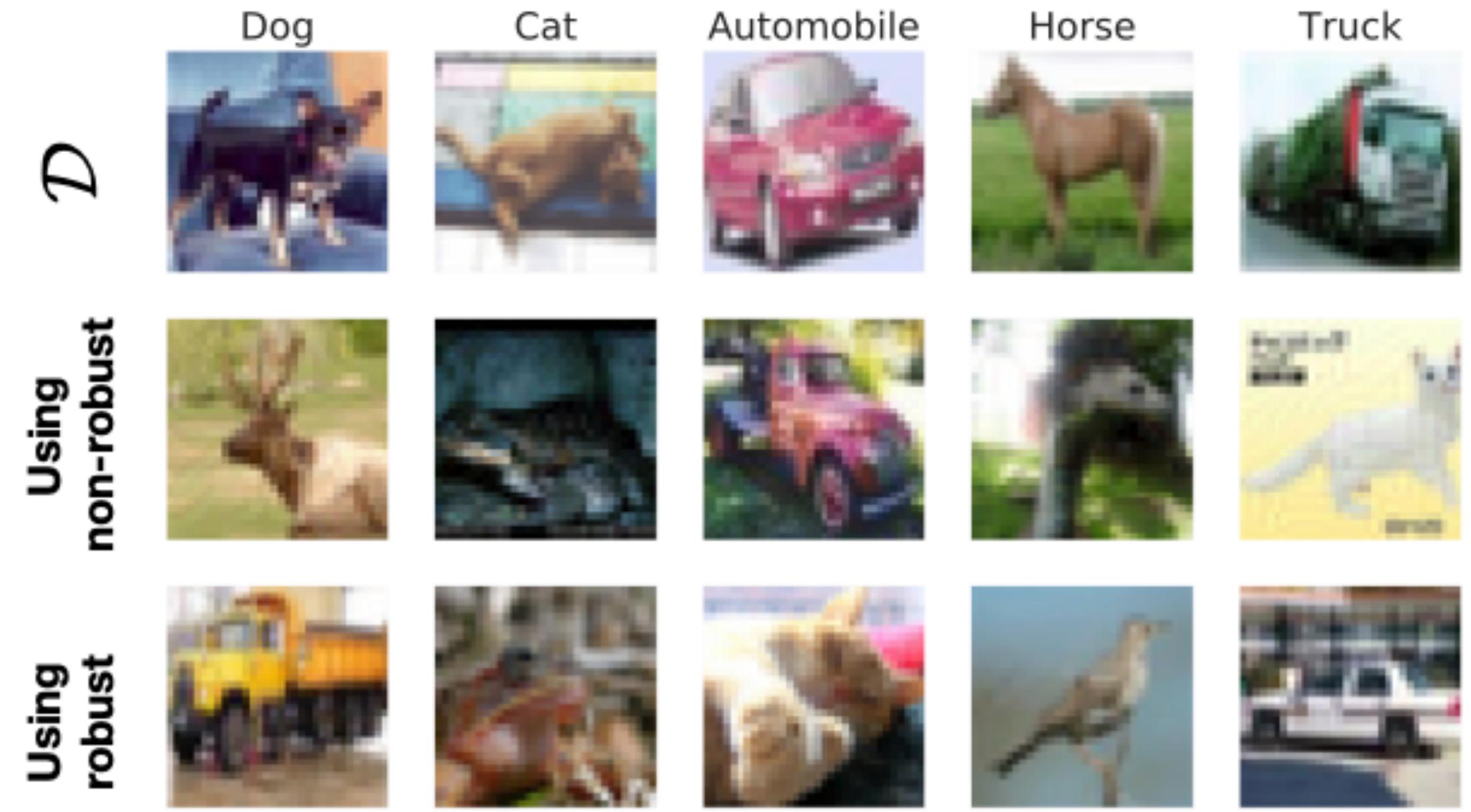
$$t \stackrel{\text{uar}}{\sim} [C] \quad \# \text{ or } t \leftarrow (y + 1) \bmod C$$

$$x_{NR} \leftarrow \min_{||x' - x|| \leq \varepsilon} L_C(x', t) \quad \# \text{ Solved using } \ell_2 \text{ PGD}$$

