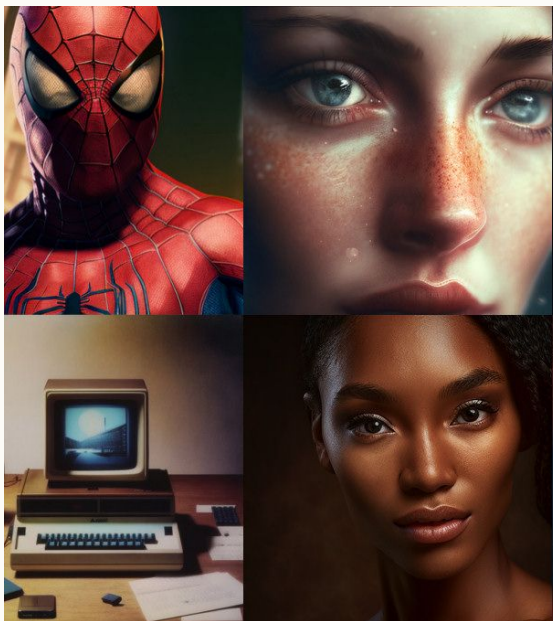


Finding DoRI : Discovery of Retained Images in Diffusion Models

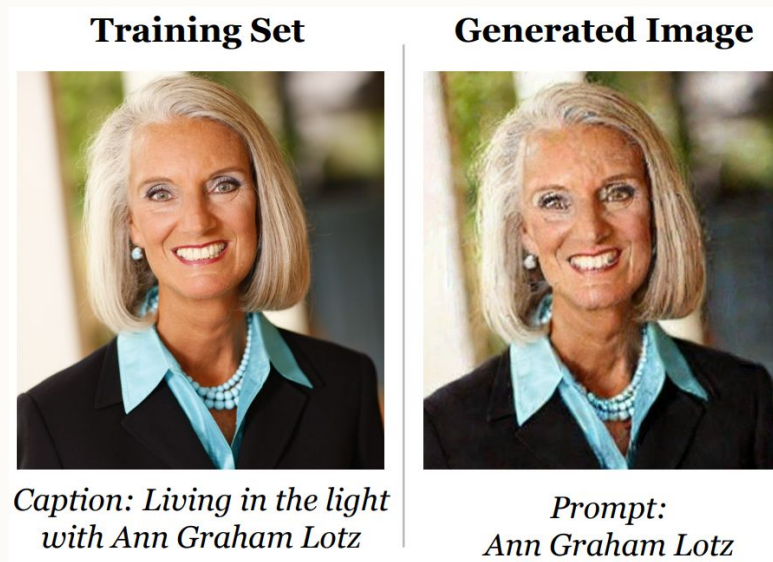
Antoni Kowalczyk*, Dominik Hintersdorf*, Lukas Struppek*,
Kristian Kersting, Adam Dziedzic, Franziska Boenisch

Diffusion Models (DMs)

