



Accelerated Linearized Laplace Approximation for Bayesian Deep Learning

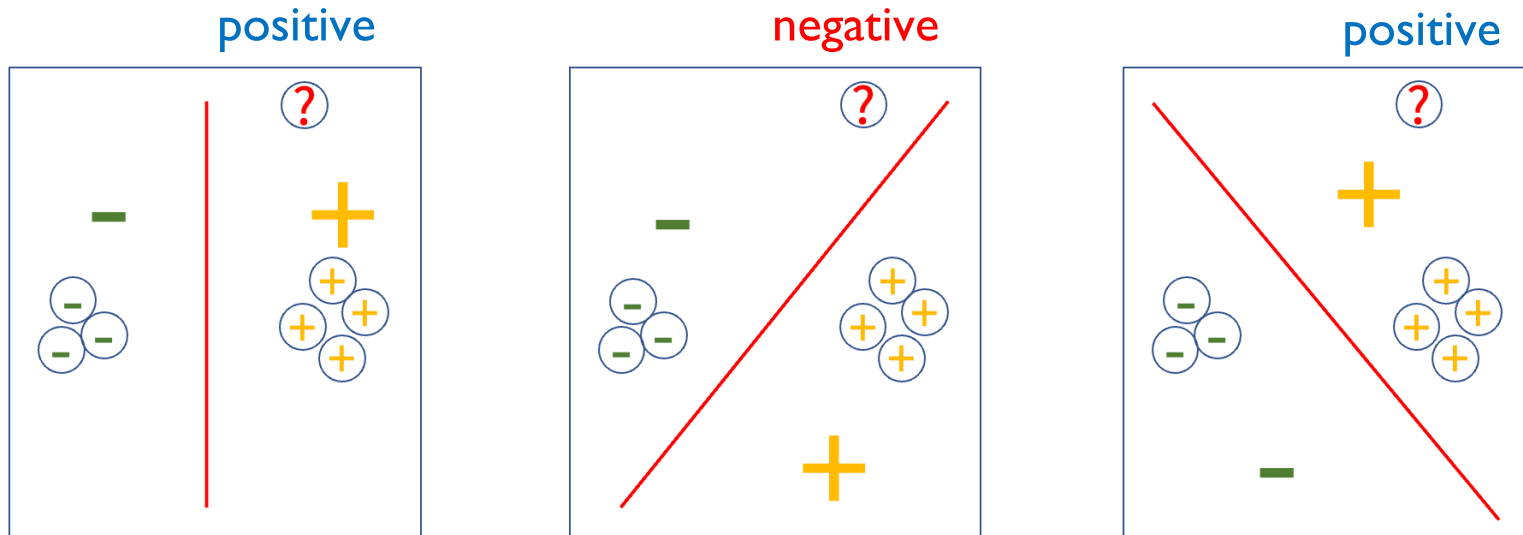
Zhijie Deng

Shanghai Jiao Tong University

Contact: zhijied@sjtu.edu.cn

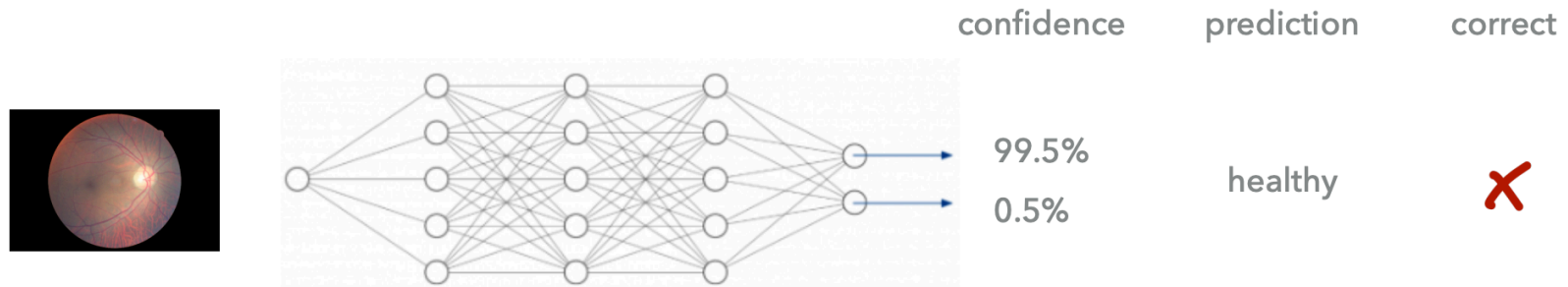
Joint work with: Feng Zhou and Jun Zhu

The issues of deterministic neural networks



Which model to select?

