

# Enhancing Diversity in Bayesian Deep Learning via Hyperspherical Energy Minimization of CKA

David Smerkous<sup>1</sup>, Qinxun Bai<sup>2</sup>, Fuxin Li<sup>1</sup>

<sup>1</sup>Oregon State University

<sup>2</sup>Horizon Robotics



地平线  
Horizon Robotics

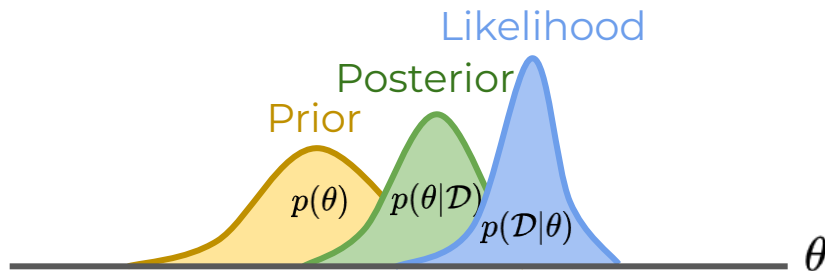
# Bayesian Deep Learning

**Bayesian Approach:** Instead of point estimates, BDL estimates the distribution of the model parameters with prior information

$$p(\theta|\mathcal{D}) \propto p(\mathcal{D}|\theta)p(\theta)$$

**Uncertainty Estimation:** Improves understanding of prediction uncertainty

**Applications:** Transfer learning, fairness, active learning, reinforcement learning, robotics



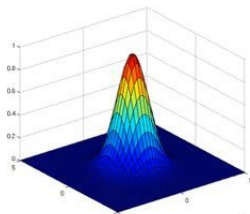
# Particle Variational Inference (ParVI)

**Basic:** MC Dropout, Stochastic Gradient Langevin Dynamics (SGLD), and Ensemble methods.

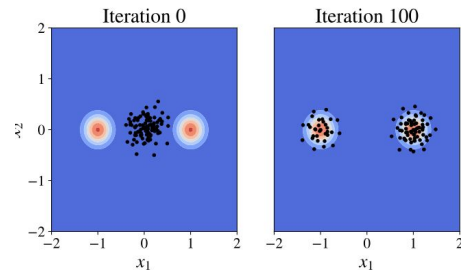
**Stein Variational Gradient Descent (SVGD):** iterative gradient-based variational inference algorithm on a finite set of particles  $\mathbf{M}$ .

$$x_i^{\ell+1} \leftarrow x_i^{\ell} + \epsilon_{\ell} \hat{\phi}^*(x_i^{\ell}) \quad \text{where} \quad \hat{\phi}^*(x) = \frac{1}{n} \sum_{j=1}^n \left[ \overbrace{k(x_j^{\ell}, x) \nabla_{x_j^{\ell}} \log p(x_j^{\ell})}^{\text{Driving}} + \overbrace{\nabla_{x_j^{\ell}} k(x_j^{\ell}, x)}^{\text{Repulsive}} \right]$$

**Kernel Selection:** Overwhelmingly popular choice is Radial Basis Function (RBF) with median heuristic.



$$K(\vec{x}, \vec{l}) = e^{-\frac{\|\vec{x} - \vec{l}\|^2}{2\sigma^2}}$$



# Issues with RBF and Ensemble Diversity

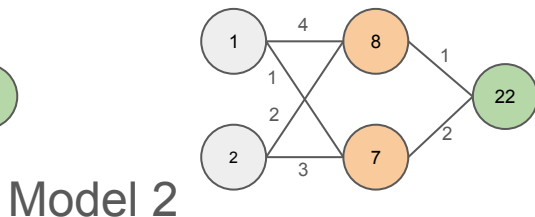
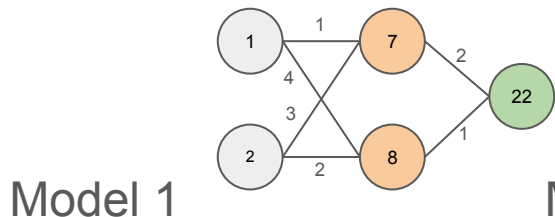
**Posterior Representation:** Particles that are more diverse better capture the posterior. Supported by methods like SVGD.