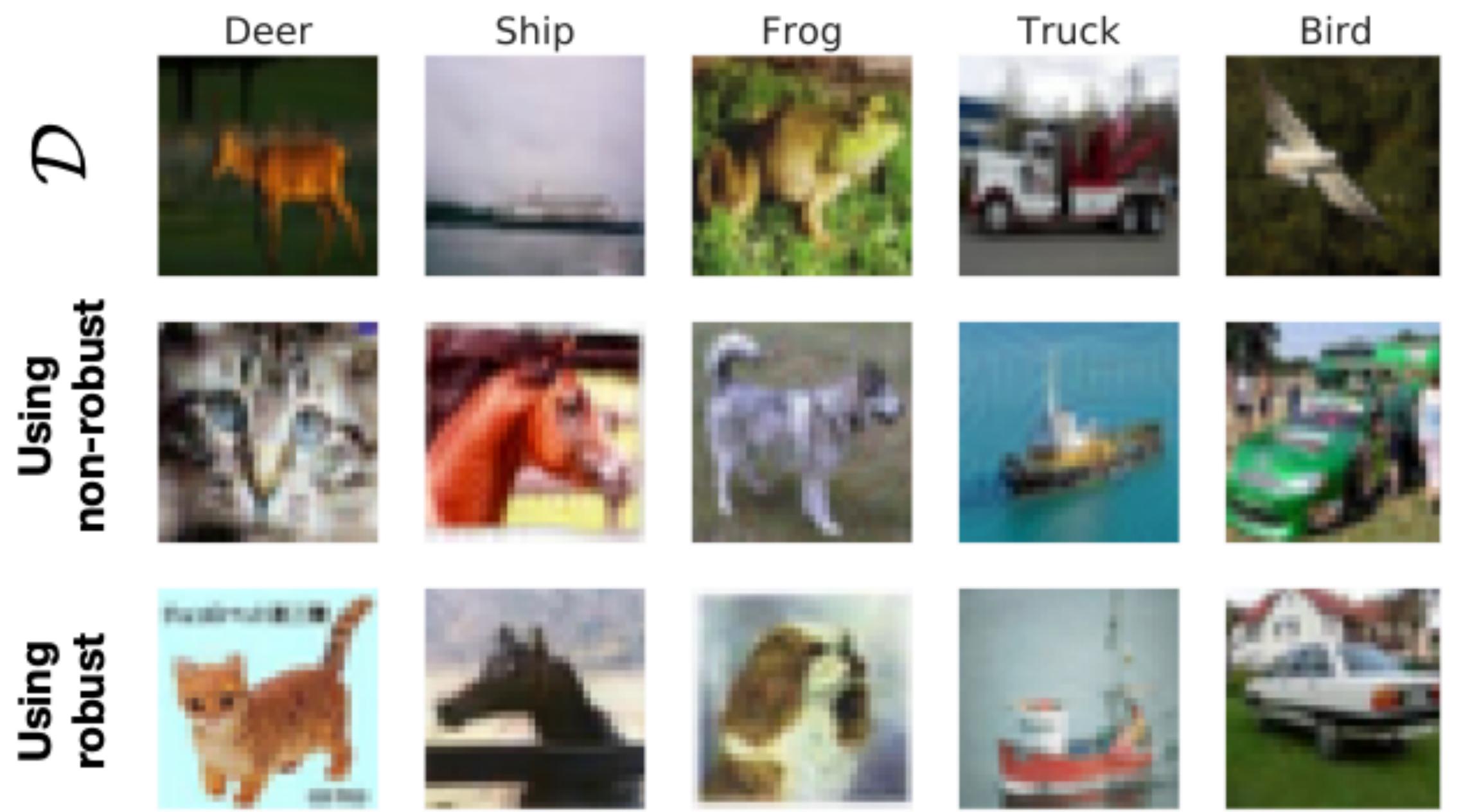
$$D_{NR} \leftarrow D_{NR} \cup \{(x_{NR}, t)\}$$

