

# Adversarial training protects the non-robust features. A trade-off emerges if those features are useful.

## A High Dimensional Statistical Model for Adversarial Training: Geometry and Trade-Offs

### Empirical Risk Minimization

$$\sum_{i=1}^n g \left( y_i \frac{\boldsymbol{\theta}^\top \mathbf{x}_i}{\sqrt{d}} - \varepsilon_t \frac{\sqrt{\boldsymbol{\theta}^\top \boldsymbol{\Sigma}_\delta \boldsymbol{\theta}}}{\sqrt{d}} \right) + r(\boldsymbol{\theta}). \quad (1)$$

0.1 Block Features

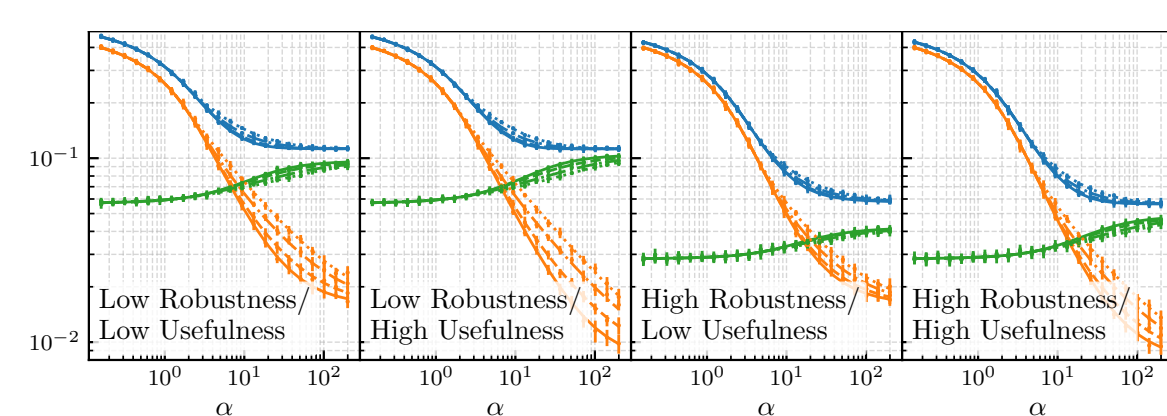
$$\begin{aligned} \boldsymbol{\Sigma}_x &= \text{blockdiag}(\psi_1 \mathbb{1}_{d_1}, \dots, \psi_k \mathbb{1}_{d_k}), \\ \boldsymbol{\Sigma}_\delta &= \text{blockdiag}(\Delta_1 \mathbb{1}_{d_1}, \dots, \Delta_k \mathbb{1}_{d_k}), \\ \boldsymbol{\Sigma}_v &= \text{blockdiag}(\Upsilon_1 \mathbb{1}_{d_1}, \dots, \Upsilon_k \mathbb{1}_{d_k}), \\ \boldsymbol{\Sigma}_\theta &= \text{blockdiag}(t_1 \mathbb{1}_{d_1}, \dots, t_k \mathbb{1}_{d_k}), \end{aligned} \quad (2)$$

