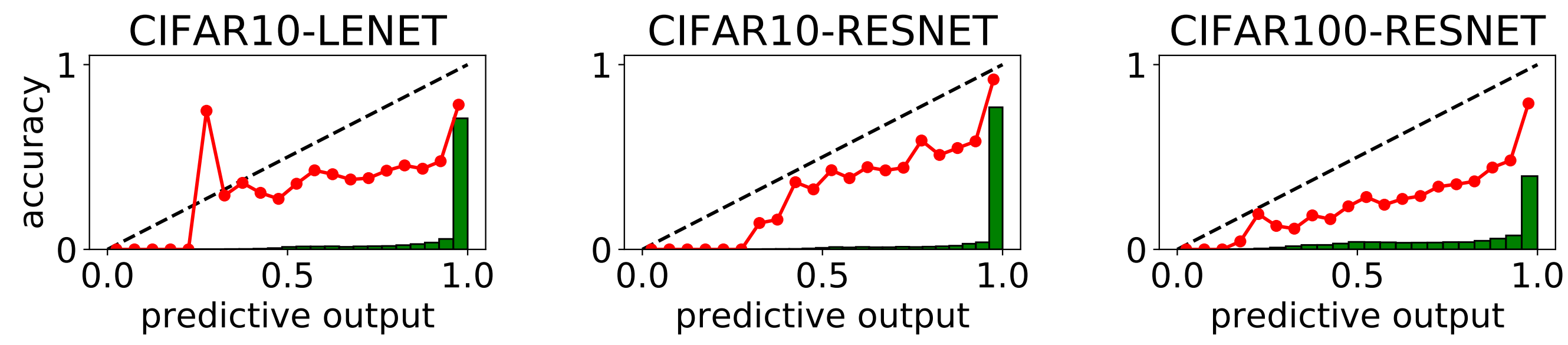
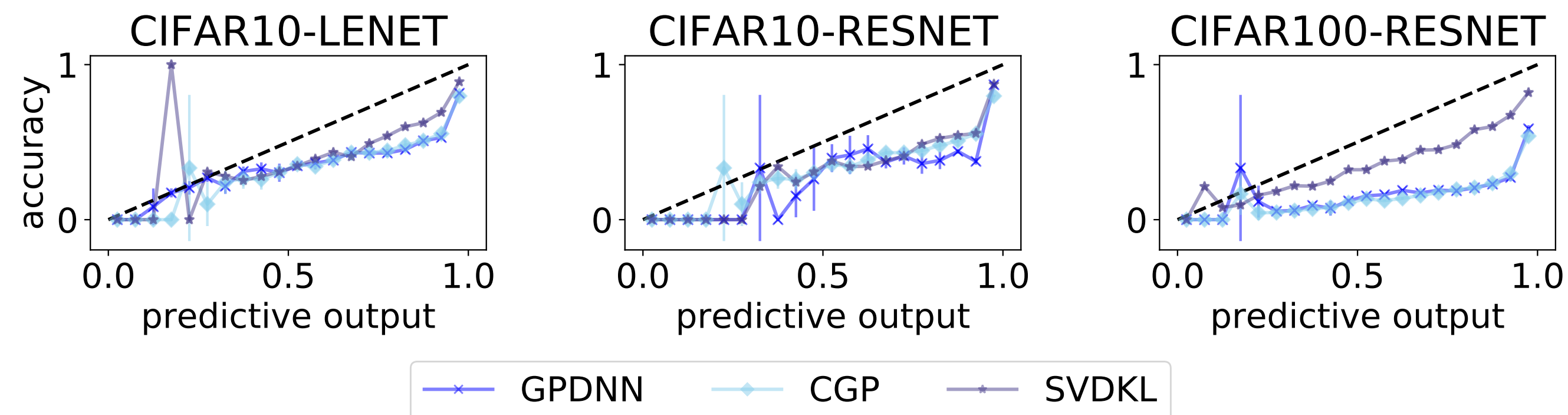


Miscalibration of Convolutional Neural Networks

- The classification networks in a decision making system must not only be accurate, but also should indicate when they are likely to be incorrect, i.e., model's calibration.
- Convolutional Neural Networks (CNNs) were found to be miscalibrated.



