

Model Immunization from a Condition Number Perspective

ICML 2025

Amber Yijia Zheng*



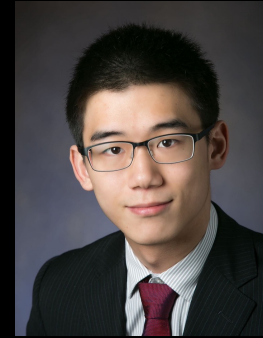
Cedar Site Bai*



Brian Bullins



Raymond A. Yeh



OPEN-WEIGHT GENERATIVE AI

Stable
Diffusion



Deep
Floyd IF



Content generation just a few prompts away!

RESPONSIBILITY



RISKS OF MISUSE

X blocks searches for Taylor Swift after explicit AI images of her go viral

28th January 2024, 12:42 EST

Share

By Nadine Yousif
BBC News



Getty Images

Social media platform X has blocked searches for Taylor Swift after explicit AI-generated images of the singer began circulating on the site.

The New York Times

People Love Studio Ghibli. But Should They Be Able to Recreate It?

An update to ChatGPT made it easy to simulate Hayao Miyazaki's style of animation, which has flooded social media with memes.



BBC

AI can be easily used to make fake election photos - report

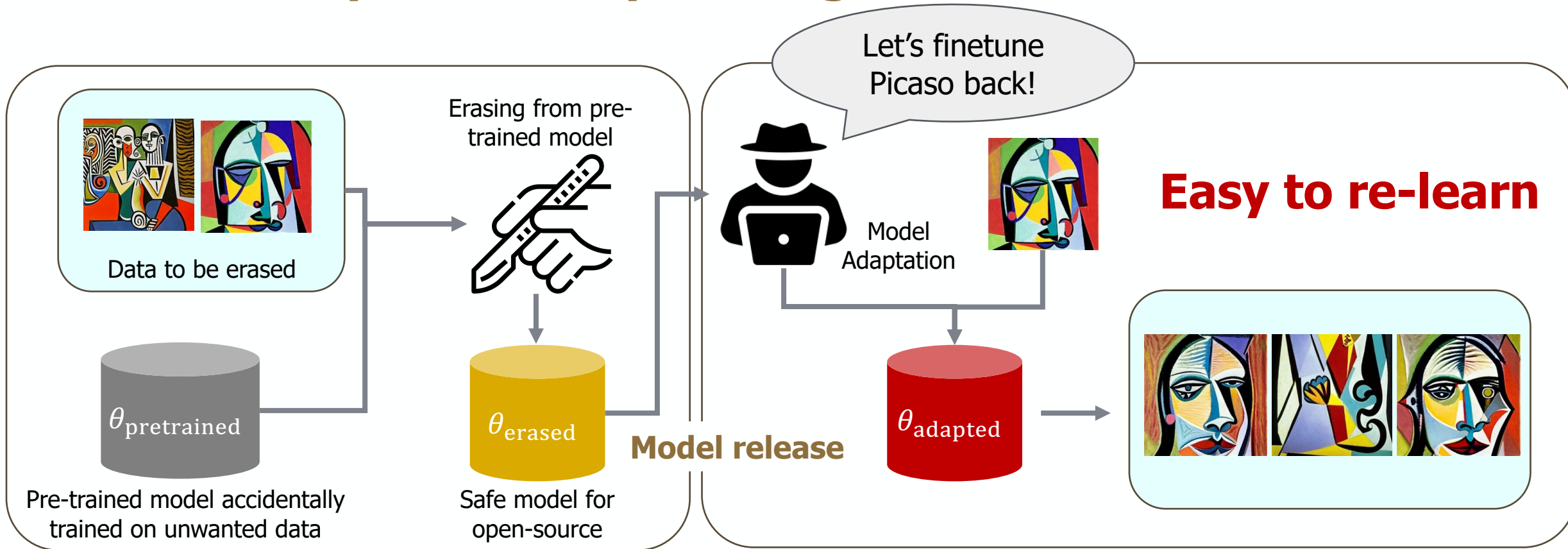
This fake image of a man lurking outside a polling place with a gun was created by artificial intelligence tool ChatGPT Plus.