**Challenge:** what should we consider diverse and how should we better measure it?

**RBF:** Works with small networks, but fails with larger ones (D'Angelo et al.)

**Current limitations:**

- Euclidean based kernels - using  $L_1, L_2$  - lack scale and permutation invariance
- Curse of dimensionality in high-dimensional parameter spaces



$$L_1(M_1, M_2) = 10$$

# Centered Kernel Alignment

**Goal:** construct a differentiable kernel that is invariant to isotropic scaling and permutations

**Solution:** Centered Kernel Alignment

## CKA between two models

$$\text{HSIC}(K^1, K^2) = \frac{1}{(N-1)^2} \text{tr}(K^1 H K^2 H) \quad H = I - \frac{1}{N} \mathbf{1}\mathbf{1}^\top$$

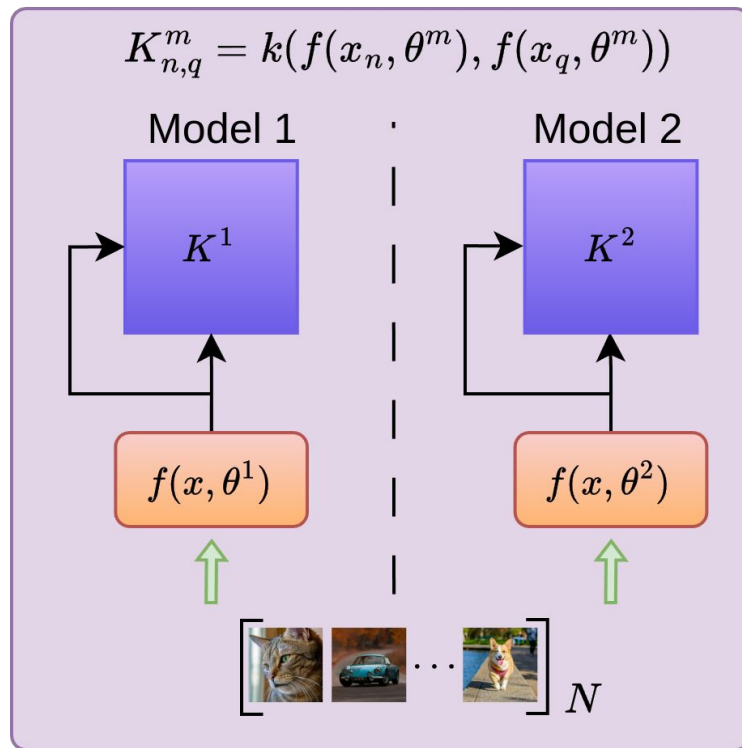
$$\text{CKA}(K^1, K^2) = \frac{\text{HSIC}(K^1, K^2)}{\sqrt{\text{HSIC}(K^1, K^1) \text{HSIC}(K^2, K^2)}}$$

## Ensemble metric: Pairwise CKA

$$\text{CKA}_{\text{pw}}(\mathcal{K}) = \frac{1}{LM(M-1)} \sum_{l=1}^L \sum_{\substack{m, m'=1 \\ m \neq m'}}^{M, M} \text{CKA}(K_l^m, K_l^{m'})$$