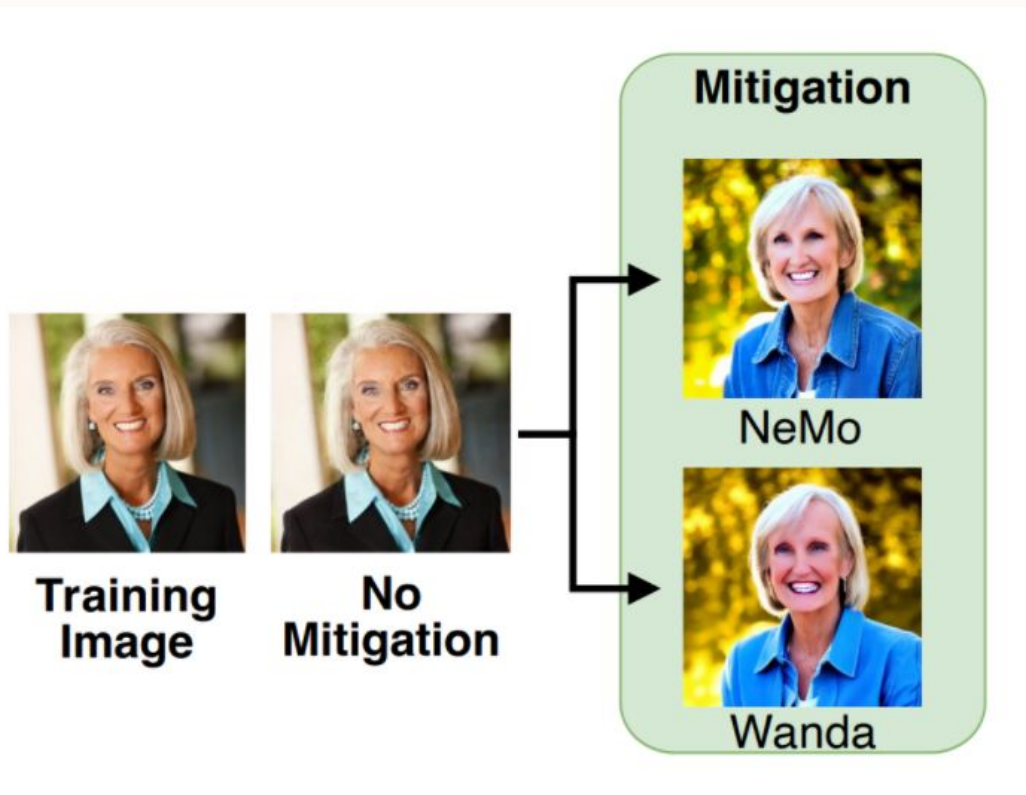
High-Quality Images



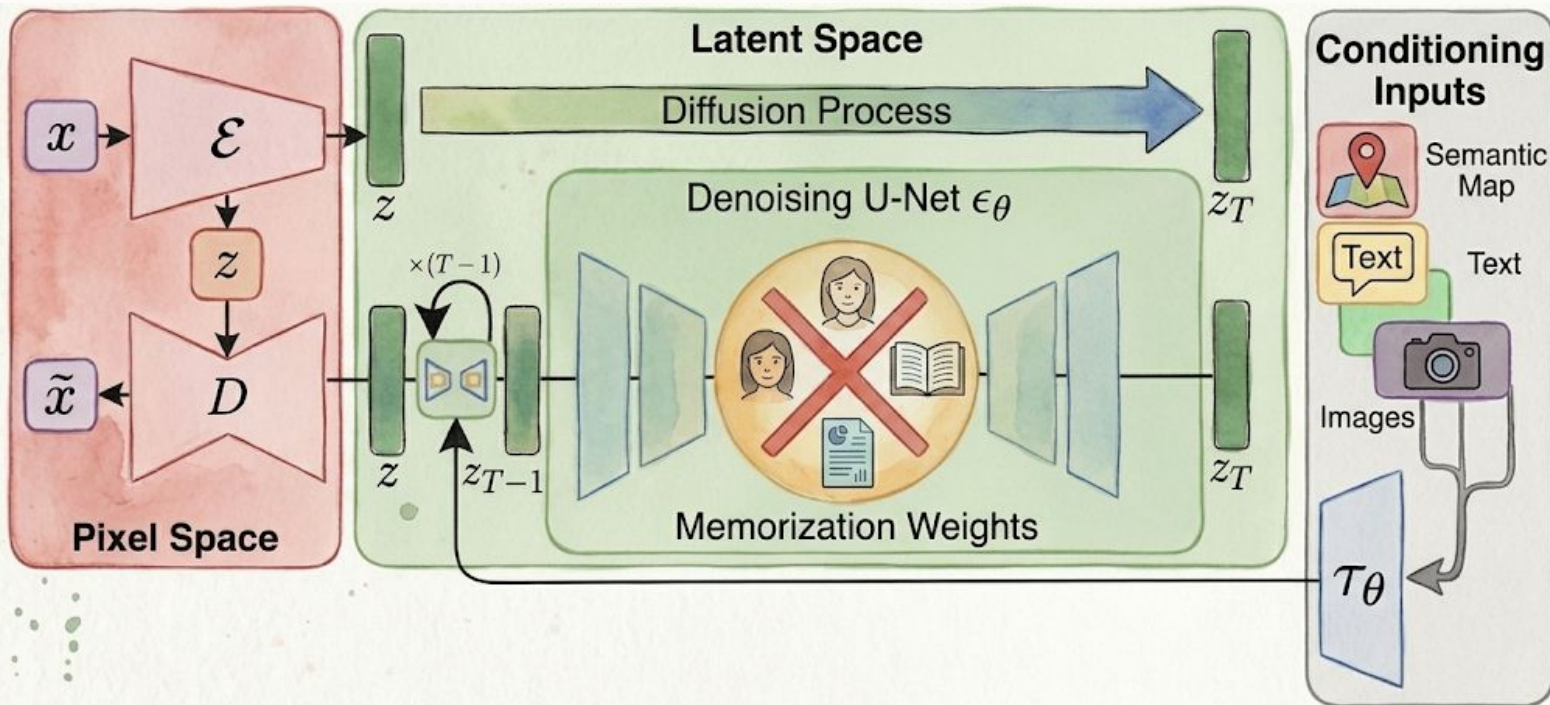
Memorize Training Data



Mitigation Methods



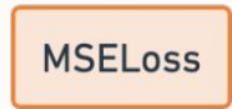
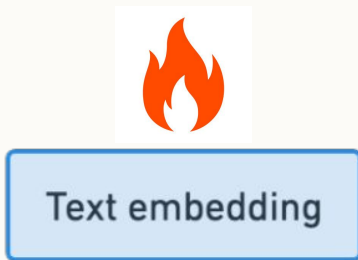
Pruning