Mar 6, 2024



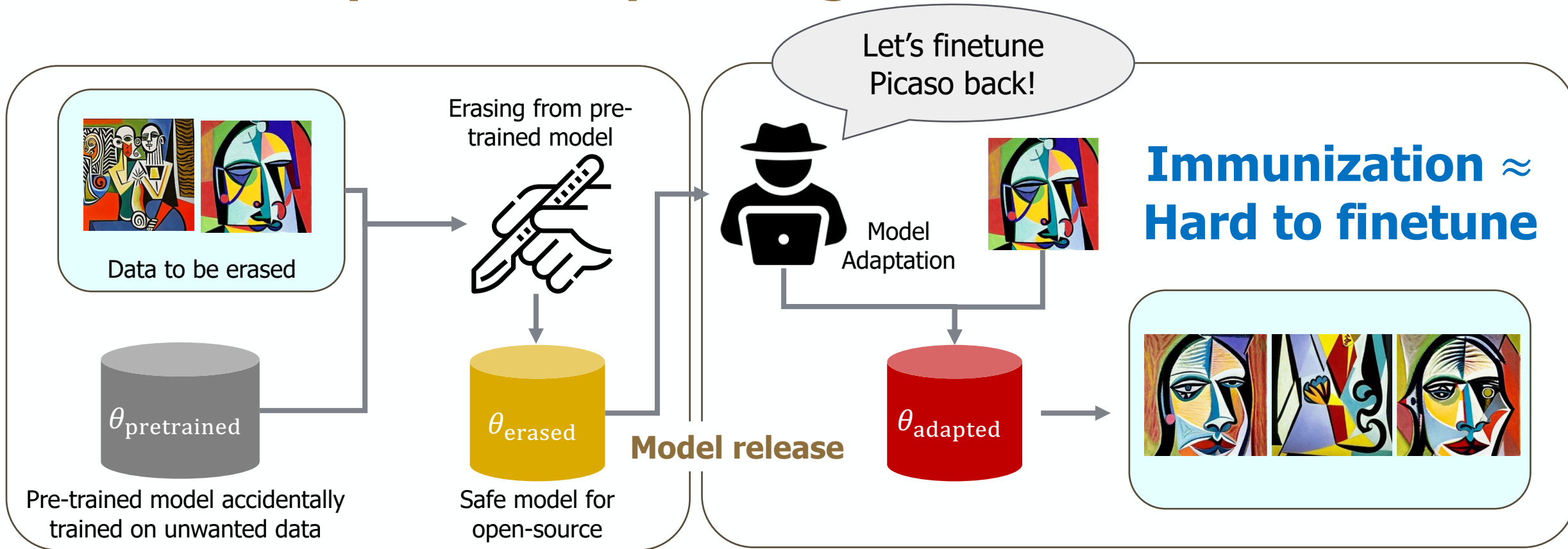
THREAT MODEL

Malicious adaptation of open-weight models



THREAT MODEL

Malicious adaptation of open-weight models



IMMUNIZATION FRAMEWORK

Goal: learn a pre-trained feature extractor, such that fine-tuning an adapter $h_{\mathbf{w}}$ on features of the **harmful task is difficult** but not for other tasks.

Setting: Linear feature extractor $f_{\theta} \triangleq \mathbf{x}^{\top} \theta$

- Linear adapter g_{ω} trained on Pre-training task $\mathcal{D}_{\text{P}} = \{(\mathbf{x}, \mathbf{y})\}$
- Linear adapter $h_{\mathbf{w}}$ trained on Harmful task $\mathcal{D}_{\text{H}} = \{(\mathbf{x}, \tilde{\mathbf{y}})\}$
- Regression task:

$$\mathcal{L}(\mathcal{D}, \mathbf{w}, \theta) \triangleq \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} \ell_2(h_{\mathbf{w}} \circ f_{\theta}(\mathbf{x}), \mathbf{y})$$

- The bad actor performs **linear probing** on \mathcal{D}_{H} by $\min_{\mathbf{w}} \mathcal{L}(\mathcal{D}, \mathbf{w}, \theta)$
 - Only training the last linear layer.

CONDITION NUMBER

Condition number of a general matrix \mathbf{S}

- The ratio between max / min singular values

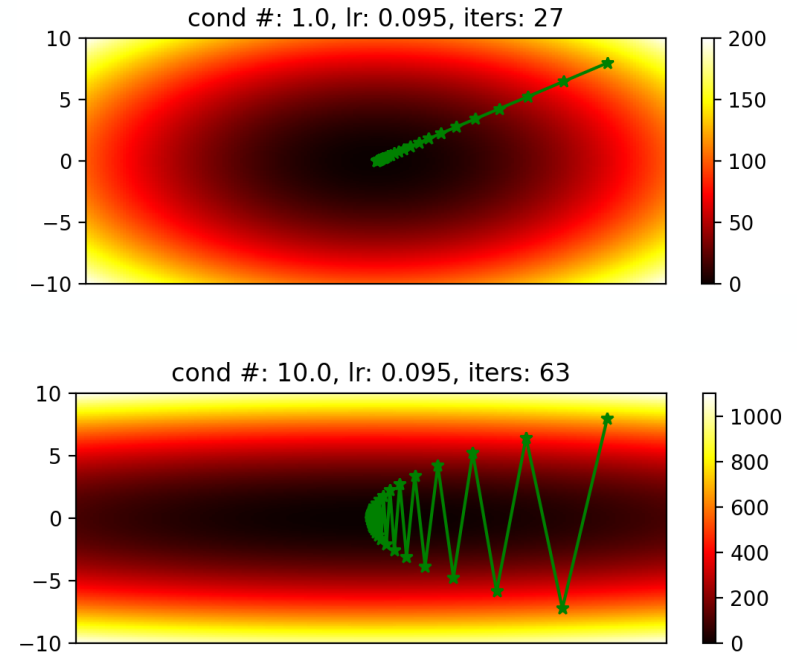
$$\kappa(\mathbf{S}) \triangleq \|\mathbf{S}\|_2 \|\mathbf{S}^\dagger\|_2 = \sigma_{\mathbf{S}}^{\max} / \sigma_{\mathbf{S}}^{\min}$$