4. Return D_{NR}

Source Dataset	Dataset	
	CIFAR-10	ImageNet _R
\mathcal{D}	95.3%	96.6%
$\widehat{\mathcal{D}}_{rand}$	63.3%	87.9%
$\widehat{\mathcal{D}}_{det}$	43.7%	64.4%

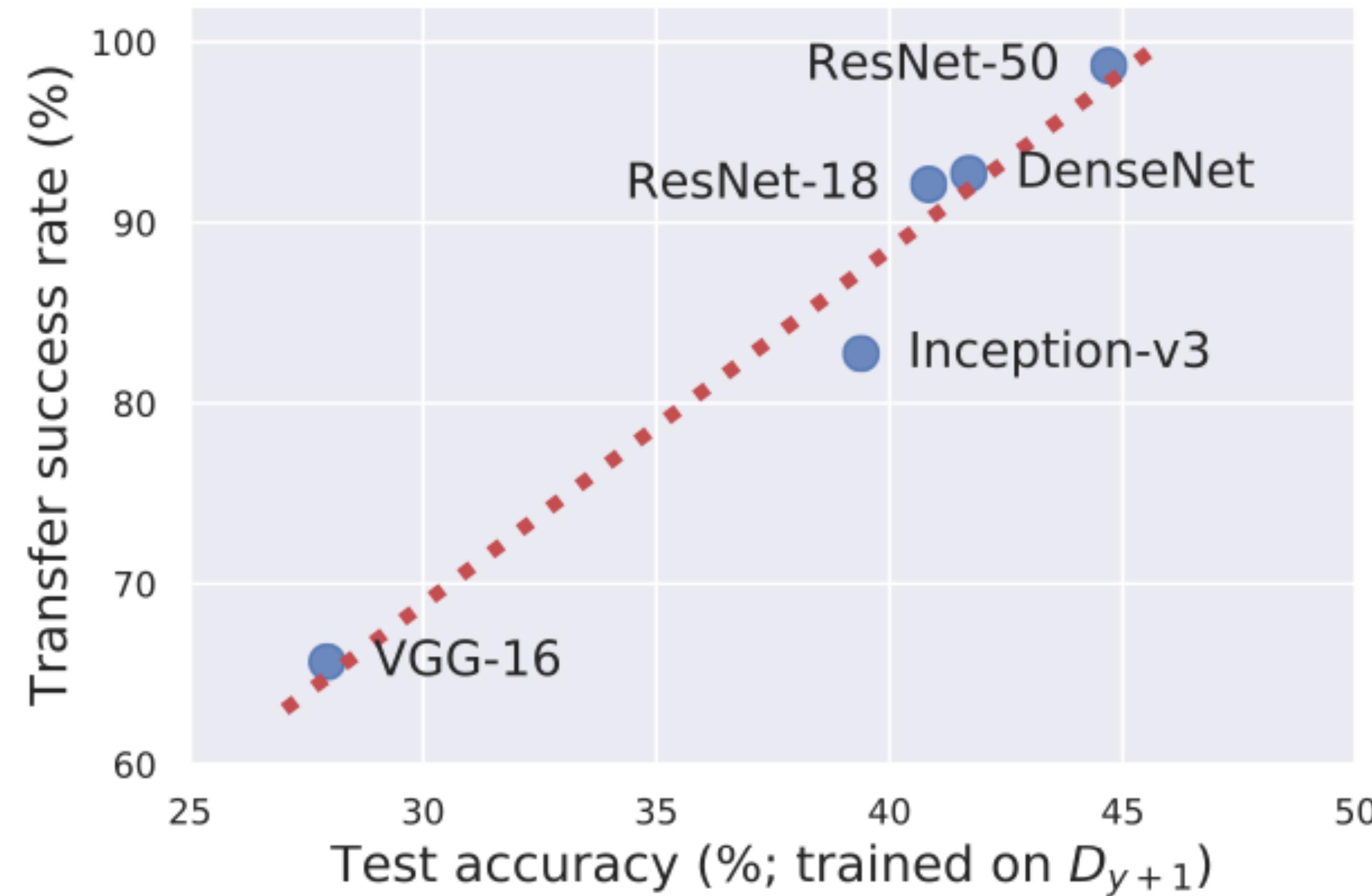


(a) $\widehat{\mathcal{D}}_{rand}$



(b) $\widehat{\mathcal{D}}_{det}$

Transferability can arise from non-robust features



A Theoretical Framework for Studying (Non)-Robust Features

Setup

- maximum likelihood classification between two Gaussian distributions

- $y \stackrel{\text{u.a.r.}}{\sim} \{-1, +1\}, \quad x \sim \mathcal{N}(y \cdot \mu_*, \Sigma_*)$

- Goal: $\Theta = \arg \min_{\mu, \Sigma} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell(x; y \cdot \mu, \Sigma)]$

- $y = \arg \max_y \ell(x; y \cdot \mu, \Sigma) = \text{sign} (x^\top \Sigma^{-1} \mu)$

- $\Theta_r = \arg \min_{\mu, \Sigma} \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\max_{\|\delta\|_2 \leq \varepsilon} \ell(x + \delta; y \cdot \mu, \Sigma) \right]$

Vulnerability from metric misalignment (non-robust features)

Theorem 1 (Adversarial vulnerability from misalignment). *Consider an adversary whose perturbation is determined by the “Lagrangian penalty” form of (12), i.e.*

$$\max_{\delta} \ell(x + \delta; y \cdot \mu, \Sigma) - C \cdot \|\delta\|_2,$$

where $C \geq \frac{1}{\sigma_{\min}(\Sigma_)}$ is a constant trading off NLL minimization and the adversarial constraint¹⁴. Then, the adversarial loss \mathcal{L}_{adv} incurred by the non-robustly learned (μ, Σ) is given by:*

$$\mathcal{L}_{adv}(\Theta) - \mathcal{L}(\Theta) = \text{tr} \left[\left(I + (C \cdot \Sigma_* - I)^{-1} \right)^2 \right] - d,$$