- The current combinations of CNNs and Gaussian Processes (GPs) are miscalibrated.



Gaussian Processes with Random Features

- Convolutional features $C(\mathbf{X}|\Psi)$ are fed into a GPs.
- Low-rank approximation for kernel matrix:

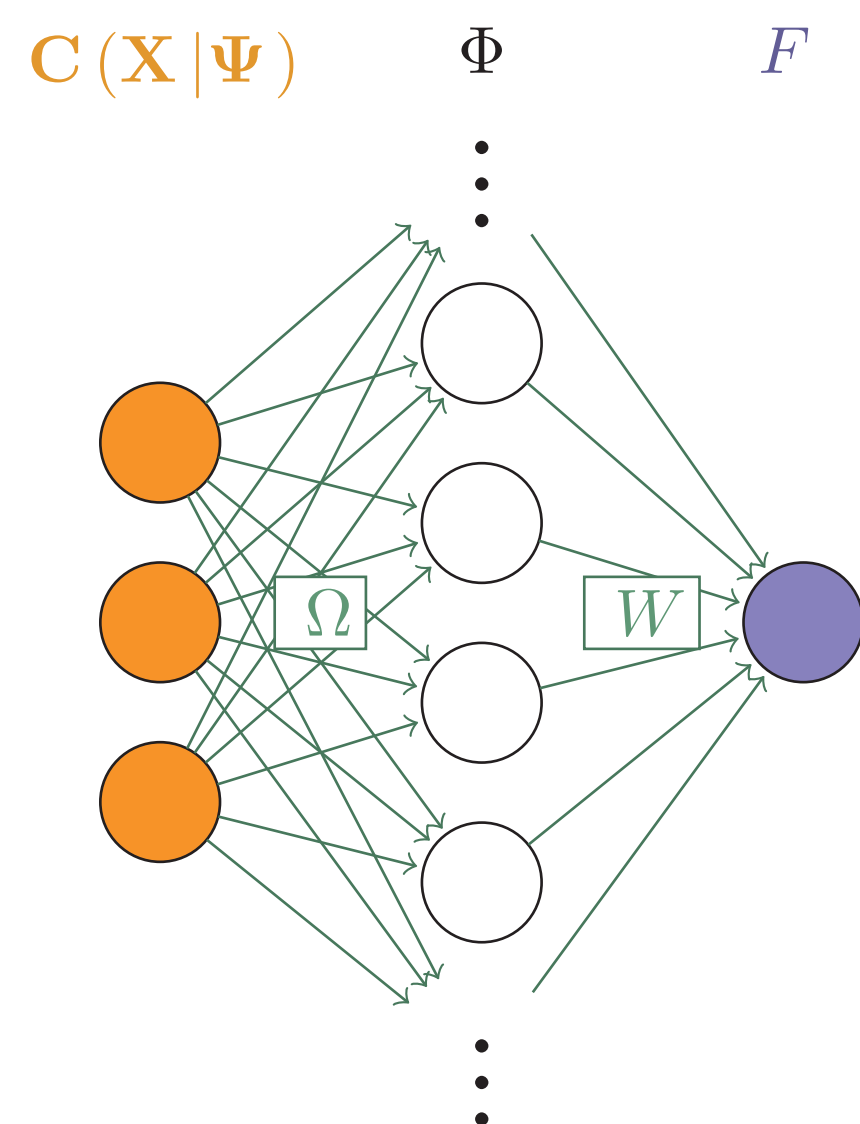
$$\mathbf{K} \approx \Phi \Phi^T$$

- Taking a weight-space view of a GP:

$$\mathbf{F} = \Phi \mathbf{W}$$

- The priors over the weights are:

$$p(\mathbf{W}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$



- The order-one ARC-COSINE kernel:

$$k_{\text{arc}}^{(1)}(\mathbf{x}_i, \mathbf{x}_j | \Psi, \theta) = \frac{\sigma^2}{\pi} \left\| \Lambda^{-\frac{1}{2}} \mathbf{c}(\mathbf{x}_i | \Psi) \right\| \left\| \Lambda^{-\frac{1}{2}} \mathbf{c}(\mathbf{x}_j | \Psi) \right\| [\sin(\alpha) + (\pi - \alpha) \cos(\alpha)], \text{ where}$$

