



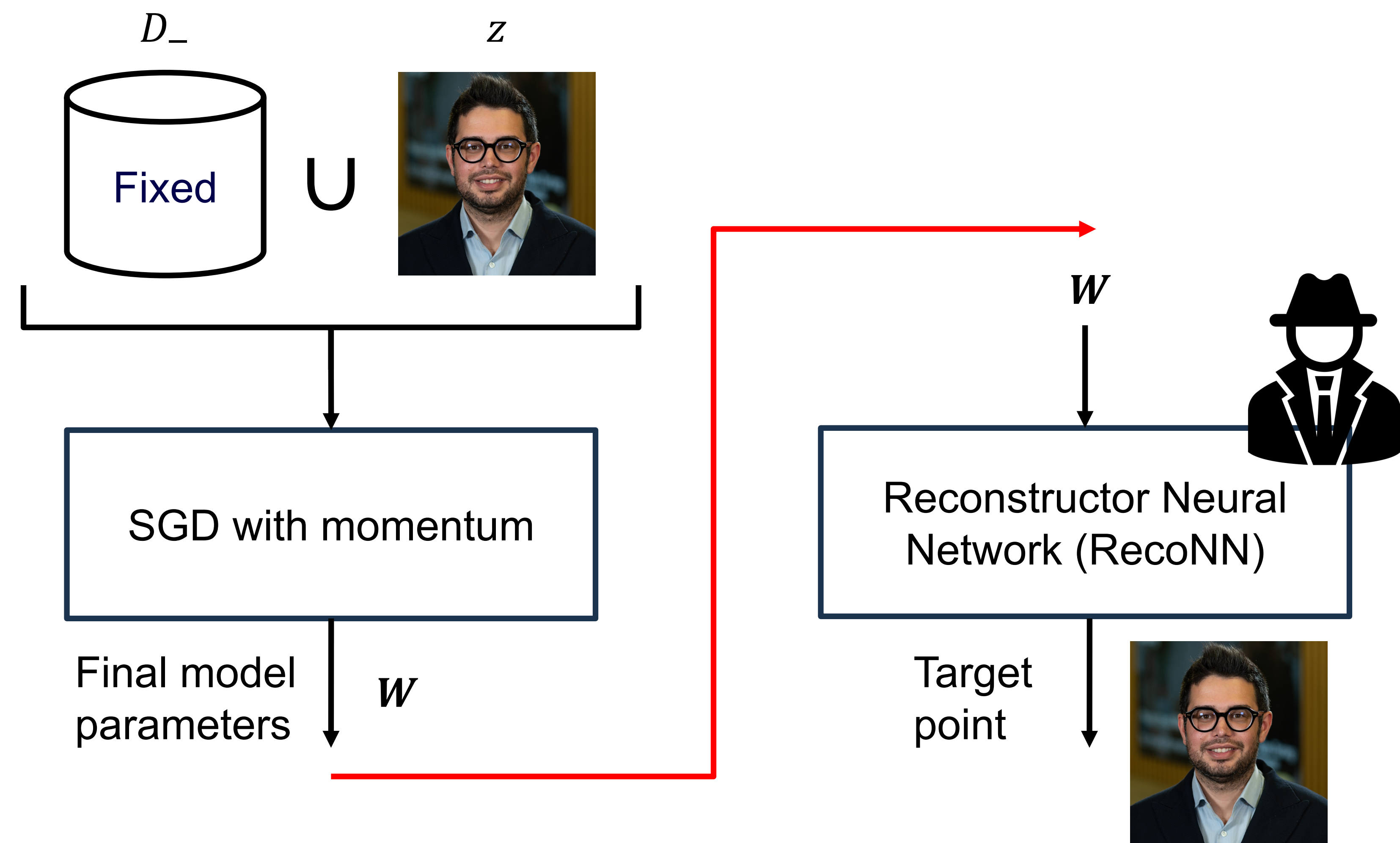
# Mnemonist: Locating Model Parameters that Memorize Training Examples