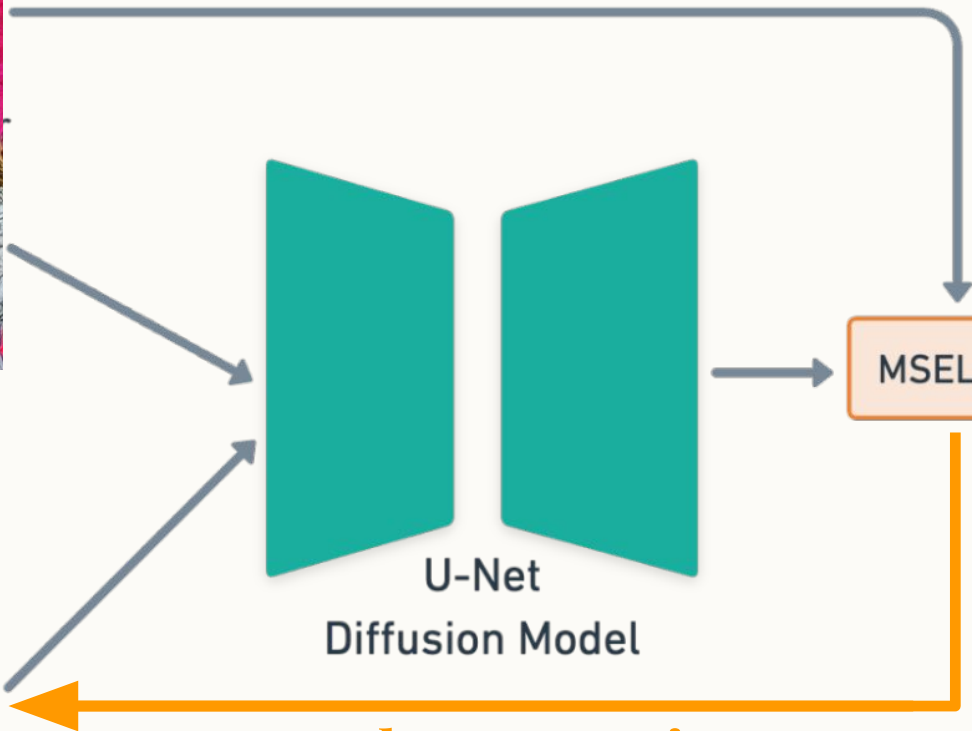
This Work

1. DoRI 
2. Locality Refutation
3. Robust Mitigation

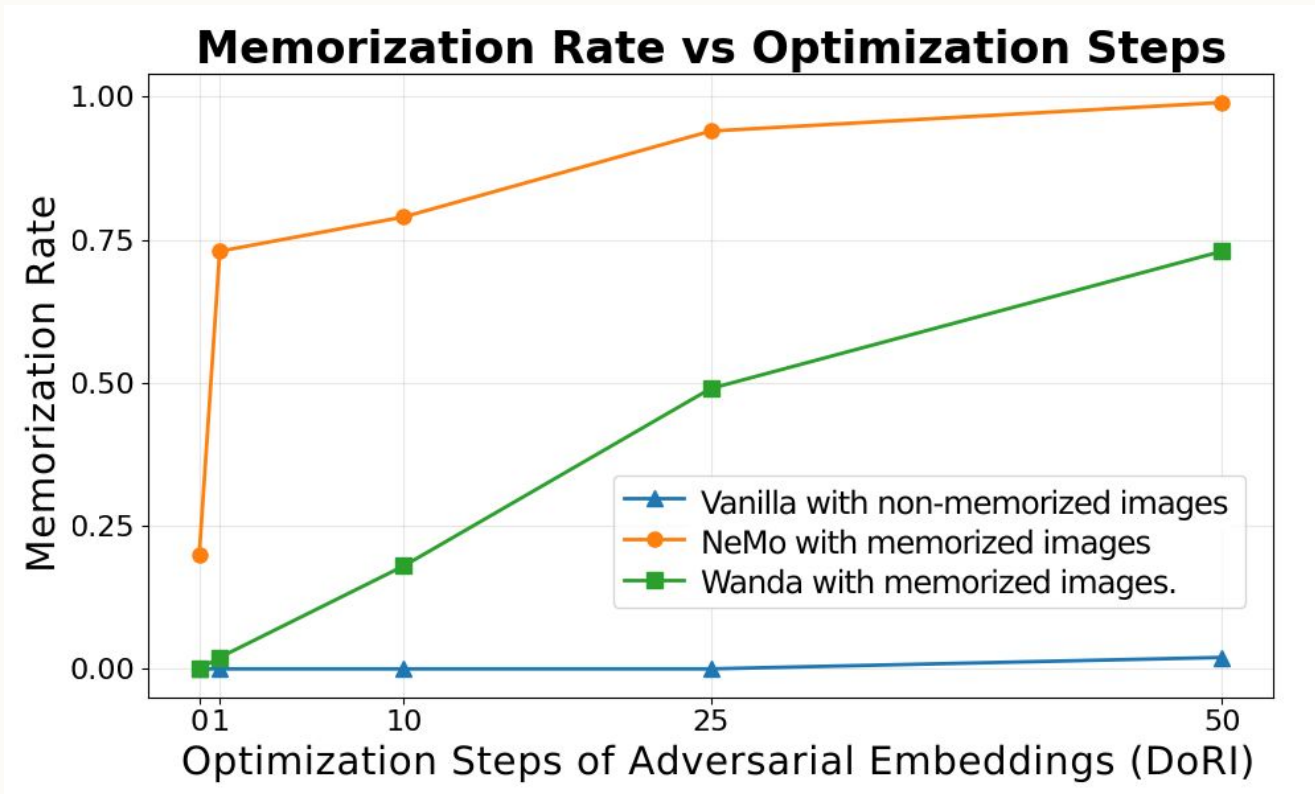
Discovery of Retained Images (DoRI)