For constant step-size steepest descent

- $\min_{\mathbf{w}} \mathcal{L}(\mathbf{w})$, \mathcal{L} is strongly convex
- Hessian $\nabla^2 \mathcal{L}$ with max/min singular values $\sigma^{\max/\min}$
- Convergence rate (Bubeck, 2015) :

$$\|\mathbf{w}_t - \mathbf{w}^*\|^2 \leq \left(1 - \frac{\sigma_{\min}}{\sigma_{\max}}\right)^t \|\mathbf{w}_0 - \mathbf{w}^*\|^2$$

- Large condition number \rightarrow slower convergence \rightarrow hard to finetune



A CONDITION NUMBER PERSPECTIVE

Definition of Immunization

(a) It is **more difficult** to apply linear probing on the harmful task \mathcal{D}_H using the immunized feature extractor f_{θ^I} than directly on the input data,

$$\kappa(\nabla_{\mathbf{w}}^2 \mathcal{L}(\mathcal{D}_H, \mathbf{w}, \theta^I)) \gg \kappa(\nabla_{\mathbf{w}}^2 \mathcal{L}(\mathcal{D}_H, \mathbf{w}, \mathbf{I}))$$

\uparrow Harmful task \uparrow Identity

(b) It is **not more difficult** to apply linear probing on other tasks. As there is only one other task \mathcal{D}_P , an immunized feature extractor should have

$$\kappa(\nabla_{\omega}^2 \mathcal{L}(\mathcal{D}_P, \omega, \theta^I)) \leq \kappa(\nabla_{\omega}^2 \mathcal{L}(\mathcal{D}_P, \omega, \mathbf{I})).$$

\uparrow Pre-training task

(c) The immunized model should maintain a competitive task performance on the pre-training dataset \mathcal{D}_P ,

$$\min_{\omega, \theta} \mathcal{L}(\mathcal{D}_P, \omega, \theta) \approx \min_{\omega} \mathcal{L}(\mathcal{D}_P, \omega, \theta^I).$$

WHEN IMMUNIZATION IS POSSIBLE

Everything is linear → nice analytical forms

Linear Probing:

$$\mathcal{L}(\mathcal{D}_H, \mathbf{w}, \theta) = \min_{\mathbf{w}} \|(\mathbf{X}_H \theta) \mathbf{w} - \mathbf{Y}\|_2^2$$

The Hessian Matrix:

$$\mathbf{H}_H(\theta) = \nabla_{\mathbf{w}}^2 \mathcal{L}(\mathcal{D}_H, \mathbf{w}, \theta) = \theta^\top \mathbf{K}_H \theta$$

$$\text{with } \mathbf{K}_H = \mathbf{X}_H^\top \mathbf{X}_H$$

Singular value is given by:

$$\sigma_i = \sum_{j=1}^{D_{\text{in}}} \left(\sigma_{\theta,i} (\mathbf{u}_{\theta,i}^\top \mathbf{q}_j) \sqrt{\gamma_j} \right)^2, \quad \forall i \in \{1, \dots, D^{\text{in}}\}$$

Here, $\sigma_{\theta,i}$ and $\mathbf{u}_{\theta,i}$ correspond to the i -th singular value and vector of θ . Next, γ_j and \mathbf{q}_j correspond to the j -th singular value and vector of the covariance \mathbf{K} .

Immunization depends on the **“relative angle”** between singular vectors of the covariance matrix for pre-training/harmful task.

How to train such a feature extractor θ ?

OPTIMIZING CONDITION NUMBER

$$\min_{\omega, \theta} \mathcal{R}_{\text{ill}}(\mathbf{H}_{\text{H}}(\theta)) + \mathcal{R}_{\text{well}}(\mathbf{H}_{\text{P}}(\theta)) + \mathcal{L}(\mathcal{D}_{\text{P}}, \omega, \theta)$$

$$\mathcal{R}_{\text{ill}}(\mathbf{S}) \triangleq \frac{1}{\frac{1}{2k} \|\mathbf{S}\|_F^2 - \frac{1}{2} \sigma_{\mathbf{S}}^{\min 2}}$$

immunizes the model on harmful task

$$\kappa(\nabla_{\mathbf{w}}^2 \mathcal{L}(\mathcal{D}_{\text{H}}, \mathbf{w}, \theta^{\text{I}})) \gg \kappa(\nabla_{\mathbf{w}}^2 \mathcal{L}(\mathcal{D}_{\text{H}}, \mathbf{w}, \mathbf{I}))$$

$$\mathcal{R}_{\text{well}}(\mathbf{S}) \triangleq \frac{1}{2} \|\mathbf{S}\|_2^2 - \frac{1}{2p} \|\mathbf{S}\|_F^2$$