## Gram Matrix Construction

$$K_{n,q}^m = k(f(x_n, \theta^m), f(x_q, \theta^m))$$

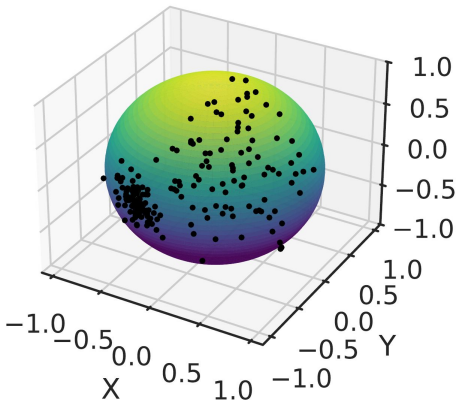


# Hyperspherical Energy

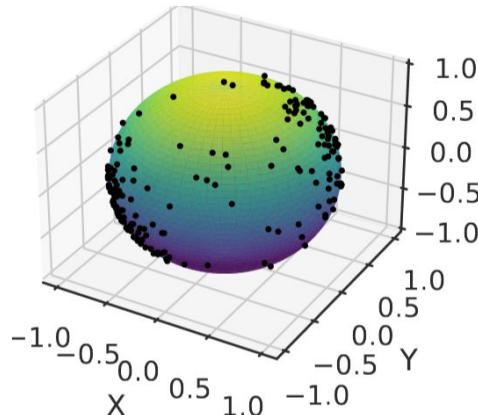
**Alternative perspective of pairwise CKA:** projection of feature gram matrices vectorized and projected onto unit *hypersphere* and minimization of pairwise cosine similarity.

**Potential deficiency:** small gradients when model CKAs are similar as gradient of  $\cos$  is  $-\sin$

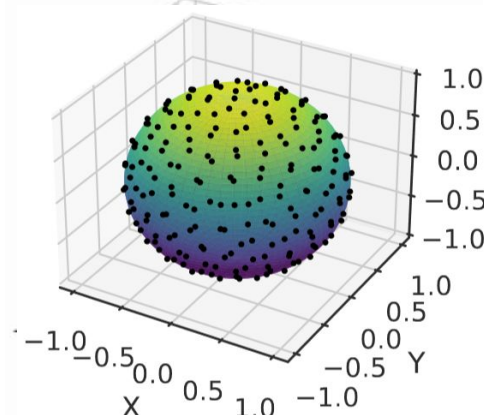
**Initial Points**



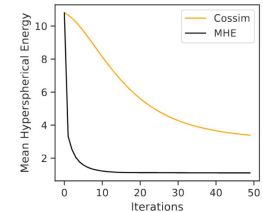
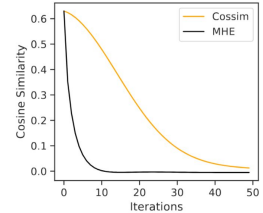
**Pairwise Cossim**