Ali Shahin Shamsabadi, Jamie Hayes, Borja Balle, Adrian Weller



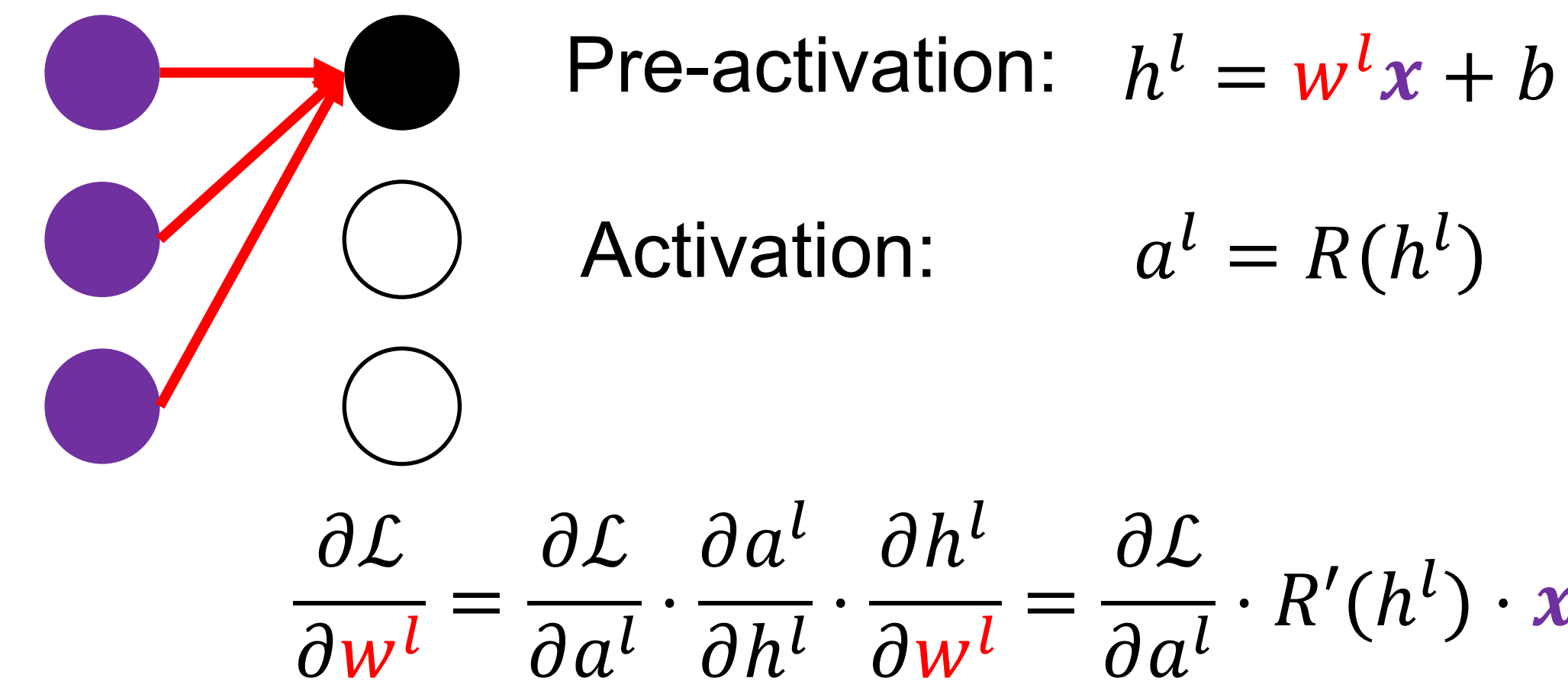
39th Conference on Uncertainty in Artificial Intelligence (UAI), Pittsburgh, PA, USA, August 2023