0.2 Usefulness and robustness

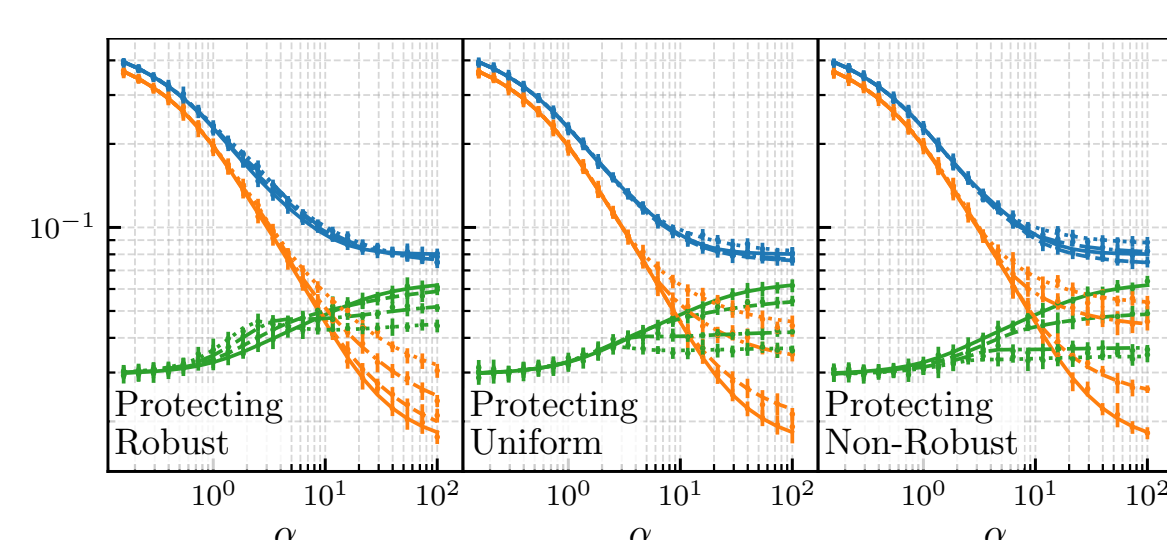
$$\begin{aligned} \mathcal{U}_{\theta_0} &= \frac{1}{\sqrt{d}} \mathbb{E}_{\mathbf{x}, y} [y \boldsymbol{\theta}_0^\top \mathbf{x}], \\ \mathcal{R}_{\theta_0} &= \frac{1}{\sqrt{d}} \mathbb{E}_{\mathbf{x}, y} \left[ \inf_{\|\boldsymbol{\delta}\|_{\boldsymbol{\Sigma}_v^{-1}} \leq \varepsilon_g} y \boldsymbol{\theta}_0^\top (\mathbf{x} + \boldsymbol{\delta}) \right]. \end{aligned} \quad (3) \quad (4)$$