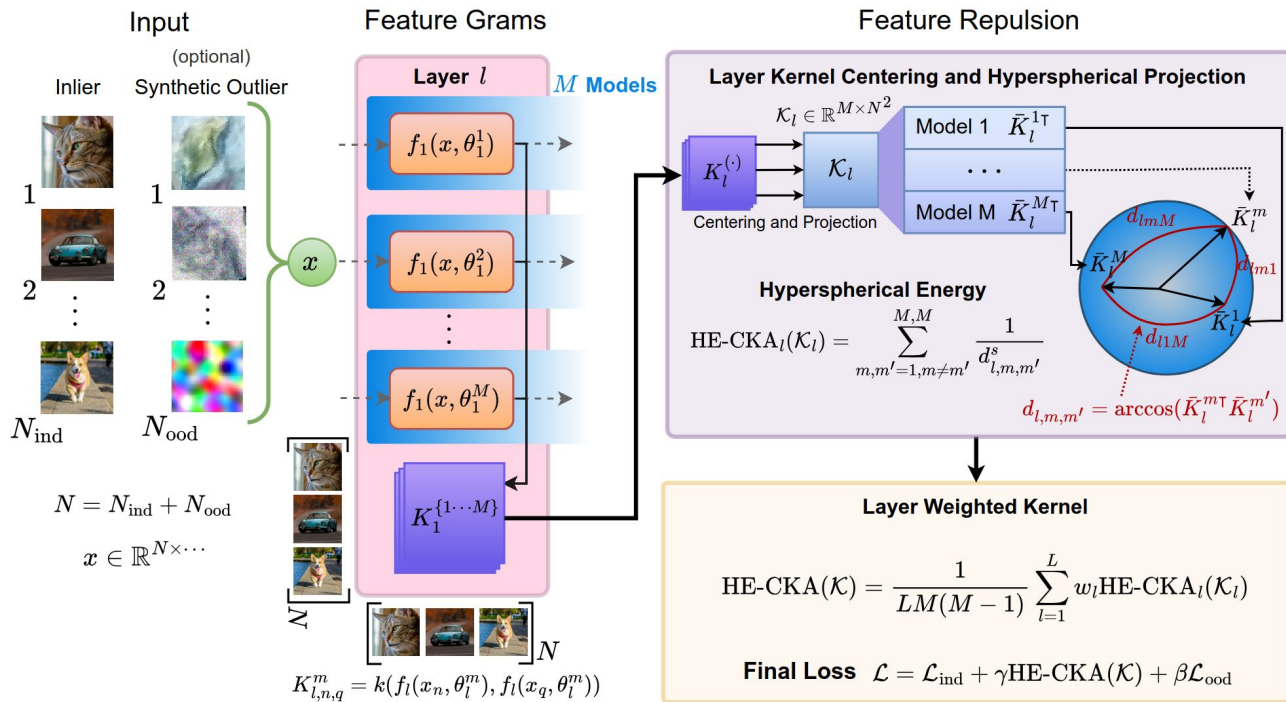
**Hyperspherical Energy**



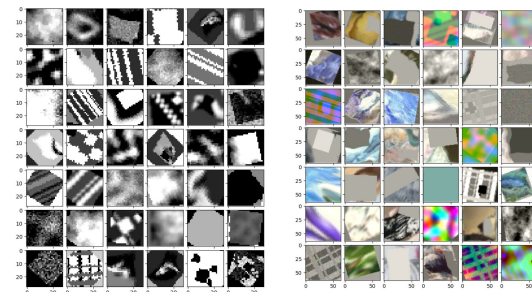
**HE vs Cossim**



# Architecture



## Synthetic OOD



**MNIST** **TinyImageNet**

**Noise:** Perlin, Simplex, Gaussian, random shapes

**ID2OOD:** augmentations such as blur, elastic transformation, cutout, etc

# Synthetic four class classification

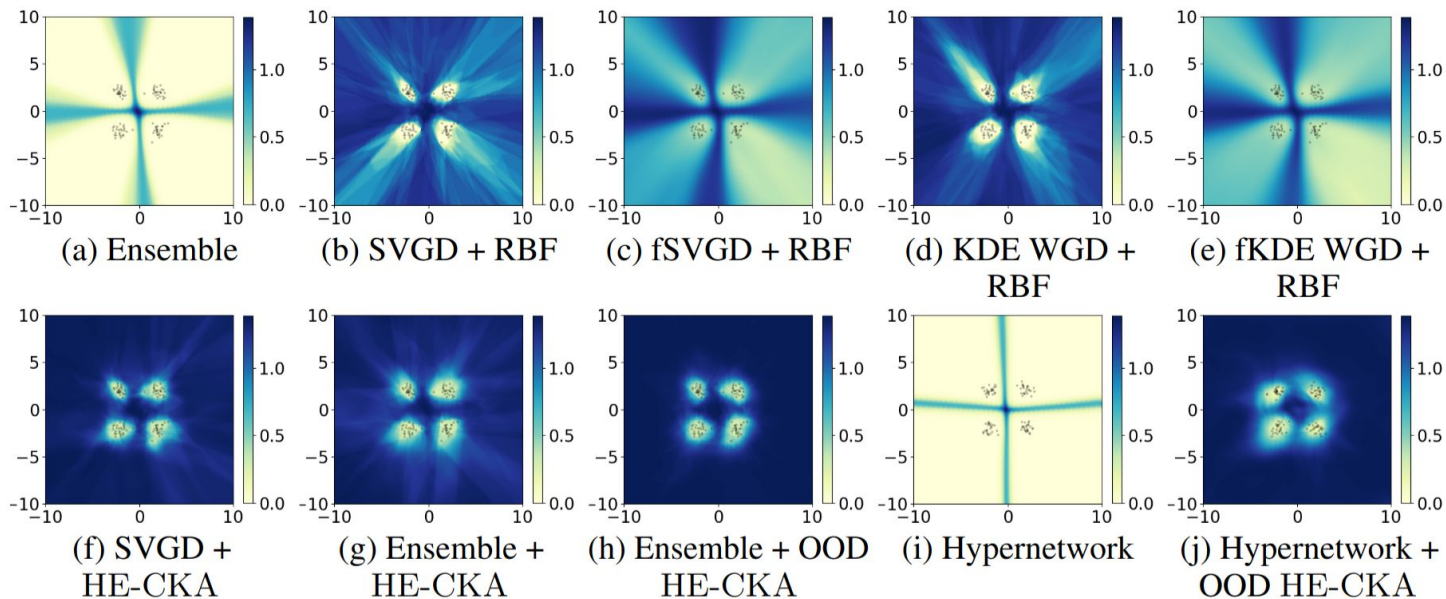


Figure 3: Predictive entropies (PE) on a four-cluster 2D classification task. Darker values indicate higher entropy, lower confidence regions, and lighter values indicate higher confidence regions. (b) and (d) use an RBF kernel on ensemble member weights, whereas (c) and (e) use an RBF kernel on ensemble member outputs. (f) and (g) use the HE-CKA, RBF feature kernel, for feature diversity on inlier points. Both (h) and (j) use HE-CKA and OOD entropy terms. All methods were trained on an ensemble of 30 four layer MLPs for 1k iterations with the same seeds.