$$\theta = (\sigma, \Lambda = \text{diag}(\ell_1^2, \dots, \ell_d^2)) \text{ and } \alpha = \cos^{-1} \left(\frac{(\Lambda^{-\frac{1}{2}} \mathbf{c}(\mathbf{x}_i | \Psi))^T (\Lambda^{-\frac{1}{2}} \mathbf{c}(\mathbf{x}_j | \Psi))}{\left\| \Lambda^{-\frac{1}{2}} \mathbf{c}(\mathbf{x}_i | \Psi) \right\| \left\| \Lambda^{-\frac{1}{2}} \mathbf{c}(\mathbf{x}_j | \Psi) \right\|} \right)$$

can be approximated using Rectified Linear Units (ReLU):

$$\Phi_{\text{arc}} = \sqrt{\frac{2\sigma^2}{N_{\text{RF}}}} \max(\mathbf{0}, C(\mathbf{X}|\Psi) \Omega), \text{ where } \Psi \text{ are convolutional parameters or filters}$$

- Ω can be structured through a series of Hadamard transformations of diagonal matrices \mathbf{D}_i with elements randomly sampled from $\{-1, +1\}$:

$$\Omega = \frac{\sqrt{d}}{l} \mathbf{H} \mathbf{D}_1 \mathbf{H} \mathbf{D}_2 \mathbf{H} \mathbf{D}_3, \text{ where } d \text{ is the dimensionality of the convolutional features}$$

Stochastic Variational Learning

- The lower bound on the log-marginal likelihood $\mathcal{L} = \log [p(\mathbf{Y}|\mathbf{X}, \theta)]$:

$$\mathcal{L} \geq E_{q(\mathbf{W}, \Omega, \Psi)} (\log [p(\mathbf{Y}|\mathbf{X}, \mathbf{W}, \Omega, \Psi, \theta)]) - \text{KL} [q(\mathbf{W}, \Omega, \Psi) \| p(\mathbf{W}, \Omega, \Psi)]$$

where $q(\mathbf{W}, \Omega, \Psi)$ approximates $p(\mathbf{W}, \Omega, \Psi|\mathbf{X}, \mathbf{Y})$.

- Factorization of approximate posterior:

$$q(\mathbf{W}, \Omega, \Psi) = q(\mathbf{W}) q(\Omega) q(\Psi) = \prod_i q(\mathbf{W}_i) \prod_i q(\Omega_i) \prod_i q(\Psi_i), \text{ where}$$

$$q(\mathbf{W}_i) = \pi_w \mathcal{N}([\mathbf{M}_w]_r, \sigma^2 \mathbf{I}) + (1 - \pi_w) \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

$$q(\Omega_i) = \pi_\Omega \mathcal{N}([\mathbf{M}_\Omega]_r, \sigma^2 \mathbf{I}) + (1 - \pi_\Omega) \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

$$q(\Psi_i) = \pi_\Psi \mathcal{N}([\mathbf{M}_\Psi]_r, \sigma^2 \mathbf{I}) + (1 - \pi_\Psi) \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

$\mathbf{M}_w, \mathbf{M}_\Omega$ and \mathbf{M}_Ψ are called as variational parameters; $\pi_w, \pi_\Omega; \pi_\Psi$ are probabilities and $\sigma^2 \approx 0$.

- Assuming a normal prior, the KL term can be approximated:

$$\text{KL} [q(\mathbf{W}, \Omega, \Psi) \| p(\mathbf{W}, \Omega, \Psi)] \approx \frac{\pi_w}{2} \|\mathbf{M}_w\|^2 + \frac{\pi_\Omega}{2} \|\mathbf{M}_\Omega\|^2 + \frac{\pi_\Psi}{2} \|\mathbf{M}_\Psi\|^2$$

