



## Background

Among the numerous advancements in the field of Deep Learning, only a handful address model resistance to privacy attacks. Current techniques predominantly rely on **Differential Privacy**. While they offer solid theoretical guarantees, they also come with significant drawbacks, notably in performance and complexity. Interestingly, also **regularization techniques** exhibit side effects that can occasionally benefit privacy, even though they are not explicitly designed for it.

## Discriminative Adversarial Privacy (DAP)

Consider a **classifier**  $C_{DAP}$  and a pre-trained **Membership Inference Attack** (MIA) model  $D_{MIA}$ . Let  $x, y \sim p(x, y)$  be a sample's features and labels extracted from the data domain, and  $t$  the current iteration index. Then, the learning procedure of DAP is

$$\min_{C_{DAP}} \max_{D_{MIA}} \mathbb{E}_{x \sim p(x)} [\log C_{DAP}(x, t)] + \beta \mathbb{E}_{x, y \sim p(x, y)} [\log(1 - D_{MIA}(C_{DAP}(x, t), y))]$$

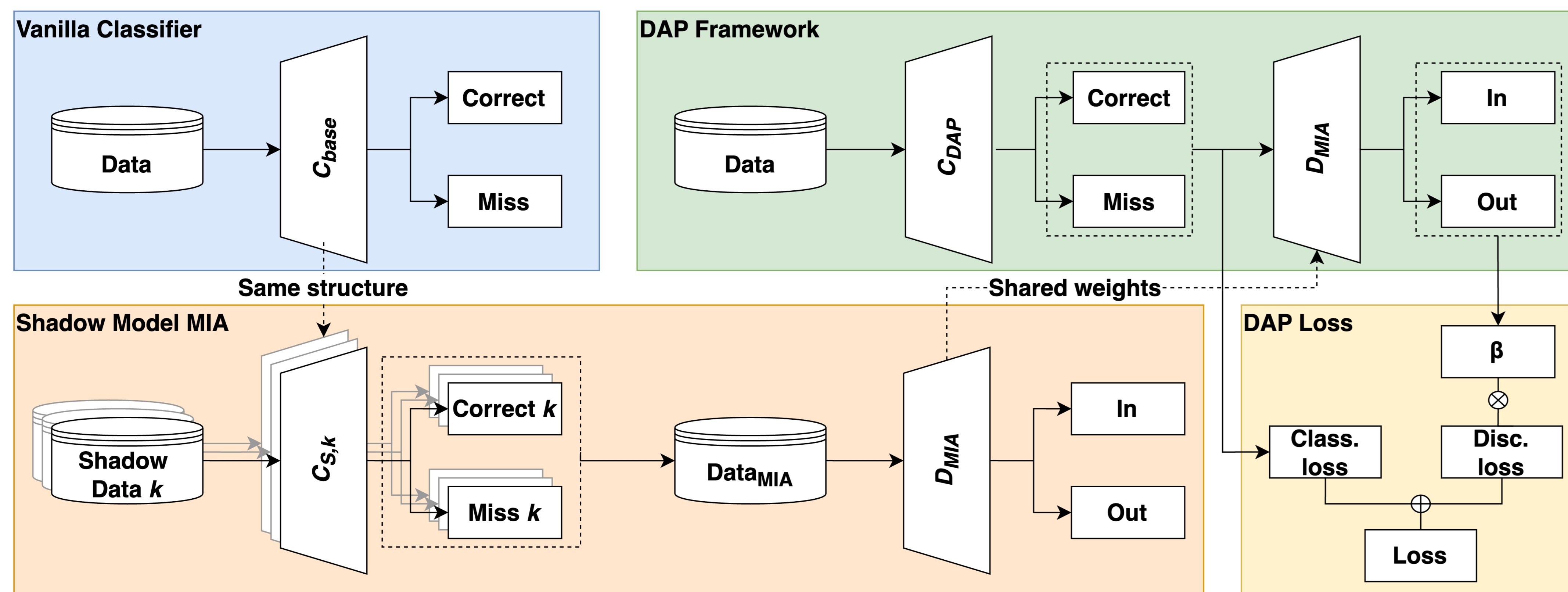
where  $\beta$  is a dynamic loss balancing parameter.  $\beta$  is a function of a pre-defined ratio factor  $r$ .  $\beta$  is defined as

$$\beta(C_{DAP}, D_{MIA}, t, r) = \begin{cases} r \cdot \frac{\mathbb{E}[\log C_{DAP}(x, t-1)]_v}{\mathbb{E}[\log(1 - D_{MIA}(C_{DAP}(x, t-1), y))]_v} & \text{if } t > 0 \\ 1 & \text{otherwise} \end{cases}$$