maintains the fine-tuning ability on other tasks

$$\kappa(\nabla_{\omega}^2 \mathcal{L}(\mathcal{D}_{\text{P}}, \omega, \theta^{\text{I}})) \leq \kappa(\nabla_{\omega}^2 \mathcal{L}(\mathcal{D}_{\text{P}}, \omega, \mathbf{I})).$$

$$\mathcal{L}(\mathcal{D}, \omega, \theta) \triangleq \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} \ell_2(g_{\omega} \circ f_{\theta}(\mathbf{x}), \mathbf{y})$$

maintains the model performance

$$\min_{\omega, \theta} \mathcal{L}(\mathcal{D}_{\text{P}}, \omega, \theta) \approx \min_{\omega} \mathcal{L}(\mathcal{D}_{\text{P}}, \omega, \theta^{\text{I}}).$$

PROPOSED OBJECTIVE IS “NICE”

[Upper Bound] $\frac{1}{\log \kappa(\mathbf{S})} \leq (\sigma_{\mathbf{S}}^{\max})^2 \mathcal{R}_{\text{i11}}(\mathbf{S})$, i.e., $\mathcal{R}_{\text{i11}}(\mathbf{S})$ upper bounds $\frac{1}{\log \kappa(\mathbf{S})}$ when $\sigma_{\mathbf{S}}^{\max}$ is reasonably away from ∞ .

[Differentiability] If $\sigma_{\mathbf{S}}^{\min} = \sigma_k < \sigma_i$ for any $i < k$, i.e., $\sigma_{\mathbf{S}}^{\min}$ is unique, then $\mathcal{R}_{\text{i11}}(\mathbf{S})$ is differentiable and

$$\nabla_{\mathbf{S}} \mathcal{R}_{\text{i11}}(\mathbf{S}) = \frac{\sigma_k \mathbf{u}_k \mathbf{v}_k^\top - \frac{1}{k} \mathbf{S}}{(\frac{1}{2k} \|\mathbf{S}\|_F^2 - \frac{1}{2} (\sigma_{\mathbf{S}}^{\min})^2)^2}.$$

[Monotonic Increase] If $\sigma_{\mathbf{S}}^{\min}$ is unique, update \mathbf{S} with $\nabla_{\mathbf{S}} \mathcal{R}_{\text{i11}}(\mathbf{S})$ such that $\mathbf{S}' = \mathbf{S} - \eta_2 \nabla_{\mathbf{S}} \mathcal{R}_{\text{i11}}(\mathbf{S})$ for $0 < \eta_2 < \frac{k}{k-1} (\frac{1}{2k} \|\mathbf{S}\|_F^2 - \frac{1}{2} (\sigma_{\mathbf{S}}^{\min})^2)^2$, then $\kappa(\mathbf{S}') > \kappa(\mathbf{S})$.

OPTIMIZING CONDITION NUMBER

$$\min_{\omega, \theta} \mathcal{R}_{\text{ill}}(\mathbf{H}_H(\theta)) + \mathcal{R}_{\text{well}}(\mathbf{H}_P(\theta)) + \mathcal{L}(\mathcal{D}_P, \omega, \theta)$$

Solve using gradient based method!

Algorithm 1 Condition number regularized gradient descent for model immunization

input Primary task $\mathcal{D}_P = (\mathbf{X}_P, \mathbf{Y}_P)$, harmful task input \mathbf{X}_H , supervised loss \mathcal{L} , learning rate η , regularizing constants $\lambda_P, \lambda_H \in \mathbb{R}_+$, model initialization θ_0, ω_0

1: $\mathbf{K}_P = \mathbf{X}_P^\top \mathbf{X}_P$

2: $\mathbf{K}_H = \mathbf{X}_H^\top \mathbf{X}_H$

3: **for** $t = 0, 1, \dots, T - 1$ **do**

4: $\omega_{t+1} = \omega_t - \eta \nabla_{\omega} \mathcal{L}(\omega_t, \theta_t; \mathcal{D}_P)$

5: $\mathbf{H}_P(\theta_t) = \theta_t^\top \mathbf{K}_P \theta_t, \mathbf{H}_H(\theta_t) = \theta_t^\top \mathbf{K}_H \theta_t$

6: $\theta_{t+1} = \theta_t - \eta \nabla_{\theta} \mathcal{L}(\omega_t, \theta_t; \mathbf{X}_1)$
 $\quad - \eta \lambda_P \mathbf{K}_P^{-1} \nabla_{\theta} \mathcal{R}_{\text{well}}(\mathbf{H}_P(\theta_t))$
 $\quad - \eta \lambda_H \mathbf{K}_H^{-1} \nabla_{\theta} \mathcal{R}_{\text{ill}}(\mathbf{H}_H(\theta_t))$

7: **end for**

output Immunized feature extractor $\theta_I \triangleq \theta_T$.