- The expectation can be unbiasedly estimated using Monte Carlo and mini-batch of size m :

$$E_{q(\mathbf{W}, \Omega, \Psi)} (\log [p(\mathbf{Y}|\mathbf{X}, \mathbf{W}, \Omega, \Psi, \theta)]) \approx \frac{n}{m} \frac{1}{N_{\text{MC}}} \sum_{i=1}^{N_{\text{MC}}} \sum_{k \in \mathcal{I}_m} \log [p(\mathbf{y}_k | \mathbf{x}_k, \tilde{\mathbf{W}}^{(i)}, \tilde{\Omega}^{(i)}, \tilde{\Psi}^{(i)}, \theta)]$$

- Due to the assumption of $\sigma^2 \approx 0$, we can sample from the above posteriors as follows:

$$\tilde{\mathbf{W}}^{(i)} = \mathbf{M}_w \text{diag} [\tilde{\mathbf{z}}_w^i] \quad \text{and} \quad \tilde{\Omega}^{(i)} = \mathbf{M}_\Omega \text{diag} [\tilde{\mathbf{z}}_\Omega^i] \quad \text{and} \quad \tilde{\Psi}^{(i)} = \mathbf{M}_\Psi \text{diag} [\tilde{\mathbf{z}}_\Psi^i]$$

where $\tilde{\mathbf{z}}_w^i, \tilde{\mathbf{z}}_\Omega^i$ and $\tilde{\mathbf{z}}_\Psi^i$ are sampled from Bernoulli(π_w), Bernoulli(π_Ω) and Bernoulli(π_Ψ).

Experiments

- Competing methods:

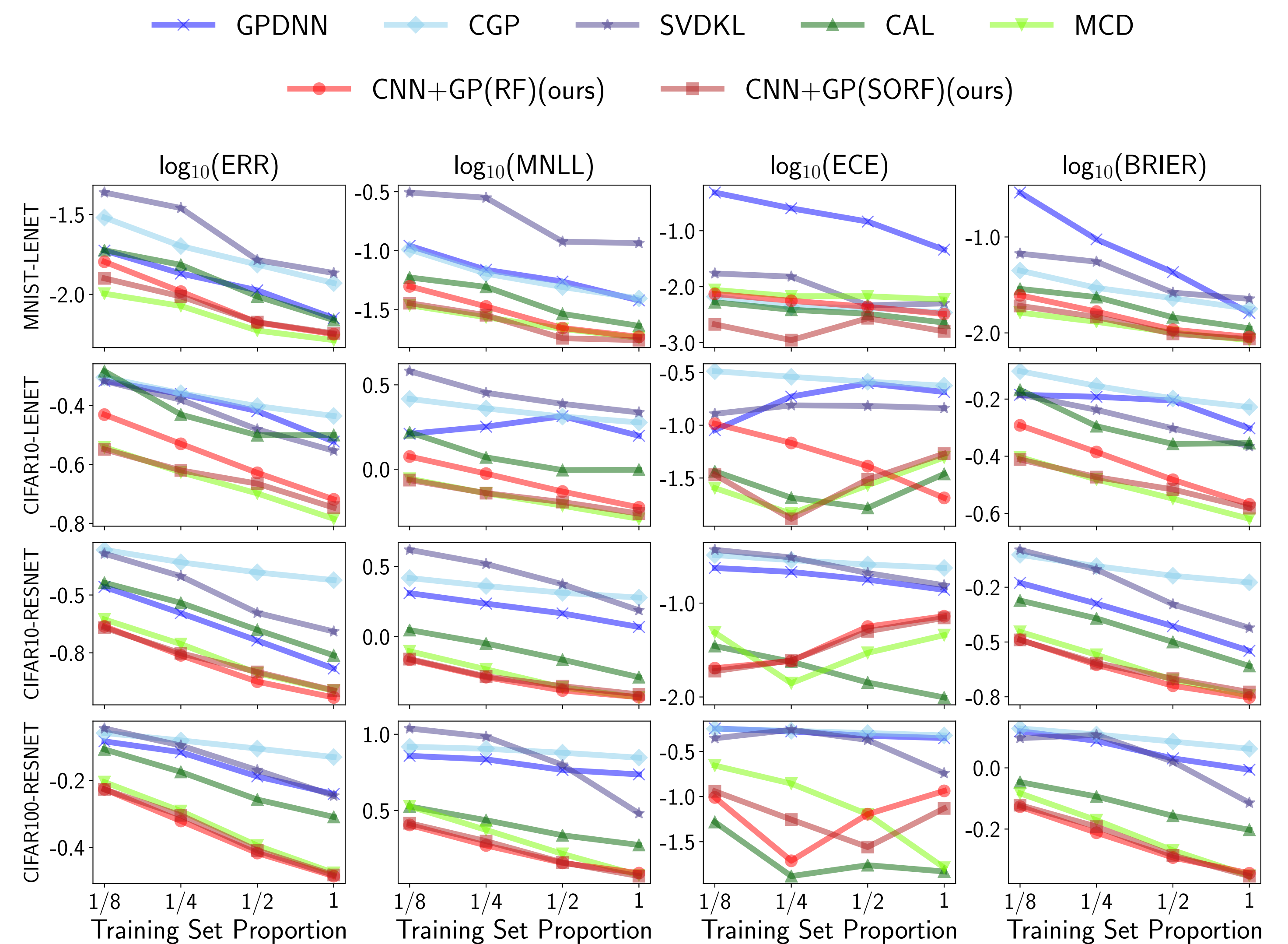


Figure 1: Comparison of our CNN+GP(RF) and CNN+GP(SORF) with competitors.

- Reliability diagram:

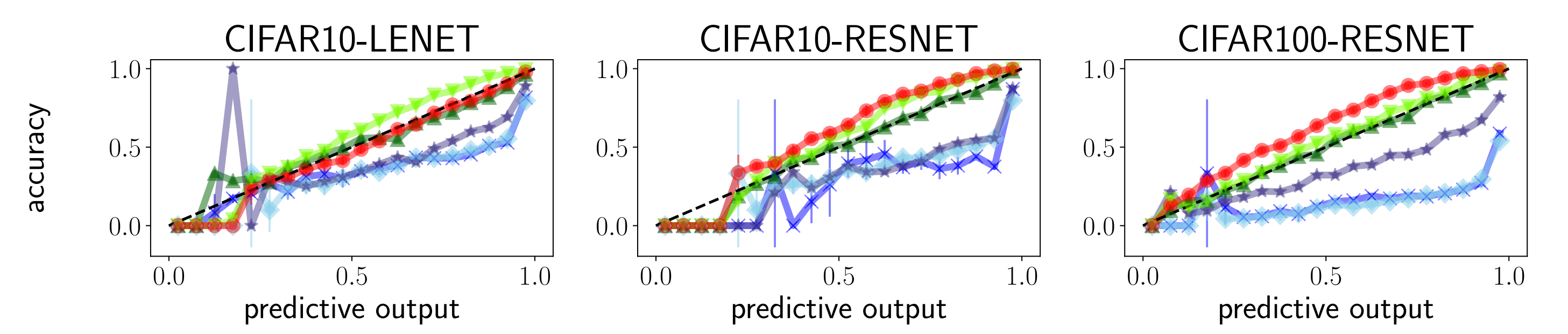


Figure 2: Reliability diagrams of our CNN+GP(RF) in comparison with competitors.

- Combination of CNNs and deep GPs:

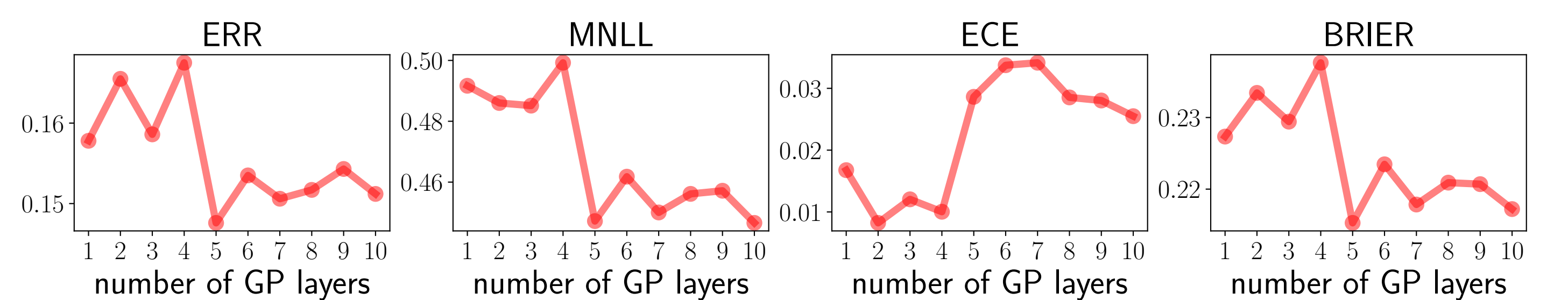


Figure 3: Performance of the proposed model when varying the number of GP layers on top of a LE NET convolutional structure on CIFAR10 dataset.

- Ability to detect out-of-distribution samples:

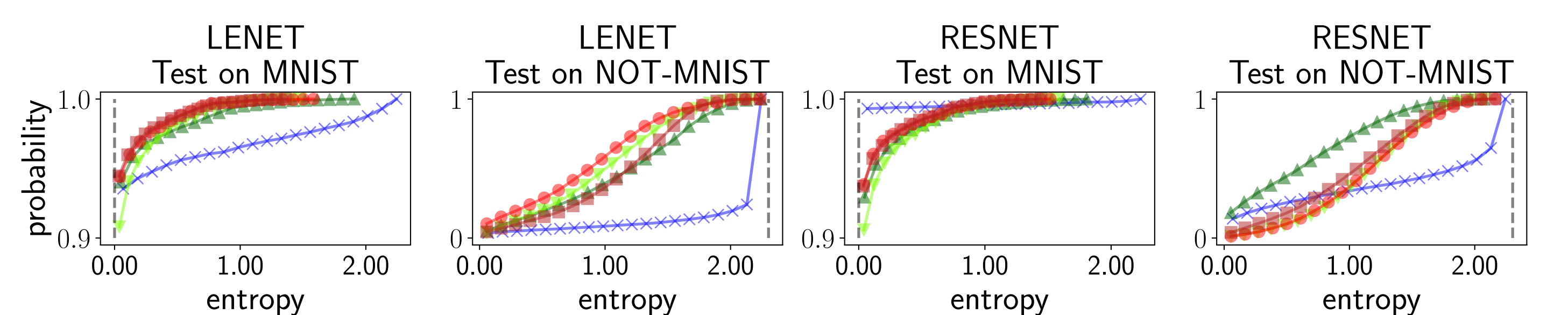


Figure 4: Cumulative distribution function plot of predictive entropies when the models trained on MNIST are tested on MNIST and NOT-MNIST.

- Substitution of fully connected layers with Structured Orthogonal Random Features (SORF):

Table 1: CIFAR10

METHOD	ERR	MNLL	ECE	BRIER
SORF	0.172	0.522	0.063	0.250
MCD	0.181	0.591	0.110	0.276

Table 2: CIFAR100

METHOD	ERR	MNLL	ECE	BRIER
SORF	0.459	1.806	0.127	0.612
MCD	0.594	2.434	0.058	0.732

Contributions

- Show that current combinations of CNNs and GPs are miscalibrated.
- Propose a novel combination of CNNs and GPs that is well-calibrated.
- Extend the model by replacing the last layer of CNNs with deep GPs.
- Obtain a compact approximation of GP by using SORF.

References

- Guo et al. *On Calibration of Modern Neural Networks*, ICML 2017.
- Bradshaw et al. *Adversarial Examples, Uncertainty, and Transfer Testing Robustness in Gaussian Process Hybrid Deep Networks*, <https://arxiv.org/abs/1707.02476>.
- Wilk et al. *Convolutional Gaussian Processes*, NeurIPS 2017.
- Wilson et al. *Stochastic Variational Deep Kernel Learning*, NeurIPS 2016.
- Cutajar et al. *Random Feature Expansions for Deep Gaussian Processes*, ICML 2017.
- Yu et al. *Orthogonal Random Features*, NeurIPS 2016.
- Gal and Ghahramani. *Dropout as a Bayesian Approximation*, ICML 2016.