Backpropagation



Evading Pruning with DoRI



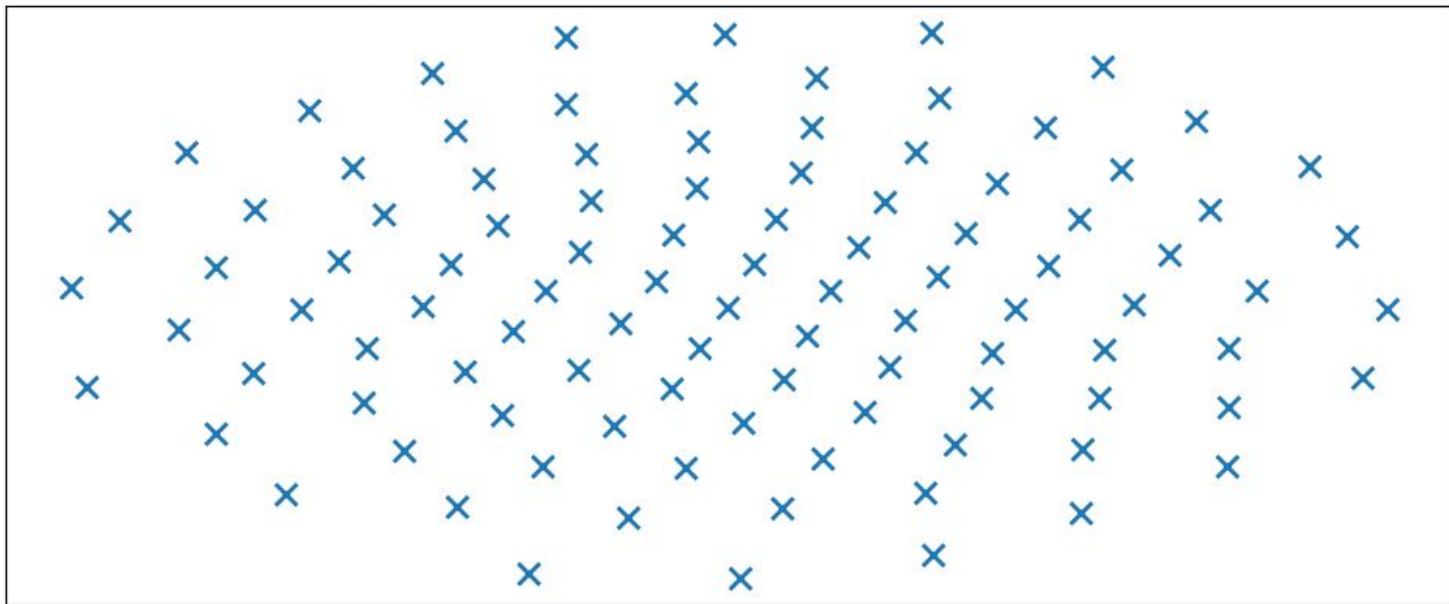
Locality Assumption

1. Input Space
2. Activation Space
3. Weights

Random Initialization

Text Embedding Type

× Initial Embedding $\mathbf{y}_{\text{adv}}^{(0)} \sim \mathcal{N}(0, \mathbf{I})$



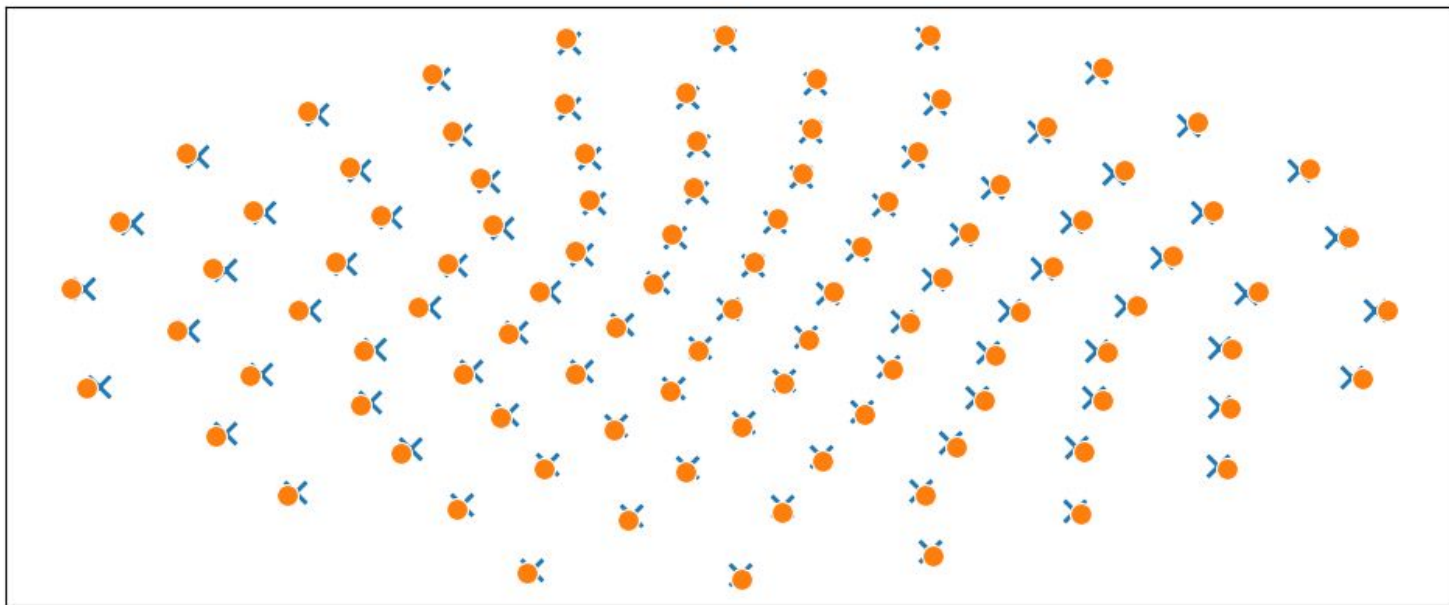
The Same Output Image After DoRI



Text Embedding Type

× Initial Embedding $\mathbf{y}_{\text{adv}}^{(0)} \sim \mathcal{N}(0, \mathbf{I})$