# Dirty-MNIST and FashionMNIST

Inlier

Outlier

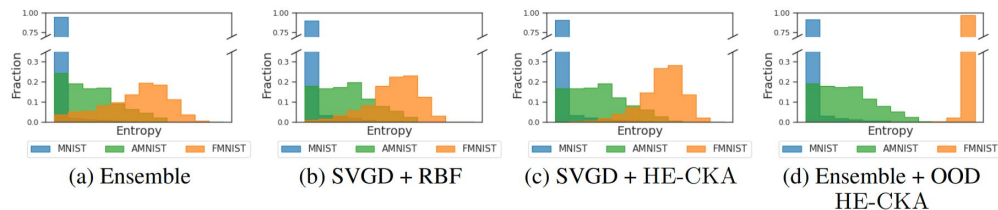
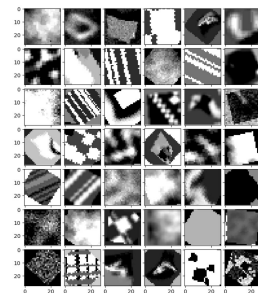


Figure 6: Predictive softmax entropy between MNIST, Dirty-MNIST (with aleatoric uncertainty), and OOD Fashion-MNIST. Utilizing an ensemble of 5 LeNets. It can be seen that HE-CKA and OOD HE-CKA better separates the inlier Dirty-MNIST from outlier Fashion-MNIST.

Table 1: OOD detection results with inlier Dirty-MNIST and outlier Fashion MNIST, over 5 runs. All models were trained on a LeNet, with HE-CKA and  $CKA_{pw}$  utilizing a cosine similarity feature kernel. One exception to predictive entropy (PE) report is DDU, which uses feature space density, indicated by a star, to calculate AUROC [Mukhoti et al. \(2023\)](#). More training details can be found in [Appendix C](#).

