

### Human-guided Models Come at a Cost

**Ongoing challenge**  
Human-guided training has proven to boost generalization capabilities of deep learning-based models. However, acquiring human annotations is costly.

**Our solution**  
We utilize models first taught by human annotations ("Teacher models") to then train "Student models" through their saliency maps using the CYBORG loss.

**CYBORG loss function**

$$\mathcal{L} = \frac{1}{K} \sum_{k=1}^K \sum_{c=1}^C \mathbf{1}_{y_k \in C_c} \left[ \underbrace{(1 - \alpha) \|s_k^{(\text{teacher})} - s_k^{(\text{model})}\|_2}_{\text{teacher saliency loss component}} - \underbrace{\alpha \log p_{\text{model}}(y_k \in C_c)}_{\text{classification loss component}} \right]$$

$\|\cdot\| = L_2$  norm  
 $y_k$  = class label for the  $k$ -th sample  
 $\mathbf{1}[\cdot]$  = 1 when  $y_k \in C_c$  (0 otherwise)  
 $C$  = total number of classes  
 $K$  = batch size  
 $\alpha$  = 0.5 (trade-off parameter)

### Research Questions

- (RQ1) Which type of training produces better Teacher models: human-guided or purely data-driven?
- (RQ2) Can the top-performing Teacher model improve the performance of Student models across different CNN architectures?
- (RQ3) What are the potential performance benefits of the Teacher-Student training paradigm over the baselines?
- (RQ4) Can this training approach be applied to domains beyond synthetic face detection?

### Once Taught By Humans, Models Can Teach Themselves

Table 2: Area Under the Curve (AUC, mean  $\pm$  std) achieved by the baselines and the optimal student model in synthetic face detection task. Optimal AI Student configurations were achieved by training the model using the optimal AI Teacher's configuration (Xception + CAM) with  $\alpha = 0.01$  (i.e. encouraging the model to focus on saliency instead of class label).

Model	Baseline 1 (small set, entire human saliency available)	Baseline 2 (larger set, no human saliency)	Optimal AI Student (this paper: large set, optimal use of human saliency)
DenseNet	0.633 $\pm$ 0.04	0.629 $\pm$ 0.039	<b>0.767 <math>\pm</math> 0.020</b>
ResNet	0.612 $\pm$ 0.05	0.555 $\pm$ 0.061	<b>0.718 <math>\pm</math> 0.012</b>
Xception	0.730 $\pm$ 0.02	0.586 $\pm$ 0.074	<b>0.743 <math>\pm</math> 0.005</b>
Inception	0.679 $\pm$ 0.03	0.610 $\pm$ 0.035	<b>0.746 <math>\pm</math> 0.019</b>

