





## Semi-supervised novelty detection using ensembles with regularized disagreement

Alexandru Țifrea, Eric Stavarache, Fanny Yang Department of Computer Science, ETH Zurich



## **KEY 1: ROLE OF REGULARIZATION Goal:** Prevent complex models from interpolating on $S \cup (U, c)$ . 🐹 ensemble predictions 🗱 identified as OOD in U • 🛚 🖊 Good func. complexity Large func. complexity **Advantages of early stopping:** We prove that there exists an **optimal stopping time** at which every model predicts: (1) the correct label on ID data; and (2) the arbitrary label on the OOD unlabeled data. Efficient model selection (requires only one training run). **KEY 2: ENSEMBLE DISAGREEMENT SCORE Prior work:** Entropy of average predictor $(H \circ Avg)$ . **Our average pairwise disagreement score:** $(\operatorname{Avg} \circ \rho)(f_1(x), ..., f_K(x)) := \frac{2}{K(K-1)} \sum_{i=1}^{2} \rho(f_i(x), f_j(x))$ $\rightarrow$ e.g. $\rho$ = total variation distance **Unlike** (**H** $\circ$ **Avg**), **our score exploits ensemble diversity**. $(H \circ Avg)$ obtains lower FPR at the same TPR. TP = correctly flagged OOD point; FP = ID flagged as OOD. FPR > 0 Equal TPR red blue Novel-disease medical image data (H ∘ Avg)





	ID support
	(2 0105565)
<b>\</b>	Ensemble models
<b>\</b>	Averaged model
	<b>T</b>
	True positives
	True negatives
	False negatives
	False positives