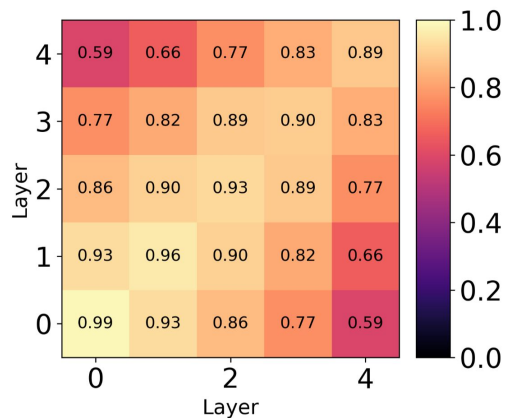
MODEL	NLL ( $\downarrow$ )	ACCURACY ( $\uparrow$ )	ECE ( $\downarrow$ )	AUROC FASHIONMNIST ( $\uparrow$ )	
				PE	MI
DDU	0.278 $\pm$ 0.001	82.177 $\pm$ 0.032	3.952 $\pm$ 0.317	94.168 $\pm$ 3.425*	—
SINGLE	<b>0.272 <math>\pm</math> 0.002</b>	82.299 $\pm$ 0.166	2.908 $\pm$ 0.129	65.935 $\pm$ 10.669	50.000 $\pm$ 0.000
ENSEMBLE	<b>0.271 <math>\pm</math> 0.001</b>	83.915 $\pm$ 0.084	<b>1.306 <math>\pm</math> 0.098</b>	86.095 $\pm$ 1.608	96.065 $\pm$ 0.798
SVGd+RBF	0.304 $\pm$ 0.001	83.560 $\pm$ 0.072	3.178 $\pm$ 0.076	91.003 $\pm$ 1.155	98.083 $\pm$ 0.516
SVGd+CKA <sub>pw</sub>	0.377 $\pm$ 0.003	82.351 $\pm$ 0.150	7.359 $\pm$ 0.211	89.195 $\pm$ 4.260	<b>99.207 <math>\pm</math> 0.160</b>
SVGd+HE-CKA	0.298 $\pm$ 0.002	83.879 $\pm$ 0.110	2.846 $\pm$ 0.153	94.380 $\pm$ 1.332	<b>99.213 <math>\pm</math> 0.147</b>
HYPERNET	0.278 $\pm$ 0.003	81.157 $\pm$ 0.174	4.253 $\pm$ 0.092	46.393 $\pm$ 3.545	64.856 $\pm$ 2.768
HYPERNET+OOD HE-CKA	0.325 $\pm$ 0.014	82.398 $\pm$ 0.628	3.058 $\pm$ 0.665	98.073 $\pm$ 0.951	77.548 $\pm$ 9.526
ENSEMBLE+HE-CKA	0.306 $\pm$ 0.001	83.684 $\pm$ 0.029	3.174 $\pm$ 0.120	94.656 $\pm$ 1.095	98.866 $\pm$ 0.148
ENSEMBLE+OOD HE-CKA	0.277 $\pm$ 0.001	<b>84.090 <math>\pm</math> 0.049</b>	1.712 $\pm$ 0.061	<b>99.996 <math>\pm</math> 0.001</b>	<b>99.742 <math>\pm</math> 0.506</b>