## Reconstructing training data with an informed adversary



## Why is training data reconstruction from ReLU activated models hard?

Characterising the fundamental property of the existence of redundant parameters in models with ReLU activations:

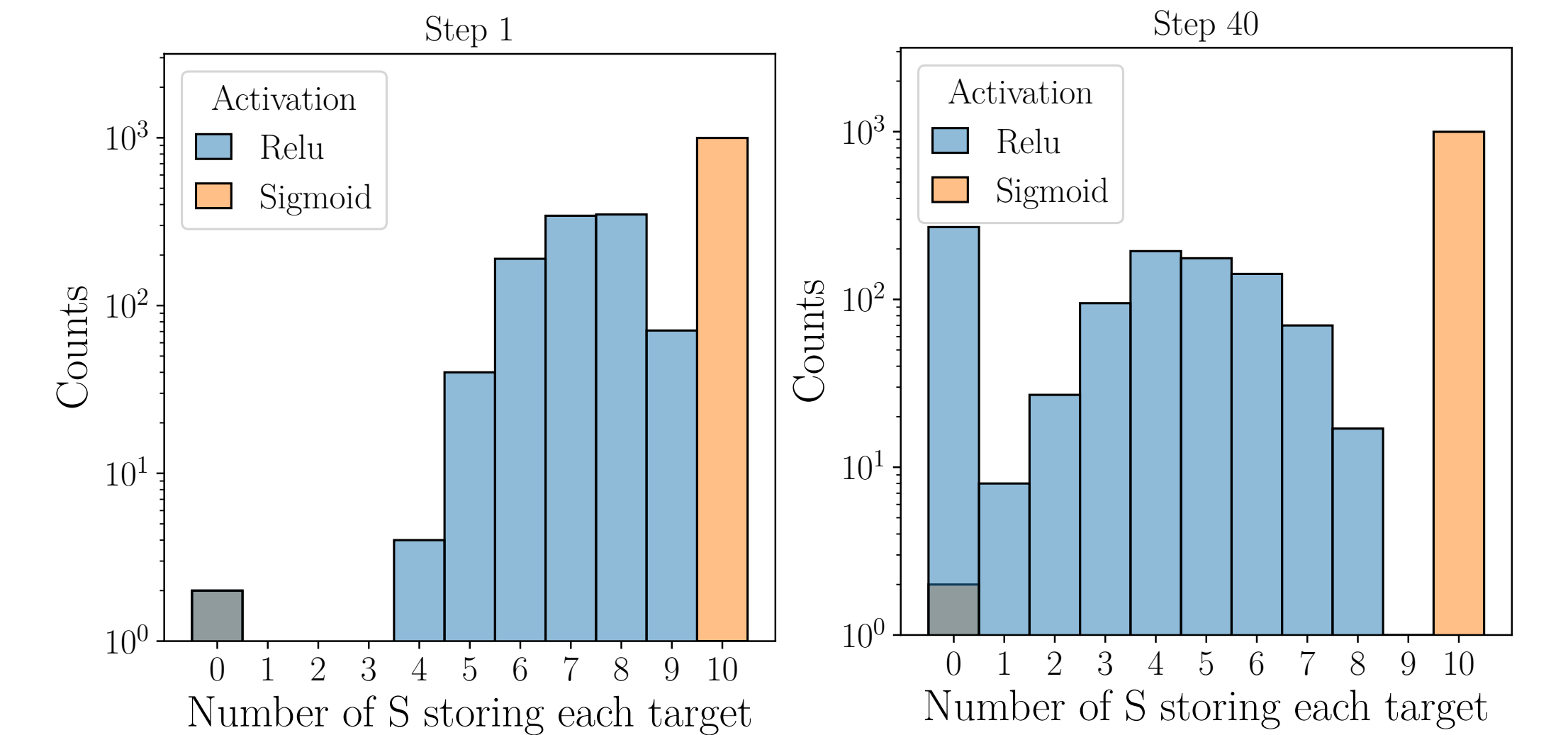


$$\frac{\partial \mathcal{L}}{\partial w^l} = \begin{cases} 0 & \text{if } a^l = 0 \\ \text{scale} \cdot x & \text{otherwise} \end{cases}$$