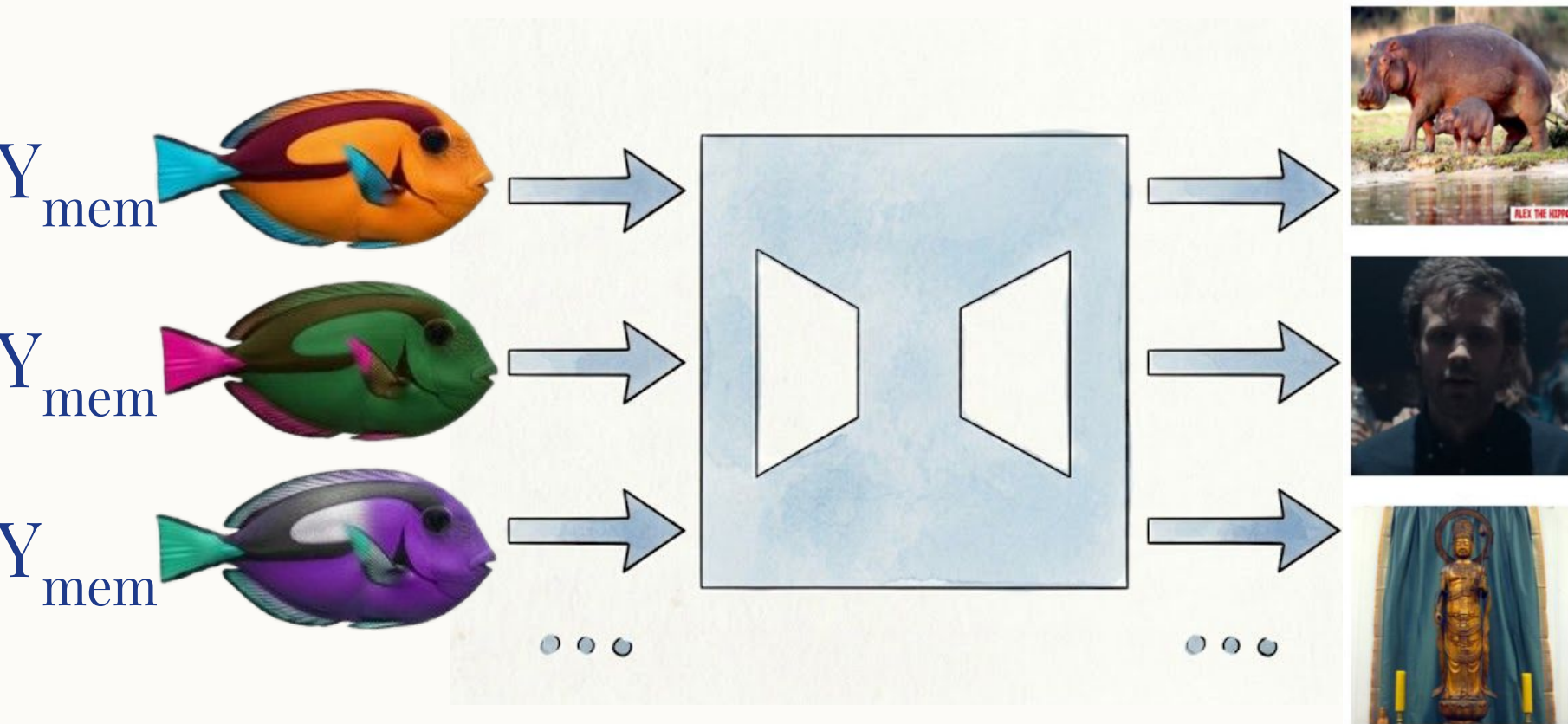
● Adversarial Embedding \mathbf{y}_{adv}



