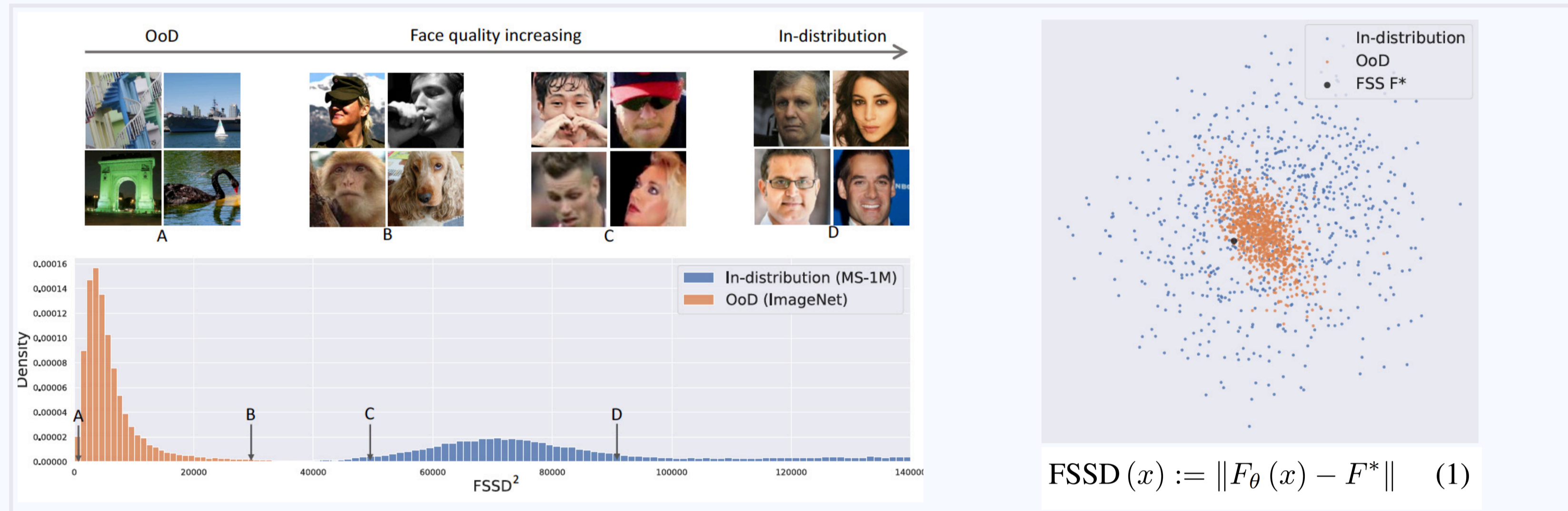


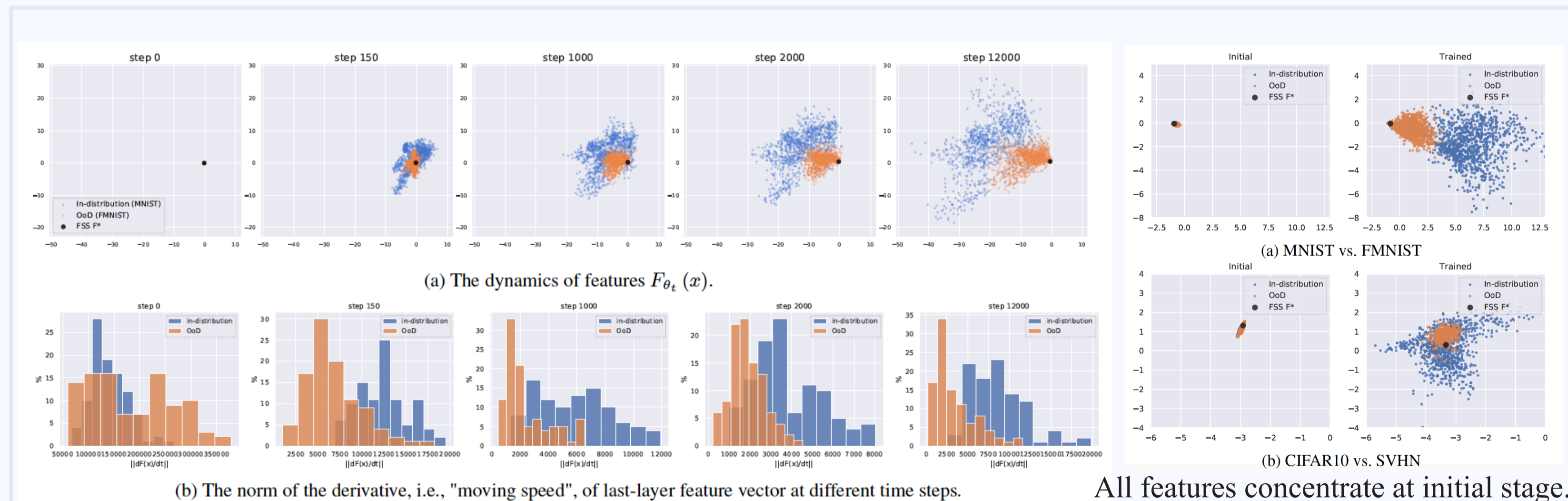
# Feature Space Singularity for Out-of-Distribution Detection

Haiwen Huang, Zhihan Li, LuluWang, Sishuo Chen, Bin Dong, Xinyu Zhou  
Corresponding: haiwen.huang2@cs.ox.ac.uk

**Observation** In a trained NN, OoD samples concentrate in the feature space.



**Understanding** OoD features move slowly during training.



"Moving speed" of the feature vector

$$\frac{dF_{\theta_t}(x)}{dt} = \frac{\partial F_{\theta_t}(x)}{\partial \theta_t} \frac{d\theta_t}{dt} = - \sum_{m=1}^M \frac{\partial F_{\theta_t}(x)}{\partial \theta_t} \frac{\partial F_{\theta_t}(x_m)^T}{\partial \theta_t} \partial_m \mathcal{L}_\phi.$$

Empirical Neural Tangent Kernel

Integrate

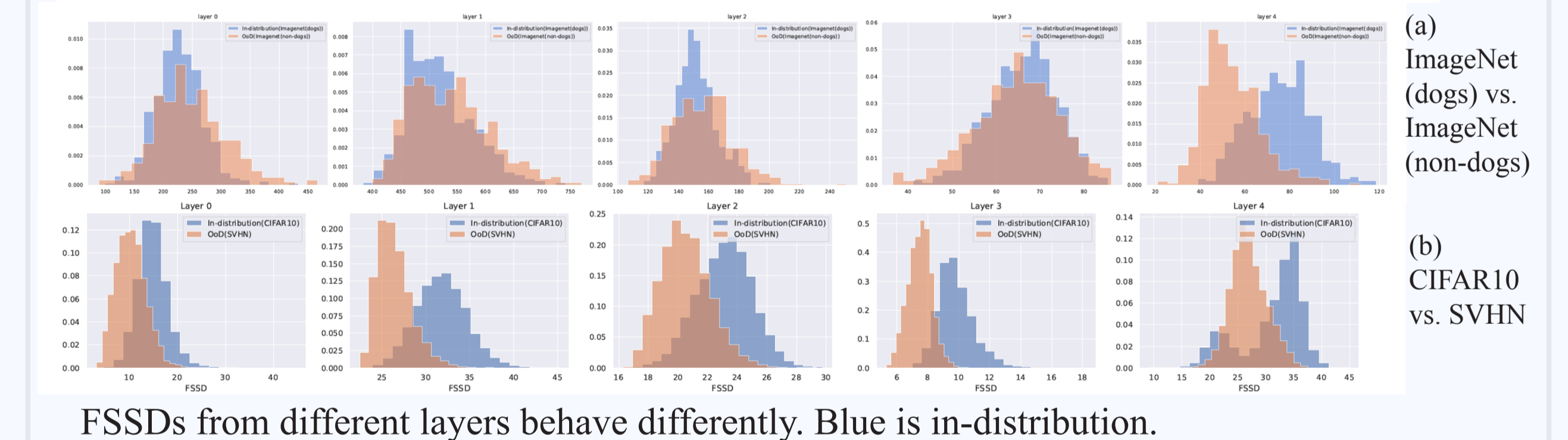
$$\begin{aligned} \text{FSSD}(x) &\stackrel{\text{Equation (1)}}{\approx} \left\| F_{\theta_0}(x) - F_{\theta_t}(x) \right\| \\ &= \left\| \sum_{m=1}^M \int_0^T \Theta_t(x, x_m) \partial_m \mathcal{L}_\phi dt \right\| \\ &\approx \left\| \sum_{m=1}^M \Theta(x, x_m) \nu_m \right\|, \end{aligned}$$