EXPERIMENTS

Evaluation metrics: relative immunization ratio

$$\text{RIR} = \underbrace{\left(\frac{\kappa(\mathbf{H}_H(\theta_I))}{\kappa(\mathbf{H}_H(I))} \right)}_{\text{(i)}} \bigg/ \underbrace{\left(\frac{\kappa(\mathbf{H}_P(\theta_I))}{\kappa(\mathbf{H}_P(I))} \right)}_{\text{(ii)}}$$

- A successful immunization has large (i) and small (ii)

Baselines:

1. \mathcal{R}_{ill} **Only** immunizes the model by minimizing only the regularizer $\mathcal{R}_{\text{ill}}(\mathbf{H}_H)$
2. **IMMA** [Zheng and Yeh, 2024]
3. **Opt** κ directly minimizes $\kappa(\mathbf{H}_P(\theta)) - \kappa(\mathbf{H}_H(\theta))$ via gradient descent instead of using our proposed regularizers.

REGRESSION TASK (House Price Dataset)

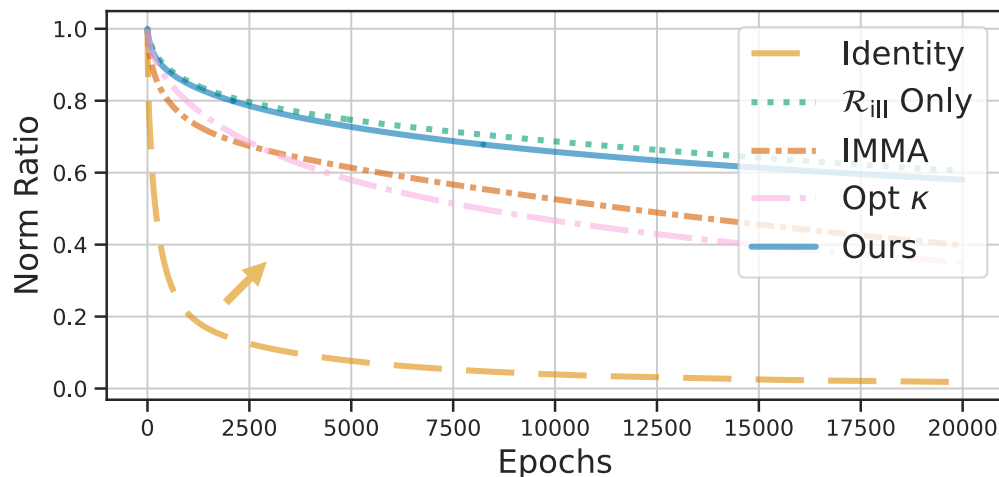
The feature extractor is a linear model

$$\text{RIR} = \underbrace{\left(\frac{\kappa(\mathbf{H}_H(\theta_I))}{\kappa(\mathbf{H}_H(\mathbf{I}))} \right)}_{(i)} / \underbrace{\left(\frac{\kappa(\mathbf{H}_P(\theta_I))}{\kappa(\mathbf{H}_P(\mathbf{I}))} \right)}_{(ii)}$$

$$\text{Norm Ratio} = \|\mathbf{w}_t - \mathbf{w}^*\|_2^2 / \|\mathbf{w}_0 - \mathbf{w}^*\|_2^2$$

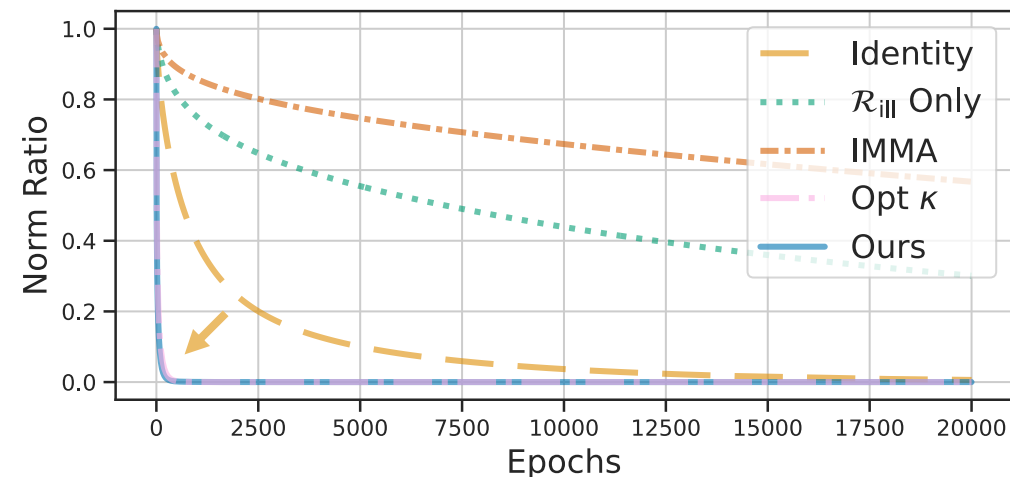
Method	RIR (i)↑	RIR (ii)↓	RIR ↑
\mathcal{R}_{ill} Only	90.02 ± 3.773	72.415 ± 3.545	1.244 ± 0.021
IMMA	7.053 ± 1.662	3.545 ± 0.880	2.001 ± 0.187
Opt κ	1.518 ± 0.027	0.016 ± 0.001	92.58 ± 4.492
Ours	18.92 ± 2.056	0.053 ± 0.002	356.20 ± 5.491