Given the accuracy ACC from  $C_{DAP}$ , the area under the curve  $AUC_{MIA}$  from  $D_{MIA}$  and the privacy parameter  $\lambda$ , the training step with the optimal model parameters is found by maximising the AOP metric:

$$AOP(\lambda) = \frac{ACC}{(2 \max(AUC_{MIA}, 0.5))^\lambda}$$

## Architecture



## Contributions

- A novel learning technique called **Discriminative Adversarial Privacy** (DAP), that combines adversarial learning and membership inference attacks to ensure an optimal balance between model performance, complexity, and privacy.
- A novel **loss function** that is tailored to simultaneously minimise the model prediction error while maximising the attacker's error.
- A novel metric called **Accuracy Over Privacy** (AOP), to capture the performance-privacy trade-off effectively.

## Results

Dataset	Baseline	Accuracy (ACC)						$DAP_t$	$DAP_v$
		Reg	$\epsilon = 0.5$	$\epsilon = 1$	$\epsilon = 2$	$\epsilon = 4$			
<b>Cifar-10</b>	0.784	<b>0.811</b>	0.313	0.374	0.417	0.418	0.624	0.613	
<b>Cifar-100</b>	0.481	<b>0.532</b>	0.039	0.083	0.090	0.072	0.315	0.276	
<b>FMNIST</b>	0.932	<b>0.926</b>	0.605	0.701	0.736	0.774	0.866	0.871	
<b>EuroSAT</b>	0.958	<b>0.950</b>	0.308	0.588	0.681	0.646	0.900	0.893	
<b>TinyImagenet</b>	0.365	<b>0.378</b>	0.031	0.032	0.032	0.025	0.260	0.217	
<b>OxfordFlowers</b>	0.566	<b>0.659</b>	0.031	0.051	0.087	0.139	0.290	0.257	
<b>STL-10</b>	0.655	<b>0.650</b>	0.084	0.142	0.250	0.289	0.480	0.384	
<b>Cinic-10</b>	0.673	<b>0.709</b>	0.280	0.341	0.391	0.405	0.577	0.586	
<b>Average</b>	0.677	<b>0.702</b>	0.211	0.289	0.336	0.346	0.539	0.512	

## Membership Inference Attack Area-Under-the-Curve ( $AUC_{MIA}$ )

Dataset	Baseline	Reg	$AUC_{MIA}$				$DAP_t$	$DAP_v$
			$\epsilon = 0.5$	$\epsilon = 1$	$\epsilon = 2$	$\epsilon = 4$		
<b>Cifar-10</b>	0.648	0.631	0.505	0.526	0.519	<b>0.503</b>	0.507	0.505
<b>Cifar-100</b>	0.603	0.621	<b>0.501</b>	0.515	0.507	0.506	0.516	0.506
<b>FMNIST</b>	0.552	0.562	<b>0.502</b>	<b>0.502</b>	0.504	0.505	0.507	0.506
<b>EuroSAT</b>	0.544	0.528	0.505	0.502	<b>0.500</b>	0.502	0.501	0.501
<b>TinyImagenet</b>	0.603	0.592	0.514	<b>0.501</b>	0.521	0.504	0.516	0.509
<b>OxfordFlowers</b>	0.761	0.765	0.543	0.537	0.526	0.532	0.538	<b>0.521</b>
<b>STL-10</b>	0.604	0.563	0.502	0.524	0.505	<b>0.501</b>	0.508	0.506
<b>Cinic-10</b>	0.572	0.614	<b>0.501</b>	0.514	0.511	0.507	0.513	0.507
<b>Average</b>	0.611	0.609	0.509	0.514	0.511	<b>0.507</b>	0.513	0.508

## Accuracy Over Privacy (AOP)

$\lambda$	Baseline	Reg	$\epsilon = 0.5$	$\epsilon = 1$	$\epsilon = 2$	$\epsilon = 4$	$DAP_t$	$DAP_v$
<b>1</b>	0.567	<b>0.587</b>	0.210	0.284	0.331	0.343	0.529	0.506
<b>2</b>	0.479	0.497	0.208	0.280	0.327	0.340	<b>0.519</b>	0.500
<b>5</b>	0.301	0.316	0.204	0.267	0.316	0.330	<b>0.492</b>	0.483
<b>10</b>	0.154	0.168	0.197	0.249	0.299	0.317	0.451	<b>0.456</b>
<b>20</b>	0.049	0.064	0.184	0.222	0.271	0.292	0.386	<b>0.409</b>
<b>50</b>	0.003	0.008	0.155	0.175	0.217	0.238	0.262	<b>0.305</b>