No information about the input is stored in incoming parameters to non-activate neurons

Incoming parameters to active neurons store a copy of input data

Demonstrating the existence of redundant parameters in practice by studying the training dynamics:



Histogram of the summation of

$$\text{Binary scale} = \begin{cases} 0 & \text{if scale} = 0 \\ 1 & \text{Otherwise} \end{cases}$$

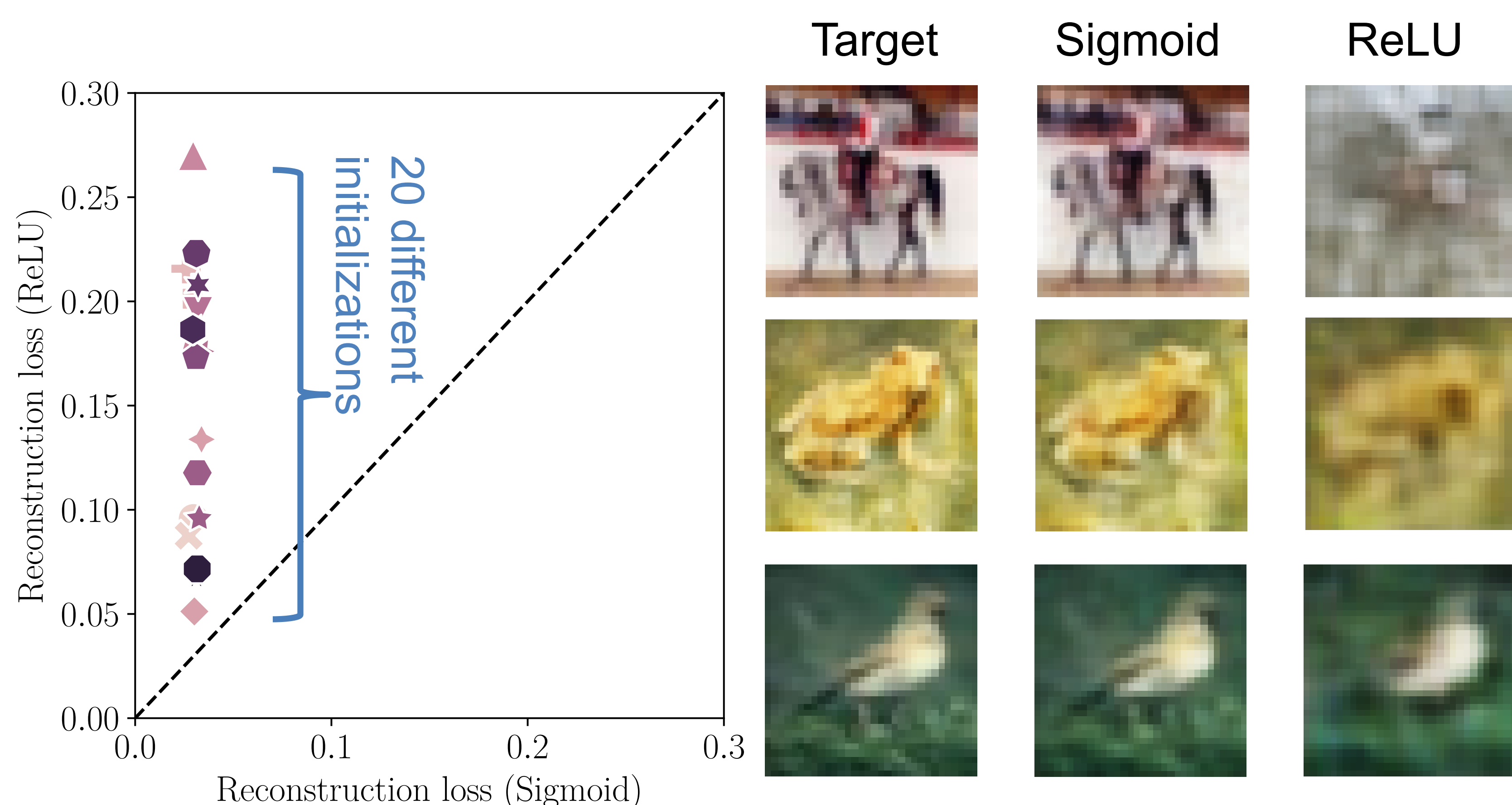
of all rows per target point. [1k target points]