Synthetic OOD

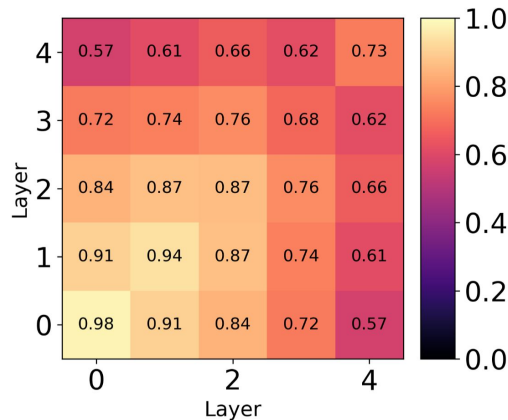


# Dirty MNIST Ensemble CKA Plots

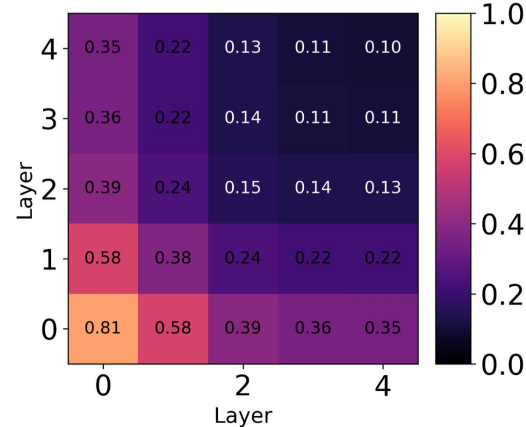
## Deep Ensemble



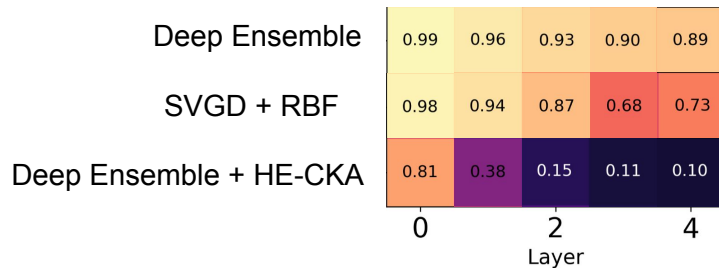
## SVGD + RBF



## Deep Ensemble + HE-CKA



## Diagonals



# CIFAR10 / TinyImageNet

Table 3: OOD results on CIFAR-10 vs SVHN. Methods used a ResNet18 ensemble of size 5. (·) indicates ensemble size.

MODEL	OOD METHOD	NLL ( $\downarrow$ )	ACCURACY ( $\uparrow$ )	ECE ( $\downarrow$ )	AUROC PE SVHN ( $\uparrow$ )
WIDERESNET28-10+SN <sub>MUKHOTI ET AL. [2023]</sub> (1)	GMM	–	95.97	0.85	97.86
ENSEMBLE	PE	0.122	<b>96.34</b>	1.08	96.18
SVGD+RBF	PE	0.143	95.71	1.19	95.37
SVGD+CKA <sub>pw</sub>	PE	0.125	96.25	<b>0.41</b>	96.07
SVGD+HE-CKA	PE	0.124	96.23	0.63	96.01
ENSEMBLE+HE-CKA	PE	<b>0.120</b>	96.23	0.59	96.71
ENSEMBLE+OOD HE-CKA	PE	0.123	96.24	0.58	<b>99.86</b>

Table 4: Performance of a five member ResNet18 ensemble trained on TinyImageNet. All models utilized a pretrained deep ensemble with no repulsive term, then fine tuned for 30 epochs for each method (including deep ensemble). Methods utilizing CKA<sub>pw</sub> and HE-CKA utilized a linear feature kernel.

MODEL	NLL ( $\downarrow$ )	ID ACCURACY ( $\uparrow$ )	ECE ( $\downarrow$ )	AUROC PE ( $\uparrow$ )		
				SVHN	CIFAR 10/100	TEXTURES (DTD)
ENSEMBLE	0.775	62.95	8.90	89.81	66.85/67.33	68.96
SVGD+RBF	0.926	61.87	16.10	92.76	72.23/73.73	65.67
SVGD+CKA <sub>pw</sub>	0.835	60.15	8.26	94.08	78.40/79.48	66.48
SVGD+HE-CKA	<b>0.732</b>	61.36	<b>3.71</b>	94.10	72.05/72.86	70.75
ENSEMBLE+HE-CKA	0.784	<b>63.10</b>	9.82	92.65	72.13/71.68	70.69
ENSEMBLE+OOD HE-CKA	0.786	61.88	8.02	<b>99.31</b>	<b>81.56/87.64</b>	<b>90.94</b>

# CKA for Bayesian Deep Learning

**CKA:** Provides a more intuitive kernel to measure model similarity compared to RBF

**HE:** Realized pairwise CKA can be reformulated as a hyperspherical energy optimization problem.

**Synthetic Outlier:** Increase feature diversity on obvious outlier examples and provide negative signals for likelihood. *Significant improvement in outlier detection* without sacrificing inlier performance.

## Limitations:

- Larger memory requirements (feature gram matrices/gradients)
- Requires tuning of multiple hyperparameters for good performance (ex: layer weighting)
- Hyperparameters are sensitive to dataset and particle number



# Enhancing Diversity in Bayesian Deep Learning via Hyperspherical Energy Minimization of CKA

David Smerkous, Qinxun Bai, Fuxin Li

Code publicly available at



<https://github.com/Deep-Machine-Vision/he-cka-ensembles>