where  $\nu_m = \int_0^T \partial_m \mathcal{L}_\phi dt$ .

## Experiments

Table 2: Main results. All values are in %.

Datasets (Architecture)	Metrics	Base	ODIN	Maha	DE	MCD	OE	FSSD	
Small-scale benchmarks	FMNIST vs. MNIST (LeNet)	AUROC	77.3	96.9	<b>99.6</b>	83.9	81.7	<b>99.6</b>	<b>99.6</b>
		AUPRC	79.2	93.0	<b>99.7</b>	83.3	85.3	99.6	<b>99.7</b>
		FPR80	43.5	2.5	<b>0.0</b>	27.5	36.8	<b>0.0</b>	<b>0.0</b>
Large-scale benchmarks	CIFAR10 vs. SVHN (ResNet34)	AUROC	89.9	96.7	99.1	93.7	96.7	90.4	<b>99.5</b>
		AUPRC	85.4	92.5	98.1	90.6	93.9	89.8	<b>99.5</b>
		FPR80	10.1	4.7	<b>0.3</b>	3.7	2.4	12.5	0.4
Sequence benchmark	Bacteria Genome (LSTM)	AUROC	69.6	70.6	70.4	70.0	69.3	NA	<b>74.8</b>
		AUPRC	69.9	71.9	69.3	56.0	70.2	NA	<b>75.8</b>
		FPR80	57.4	55.9	53.7	<b>30.0</b>	58.3	NA	47.4



## Our Algorithm

Layer ensemble helps

### Algorithm 1: Computation of FSSD-Ensem

**Input:** Test samples  $x = \{x_n^{\text{test}}\}_{n=1}^N$ , ensemble weights  $\alpha_k$ , perturbation magnitude  $\epsilon$ , feature extractors  $\{F_{(k)}\}_{k=1}^K$   
**for each feature extractor**  $\{F_{(k)}\}_{k=1}^K$  **do**  
 1. Estimate FSS  $F_{(k)}^* = \sum_{s=1}^S F_{(k)}(x_s^{\text{noise}}) / S$ , where  $x_s^{\text{noise}} \sim \mathcal{U}[0, 1]$ ,  $s = 1, \dots, S$   
 2. Add perturbation to test sample:  $\tilde{x} = x + \epsilon \text{sign}(\nabla_x \|F_{(k)}(x) - F_{(k)}^*\|)$   
 3. Calculate  $\text{FSSD}^{(k)}(x) = \|F_{(k)}(\tilde{x}) - F_{(k)}^*\|$   
**end**  
**Return**  $\text{FSSD-Ensem}(x) = \sum_{k=1}^K \alpha_k \text{FSSD}^{(k)}(x)$

## Future Work

1. Study the phenomenon in different phases of training corresponding to the recent advances in NTK;
2. Explore the connection to the probabilistic models, e.g. Gaussian Processes.