Norm ratio curve on \mathcal{D}_H



Slow convergence is preferred

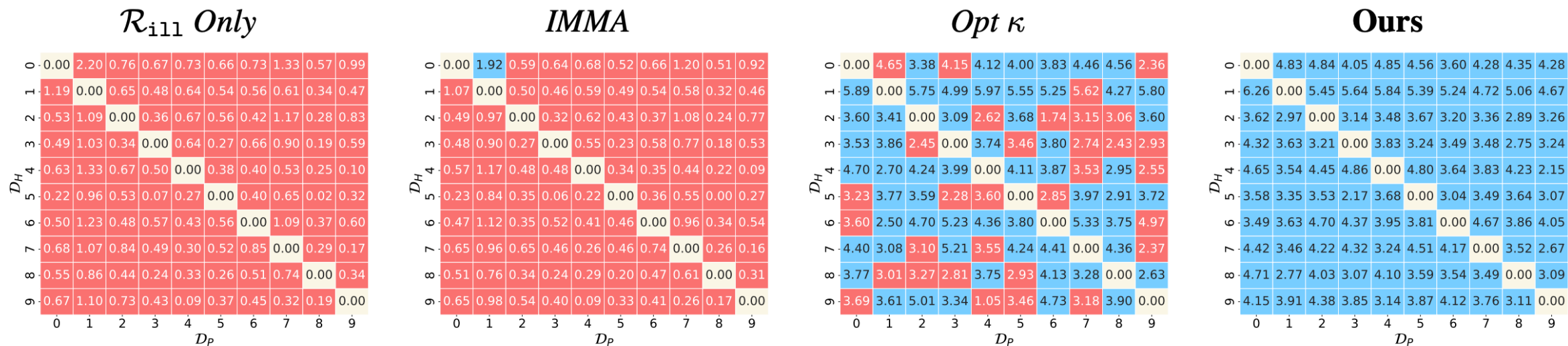
Norm ratio curve on \mathcal{D}_P



Fast convergence is preferred

CLASSIFICATION TASK (MNIST)

Binary classification on MNIST. We use different digit pairs for \mathcal{D}_P and \mathcal{D}_H .



$$\text{RIR} = \underbrace{\left(\frac{\kappa(\mathbf{H}_H(\theta_I))}{\kappa(\mathbf{H}_H(I))} \right)}_{(i)} \bigg/ \underbrace{\left(\frac{\kappa(\mathbf{H}_P(\theta_I))}{\kappa(\mathbf{H}_P(I))} \right)}_{(ii)}$$

Each element shows $\log(\text{RIR})$. Blue if successful immunization, red otherwise.

EXPERIMENTS (DEEP-NETS)

- Feature extractor: a pre-trained model (ResNet18 and ViT) with initialization θ_0
- The evaluation metric RIR is extended to comparing relative to θ_0

$$\text{RIR}_{\theta_0} \triangleq \underbrace{\left(\frac{\kappa(\tilde{\mathbf{H}}_{\text{H}}(\theta_{\text{I}}))}{\kappa(\tilde{\mathbf{H}}_{\text{H}}(\theta_0))} \right)}_{\text{(i)}} \bigg/ \underbrace{\left(\frac{\kappa(\tilde{\mathbf{H}}_{\text{P}}(\theta_{\text{I}}))}{\kappa(\tilde{\mathbf{H}}_{\text{P}}(\theta_0))} \right)}_{\text{(ii)}} \quad \text{where} \quad \tilde{\mathbf{H}}_{\text{H}}(\theta) = \nabla_{\mathbf{w}}^2 \mathcal{L}(\mathcal{D}_{\text{H}}, \mathbf{w}, \theta) = \tilde{\mathbf{X}}_{\text{H}}(\theta)^\top \tilde{\mathbf{X}}_{\text{H}}(\theta)$$

\mathcal{D}_{H}	Method	ResNet18		
		RIR $_{\theta_0}$ (i) \uparrow	RIR $_{\theta_0}$ (ii) \downarrow	RIR $_{\theta_0}$ \uparrow
	Init. θ_0	1.0	1.0	1.0
Cars	\mathcal{R}_{111} Only	1.878 \pm 0.034	1.786 \pm 0.025	1.057 \pm 0.026
	IMMA	0.866 \pm 0.002	0.889 \pm 0.001	0.974 \pm 0.002
	Opt κ	1.217 \pm 0.021	0.798 \pm 0.005	1.527 \pm 0.019
	Ours	2.386 \pm 0.442	0.699 \pm 0.062	3.467 \pm 0.358
Country211	\mathcal{R}_{111} Only	20.727 \pm 0.791	20.675 \pm 1.685	1.038 \pm 0.05
	IMMA	0.791 \pm 0.005	0.814 \pm 0.006	0.972 \pm 0.007
	Opt κ	1.538 \pm 0.155	1.053 \pm 0.091	1.472 \pm 0.043
	Ours	3.287 \pm 0.33	0.399 \pm 0.034	8.714 \pm 0.672