The attack consists of three steps:

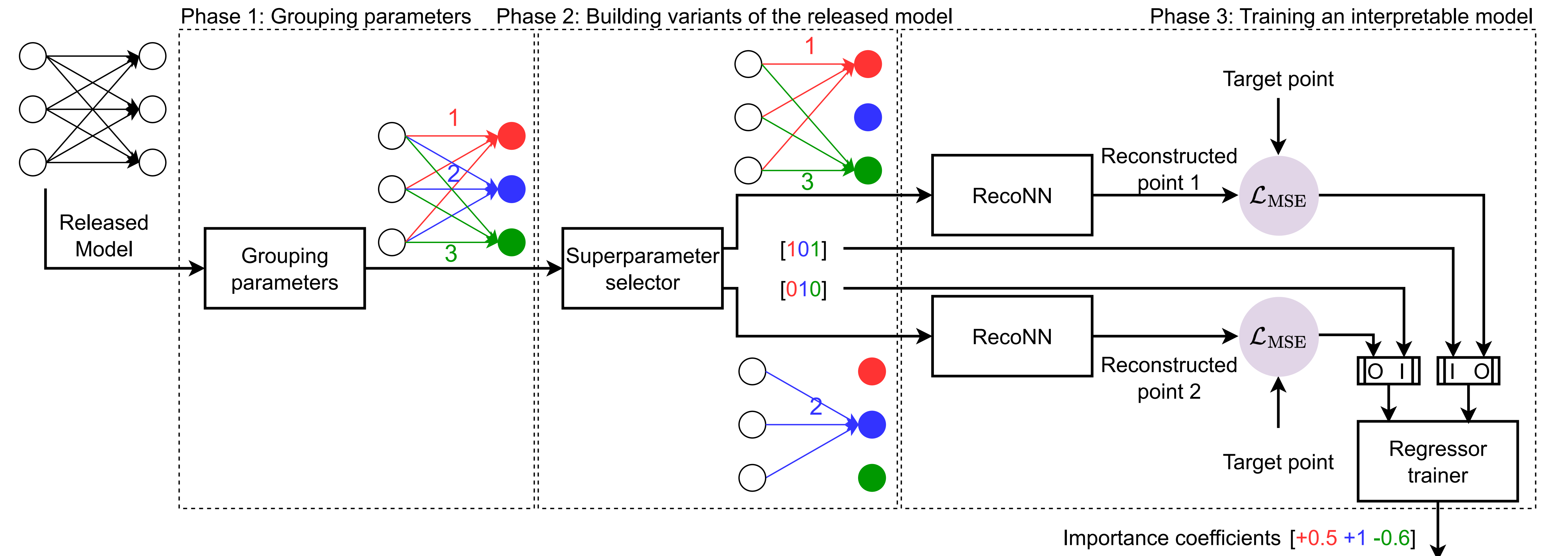
1. Training **shadow models** to collect information about the impact of training examples on model parameters
  - The adversary knows  $(D_-, W, W_{init}, S)$  + a public dataset  $\bar{Z} = \{\bar{z}_i\}_{i=1}^K$
  - The adversary trains each shadow model,  $\bar{W}_i$ , on  $D_- \cup \{\bar{z}_i\}$
2. Training a **RecoNN** to output training examples from model parameters
  - Loss function:  $\text{MSE}(\text{RecoNN}(\bar{W}_i), \bar{z}_i) + \text{MAE}(\text{RecoNN}(\bar{W}_i), \bar{z}_i)$
3. Producing a **candidate** reconstruction for the target point

## Reconstruction quality depends heavily on the activation

ReLU activations lead to higher reconstruction loss compared to Sigmoid.