### Main Theorem

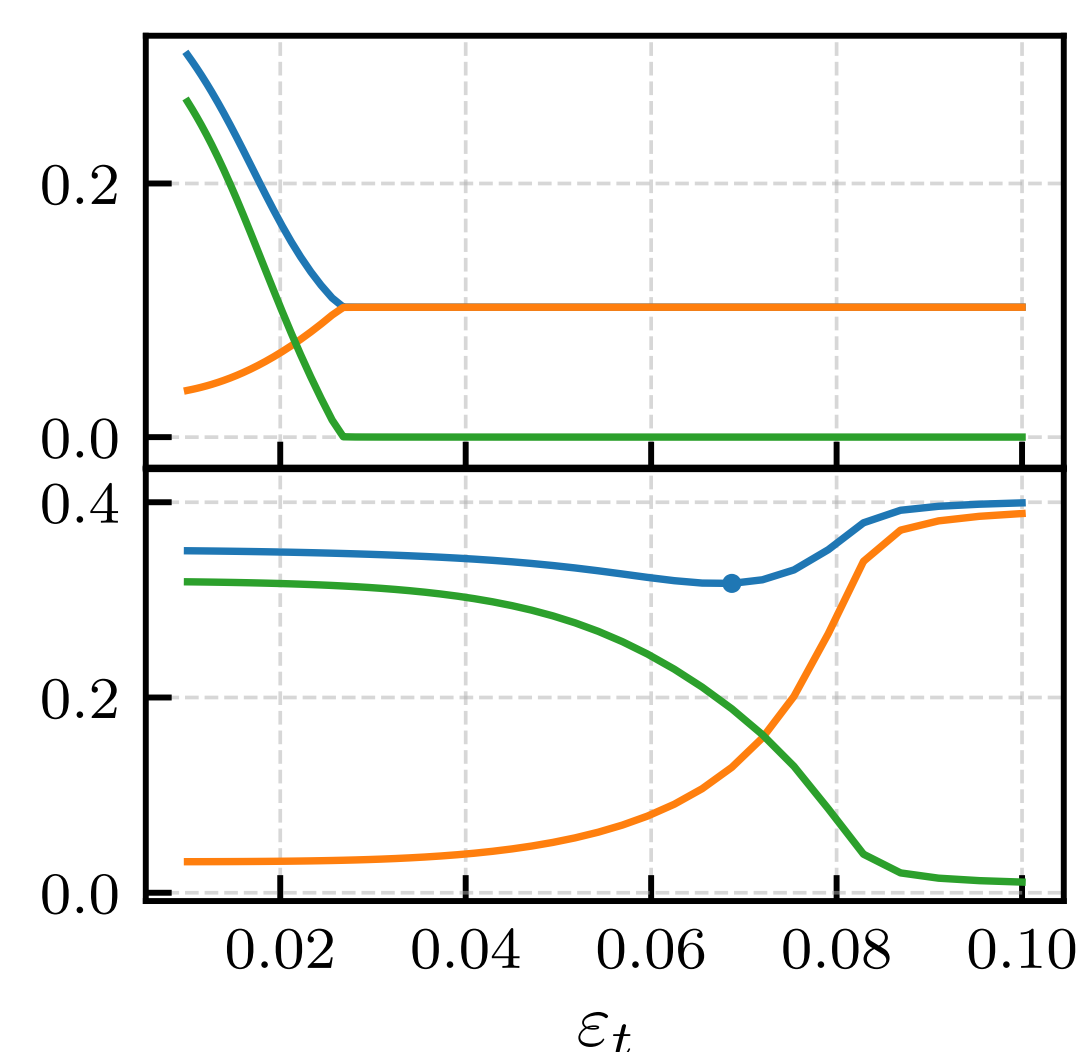
### Trade-Offs



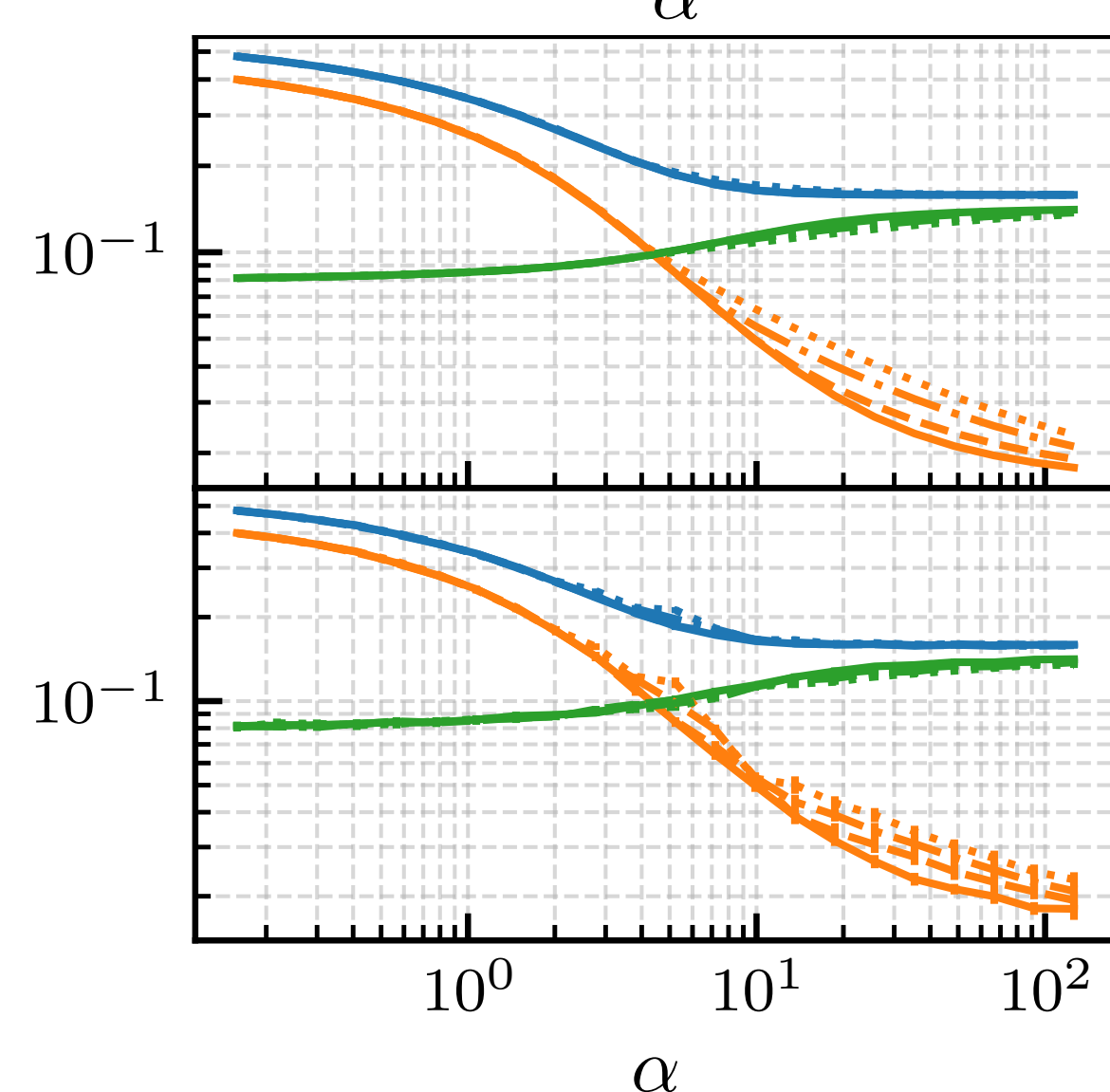