ViT		
RIR $_{\theta_0}$ (i) \uparrow	RIR $_{\theta_0}$ (ii) \downarrow	RIR $_{\theta_0}$ \uparrow
1.0	1.0	1.0
13.121 \pm 0.038	4.097 \pm 0.098	3.342 \pm 0.048
1.422 \pm 0.006	2.090 \pm 0.043	0.702 \pm 0.007
3.598 \pm 0.510	0.171 \pm 0.033	26.369 \pm 2.814
7.945 \pm 0.247	0.323 \pm 0.086	34.517 \pm 0.886
69.291 \pm 1.198	63.519 \pm 6.62	1.122 \pm 0.097
6.242 \pm 0.203	7.599 \pm 0.717	0.845 \pm 0.048
4.589 \pm 0.079	0.300 \pm 0.106	16.498 \pm 5.183
20.894 \pm 1.425	0.700 \pm 0.082	41.341 \pm 0.967

EXPERIMENTS (DEEP-NETS)

- Feature extractor: a pre-trained model (ResNet18 and ViT) with initialization θ_0
- The evaluation metric RIR is extended to comparing relative to θ_0
- We also report task performance, i.e., accuracy for image classification on \mathcal{D}_P .

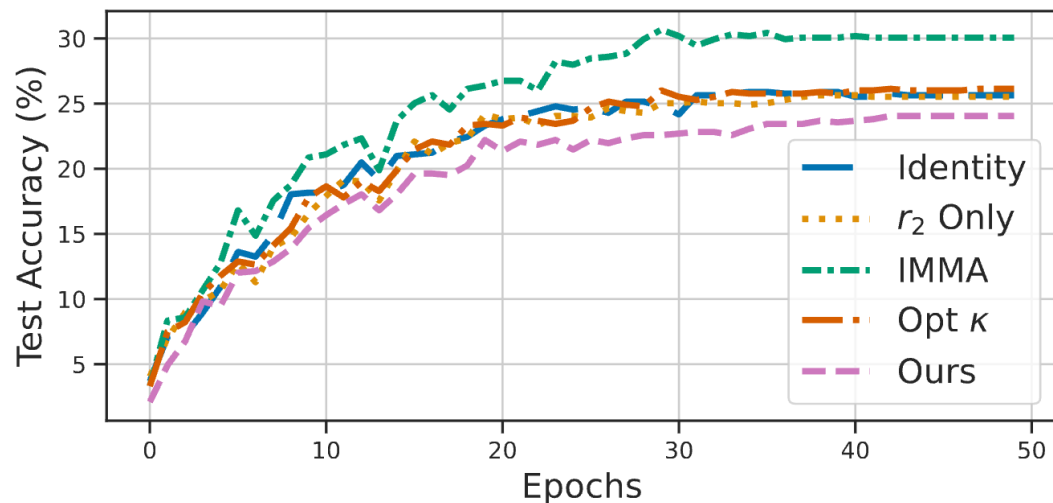
\mathcal{D}_H	Method	ResNet18			\mathcal{D}_P Test Acc. (%) \uparrow	ViT			\mathcal{D}_P Test Acc. (%) \uparrow
		RIR $_{\theta_0}$ (i) \uparrow	RIR $_{\theta_0}$ (ii) \downarrow	RIR $_{\theta_0}$ \uparrow		RIR $_{\theta_0}$ (i) \uparrow	RIR $_{\theta_0}$ (ii) \downarrow	RIR $_{\theta_0}$ \uparrow	
	Init. θ_0	1.0	1.0	1.0	68.24	1.0	1.0	1.0	81.78
Cars	\mathcal{R}_{111} Only	1.878 \pm 0.034	1.786 \pm 0.025	1.057 \pm 0.026	63.84 \pm 0.292	13.121 \pm 0.038	4.097 \pm 0.098	3.342 \pm 0.048	82.21 \pm 0.035
	IMMA	0.866 \pm 0.002	0.889 \pm 0.001	0.974 \pm 0.002	63.57 \pm 0.234	1.422 \pm 0.006	2.090 \pm 0.043	0.702 \pm 0.007	81.89 \pm 0.010
	Opt κ	1.217 \pm 0.021	0.798 \pm 0.005	1.527 \pm 0.019	63.65 \pm 0.148	3.598 \pm 0.510	0.171 \pm 0.033	26.369 \pm 2.814	82.51 \pm 0.085
	Ours	2.386 \pm 0.442	0.699 \pm 0.062	3.467 \pm 0.358	62.36 \pm 0.173	7.945 \pm 0.247	0.323 \pm 0.086	34.517 \pm 0.886	82.79 \pm 0.200
Country211	\mathcal{R}_{111} Only	20.727 \pm 0.791	20.675 \pm 1.685	1.038 \pm 0.05	62.17 \pm 1.599	69.291 \pm 1.198	63.519 \pm 6.62	1.122 \pm 0.097	80.73 \pm 0.129
	IMMA	0.791 \pm 0.005	0.814 \pm 0.006	0.972 \pm 0.007	67.03 \pm 0.146	6.242 \pm 0.203	7.599 \pm 0.717	0.845 \pm 0.048	82.47 \pm 0.036
	Opt κ	1.538 \pm 0.155	1.053 \pm 0.091	1.472 \pm 0.043	66.81 \pm 0.115	4.589 \pm 0.079	0.300 \pm 0.106	16.498 \pm 5.183	82.79 \pm 0.023
	Ours	3.287 \pm 0.33	0.399 \pm 0.034	8.714 \pm 0.672	65.01 \pm 0.143	20.894 \pm 1.425	0.700 \pm 0.082	41.341 \pm 0.967	83.17 \pm 0.075