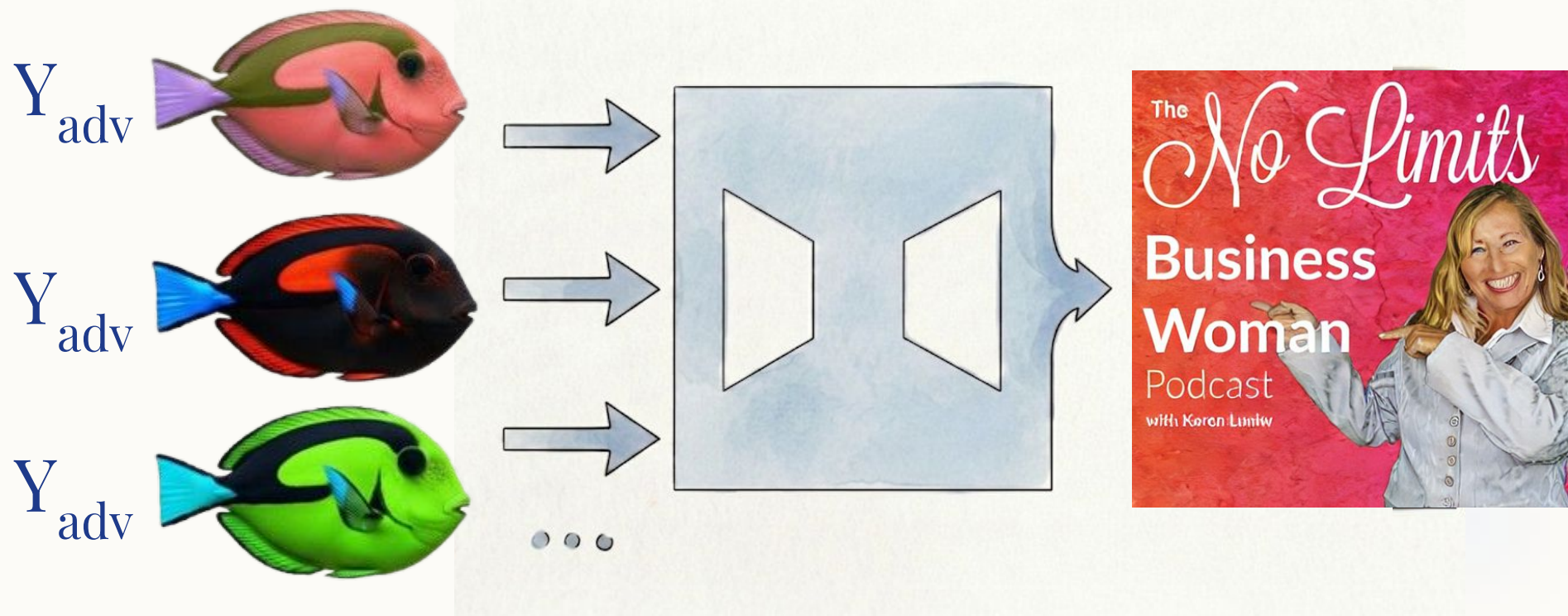
Locality Assumption

1. ~~Input Space~~ **Refuted**
2. Activation Space
3. Weights

Different Input \Rightarrow Different Output

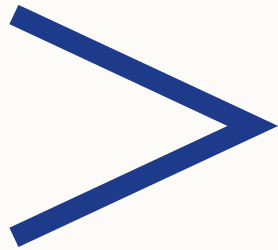


Different Input \Rightarrow Same Output



Expected if Locality Holds:

Discrepancy for Y_{mem}



Discrepancy for Y_{adv}



Actual:

