# Adversarial Debiasing Architecture



**Input X** → **Predictor Pθ** (Maximize accuracy) → **Output $\hat{Y}$**

*Gradient Reversal*

**Adversary Aφ** (Predict A) → **Pred $\hat{A}$**

**Predictor Loss:**

$$L_P = -\text{Acc}(Y, \hat{Y})$$

**Adversary Loss:**

$$L_A = -\text{Acc}(A, \hat{A})$$

**Combined:**

$$\min_{\theta}\max_{\phi} L_P - \lambda L_A$$