CLASSIFICATION TASK

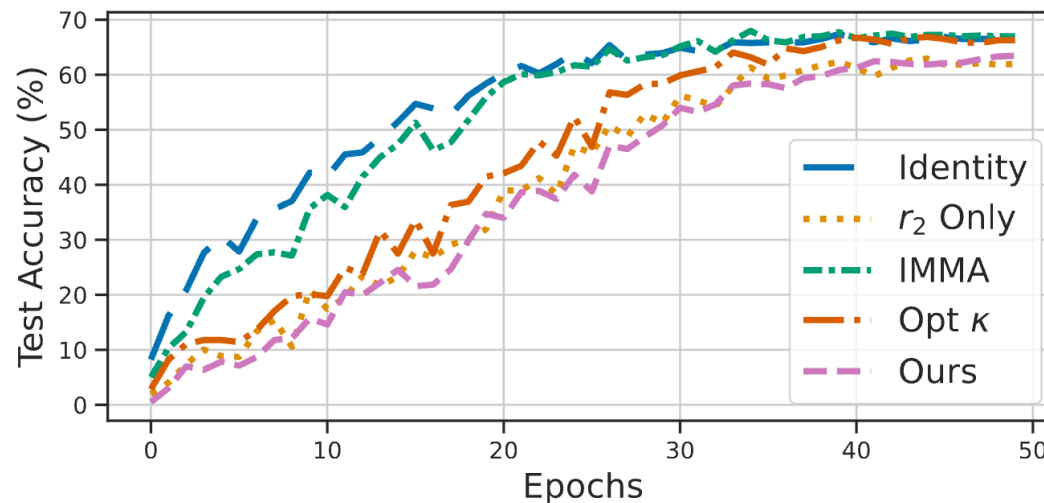
Features pre-trained on ImageNet transferring to \mathcal{D}_H

\mathcal{D}_H	Method	ResNet18				ViT			
		RIR_{θ_0} (i) \uparrow	RIR_{θ_0} (ii) \downarrow	RIR_{θ_0} \uparrow	\mathcal{D}_P Test Acc. (%) \uparrow	RIR_{θ_0} (i) \uparrow	RIR_{θ_0} (ii) \downarrow	RIR_{θ_0} \uparrow	\mathcal{D}_P Test Acc. (%) \uparrow
Cars	Init. θ_0	1.0	1.0	1.0	68.24	1.0	1.0	1.0	81.78
	\mathcal{R}_{i11} Only	1.878 ± 0.034	1.786 ± 0.025	1.057 ± 0.026	63.84 ± 0.292	13.121 ± 0.038	4.097 ± 0.098	3.342 ± 0.048	82.21 ± 0.035
	IMMA	0.866 ± 0.002	0.889 ± 0.001	0.974 ± 0.002	63.57 ± 0.234	1.422 ± 0.006	2.090 ± 0.043	0.702 ± 0.007	81.89 ± 0.010
	Opt κ	1.217 ± 0.021	0.798 ± 0.005	1.527 ± 0.019	63.65 ± 0.148	3.598 ± 0.510	0.171 ± 0.033	26.369 ± 2.814	82.51 ± 0.085
	Ours	2.386 ± 0.442	0.699 ± 0.062	3.467 ± 0.358	62.36 ± 0.173	7.945 ± 0.247	0.323 ± 0.086	34.517 ± 0.886	82.79 ± 0.200

Fine-tuning accuracy with ResNet-18



Fine-tuning accuracy with ViT



TAKE AWAYS

- A **condition number based** framework for model immunization
- Two differentiable regularizers and a gradient-based optimization algorithm.
- A first step towards principled understanding of model immunization.
- Project page / code:
 - amberyzheng.com/immu_cond_num



Closely related works:

- Rosati, Domenic, et al. "Representation noising: A defence mechanism against harmful finetuning." Proc. NeurIPS, 2024
- Zheng, Amber Yijia, et al. "Learning to obstruct few-shot image classification over restricted classes." Proc. ECCV, 2024
- Nenov, Rossen, et al. "(Almost) Smooth Sailing: Towards Numerical Stability of Neural Networks Through Differentiable Regularization of the Condition Number." ICML Workshop, 2024