and, for a fixed $\text{tr}(\Sigma_) = k$ the above is minimized by $\Sigma_* = \frac{k}{d} I$.*

Robust Learning

Theorem 2 (Robustly Learned Parameters). *Just as in the non-robust case, $\mu_r = \mu^*$, i.e. the true mean is learned. For the robust covariance Σ_r , there exists an $\varepsilon_0 > 0$, such that for any $\varepsilon \in [0, \varepsilon_0)$,*

$$\Sigma_r = \frac{1}{2}\Sigma_* + \frac{1}{\lambda} \cdot \mathbf{I} + \sqrt{\frac{1}{\lambda} \cdot \Sigma_* + \frac{1}{4}\Sigma_*^2}, \quad \text{where} \quad \Omega\left(\frac{1 + \varepsilon^{1/2}}{\varepsilon^{1/2} + \varepsilon^{3/2}}\right) \leq \lambda \leq O\left(\frac{1 + \varepsilon^{1/2}}{\varepsilon^{1/2}}\right).$$

The effect of robust optimization under an ℓ_2 -constrained adversary is visualized in Figure 4. As ε grows, the learned covariance becomes more aligned with identity. For instance, we can see that the classifier learns to be less sensitive in certain directions, despite their usefulness for natural classification.

Gradient Interpretability

Theorem 3 (Gradient alignment). *Let $f(x)$ and $f_r(x)$ be monotonic classifiers based on the linear separator induced by standard and ℓ_2 -robust maximum likelihood classification, respectively. The maximum angle formed between the gradient of the classifier (wrt input) and the vector connecting the classes can be smaller for the robust model:*

$$\min_{\mu} \frac{\langle \mu, \nabla_x f_r(x) \rangle}{\|\mu\| \cdot \|\nabla_x f_r(x)\|} > \min_{\mu} \frac{\langle \mu, \nabla_x f(x) \rangle}{\|\mu\| \cdot \|\nabla_x f(x)\|}.$$

Конец