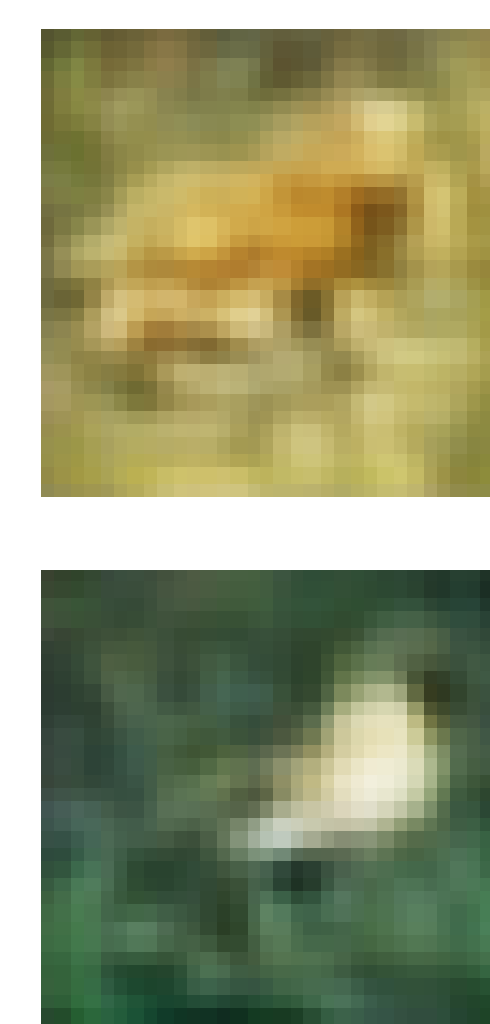


## Mnemonist: finding parameters that store target examples



## Mnemonist helps RecoNN to improve the quality of reconstructions

Mnemonist-RecoNN trains RecoNN only on informative superparameters.



Approach	run1	run2	run3	run4
RecoNN	.2040	.2705	.0908	.1315
Mnemonist-RecoNN	.1738	.2385	.0730	.1158

### Algorithm 1: Mnemonist-RecoNN

**Input:** Fixed set  $D_-$ ,  $K$  public target examples  $\{\bar{z}_k\}_{k=1}^K$ , Shadow model training Algorithm  $A(\cdot)$ , RecoNN training algorithm  $B(\cdot)$ .

**Output:** Mnemonist-guided RecoNN.

- 1: **for all**  $k \in \text{range}(K)$  **do**
- 2:  $\bar{W}_k \leftarrow A(D_- \cup \{\bar{z}_k\})$   $\triangleright$  Train  $K$  shadow models
- 3:  $\phi \leftarrow B(\{\bar{W}_k, \bar{z}_k\}_{k=1}^K)$   $\triangleright$  Train RecoNN
- 4:  $I \leftarrow \text{Mnemonist}(\phi)$   $\triangleright$  Apply Mnemonist to identify the importance of superparameters
- 5: **for all**  $k \in \text{range}(K)$  **do**
- 6:  $\bar{W}_k \leftarrow \text{Selector}(\bar{W}_k, I)$   $\triangleright$  Select only important superparameters
- 7:  $\tilde{\phi} \leftarrow B(\{\bar{W}_k, \bar{z}_k\}_{k=1}^K)$   $\triangleright$  Train RecoNN on the selected important superparameters
- 8: **return**  $\tilde{\phi}$   $\triangleright$  Mnemonist-guided RecoNN