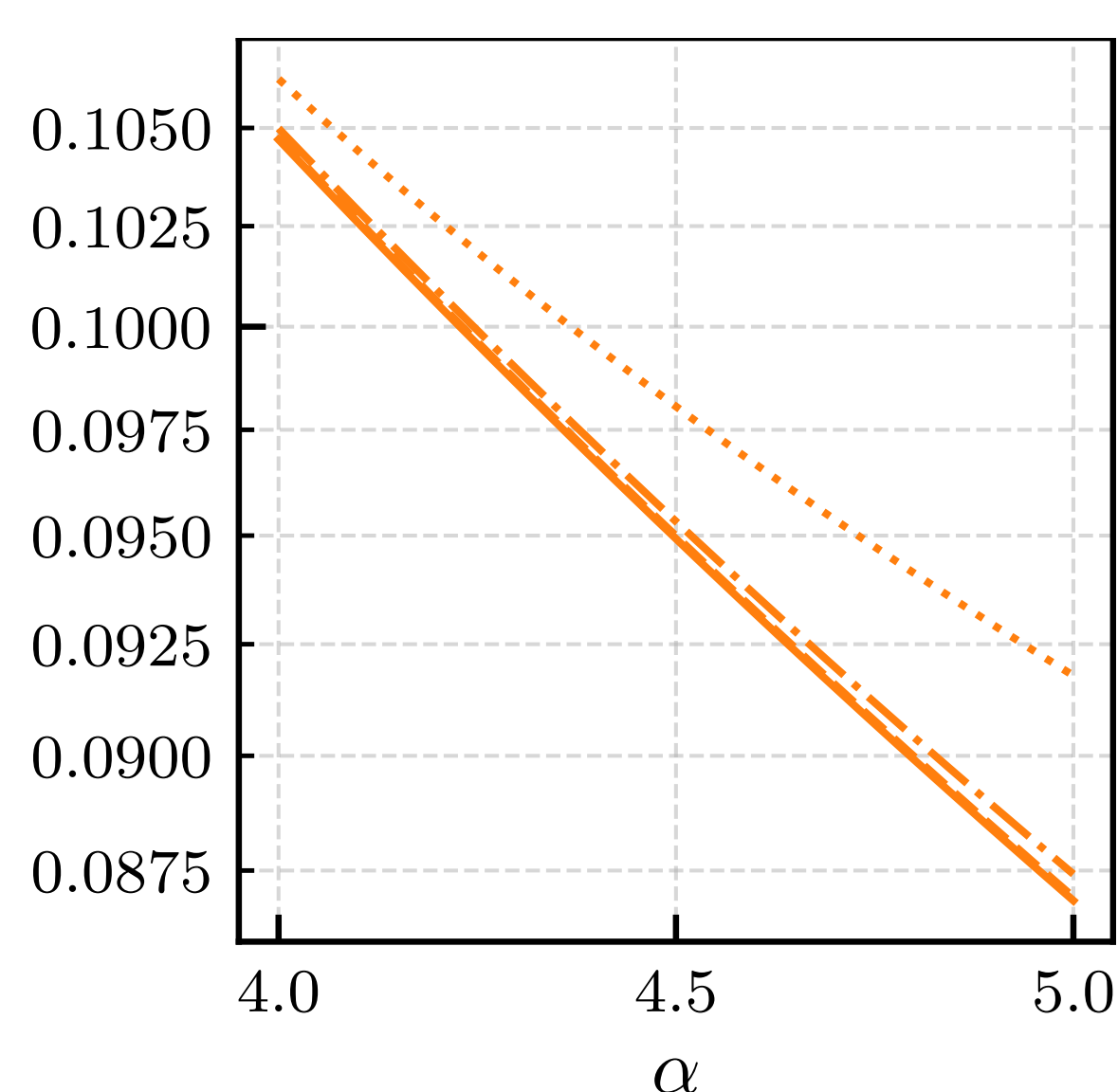
0.3 Directional Defences and structured data