Discrepancy for Y_{mem}



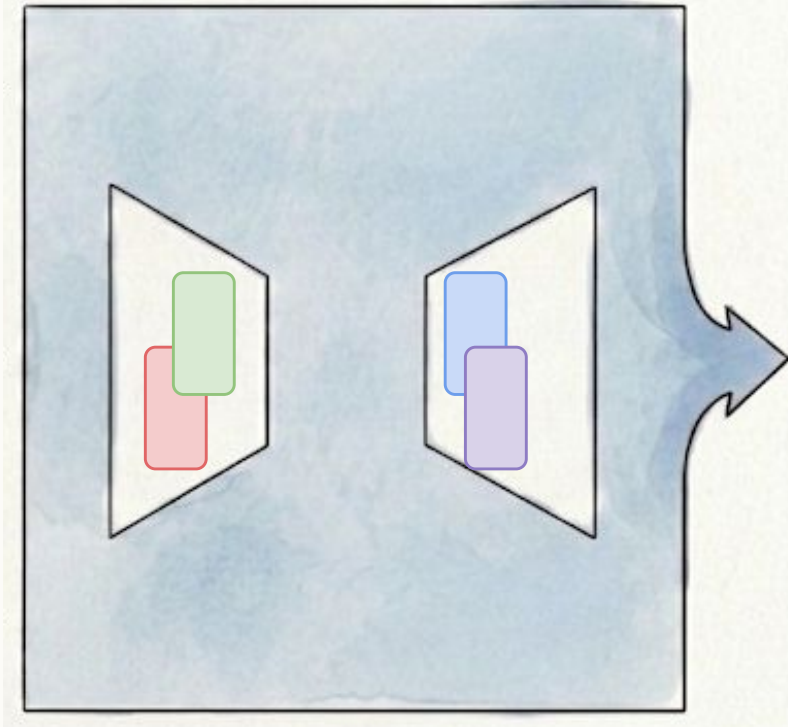
Discrepancy for Y_{adv}



Locality Assumption

1. ~~Input Space~~ **Refuted**
2. ~~Activation Space~~ **Refuted**
3. Weights

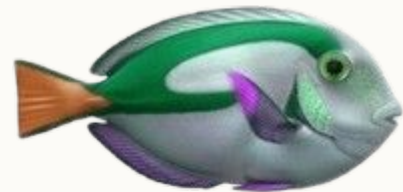
Same Output \Rightarrow Different Weights



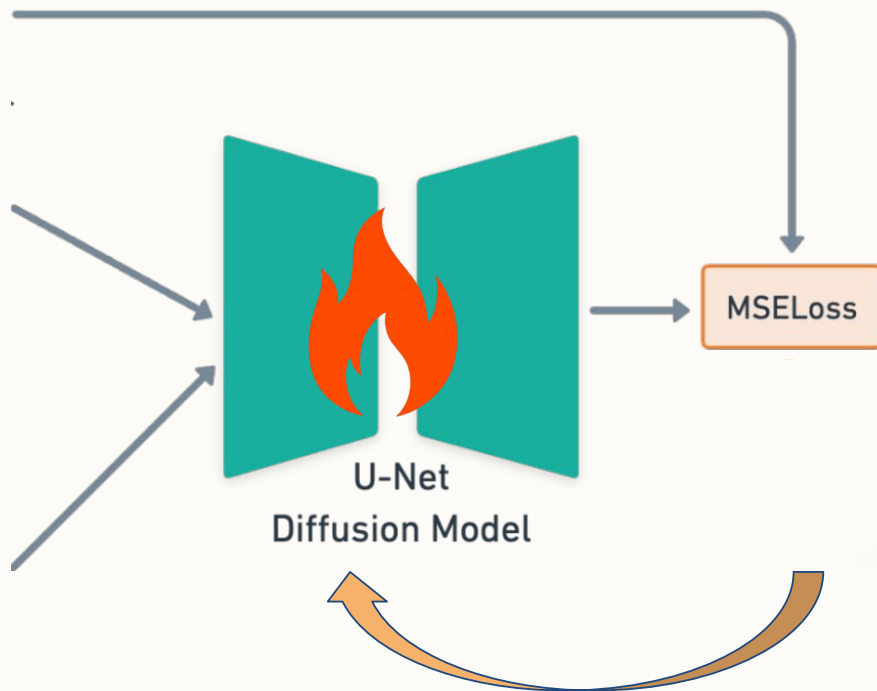
~~Locality Assumption~~ **Refuted**

1. ~~Input Space~~ **Refuted**
2. ~~Activation Space~~ **Refuted**
3. ~~Weights~~ **Refuted**

Our Mitigation



Our Mitigation



Backpropagation

Our Mitigation is Successful

	Prompt	
No Mitigation	1.00 MR	1.00 MR
Our Mitigation	0.00 MR	0.02 MR

This Work

1. DoRI 
2. Locality Refutation
3. Robust Mitigation