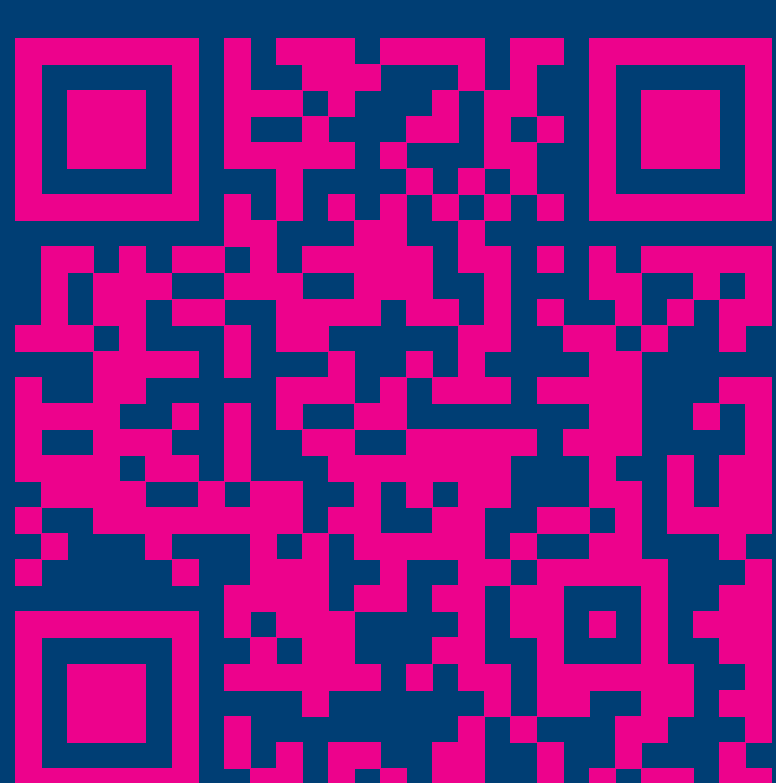
0.4 Tradeoff directions and innocuous directions



0.5 Data Dependent Regularisation



### Acknowledgements



Kasimir Tanner  
Matteo Vilucchio  
Bruno Loureiro  
Florent Krzakala

EPFL