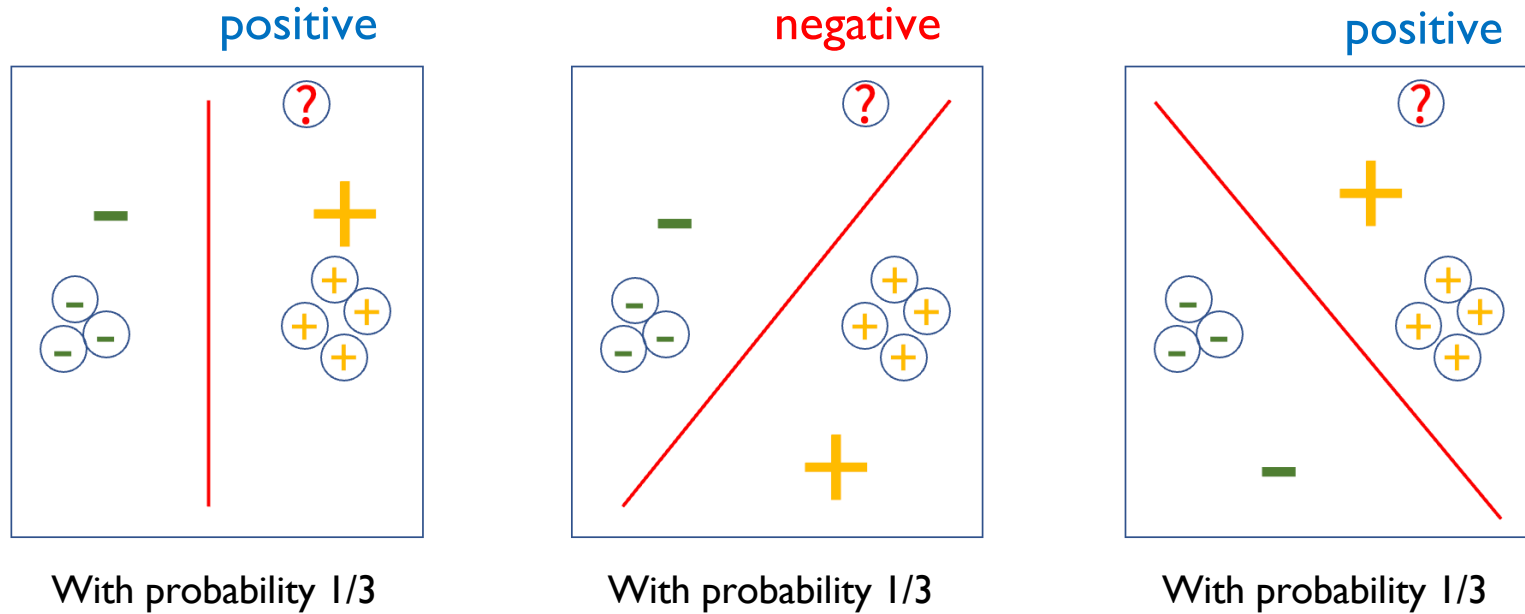
The issues of deterministic neural networks



["Benchmarking Bayesian Deep Learning with Diabetic Retinopathy Diagnosis" by Angelos Filos et al.]

Can we trust the prediction with such confidence?

Bayesian deep learning is helpful here

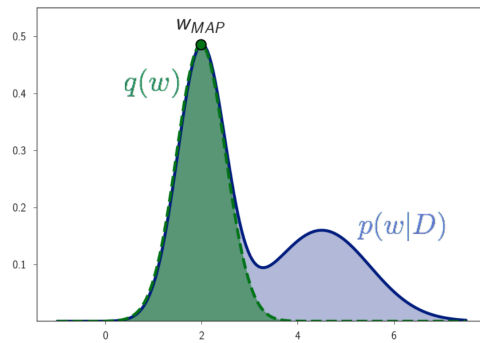


A posterior distribution over solutions: $p(\mathbf{w}|\mathcal{D}) \propto p(\mathbf{w}) \prod_{n=1}^N p(y_n|\mathbf{x}_n, \mathbf{w})$

Marginalization for prediction $p(y|x_*, \mathcal{D}) = \int p(y|x_*, w)p(w|\mathcal{D})dw$

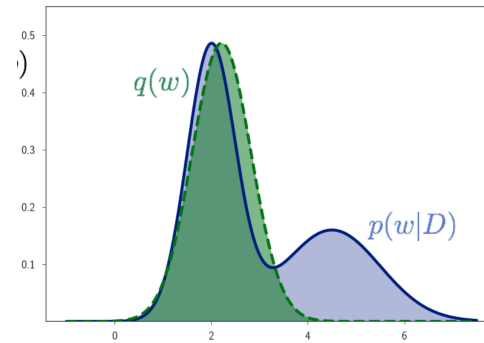
Bayesian neural networks

Approximate Bayesian inference



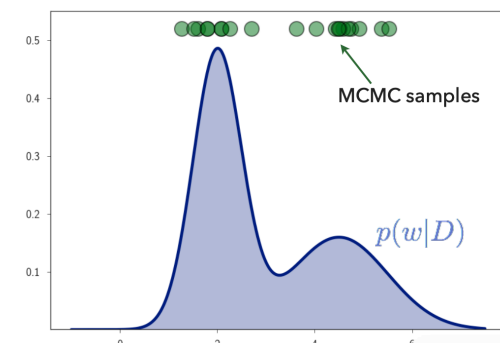
Laplace approx.

Simple yet less flexible



Variational inference

Efficient yet without the guarantee of asymptotic consistency

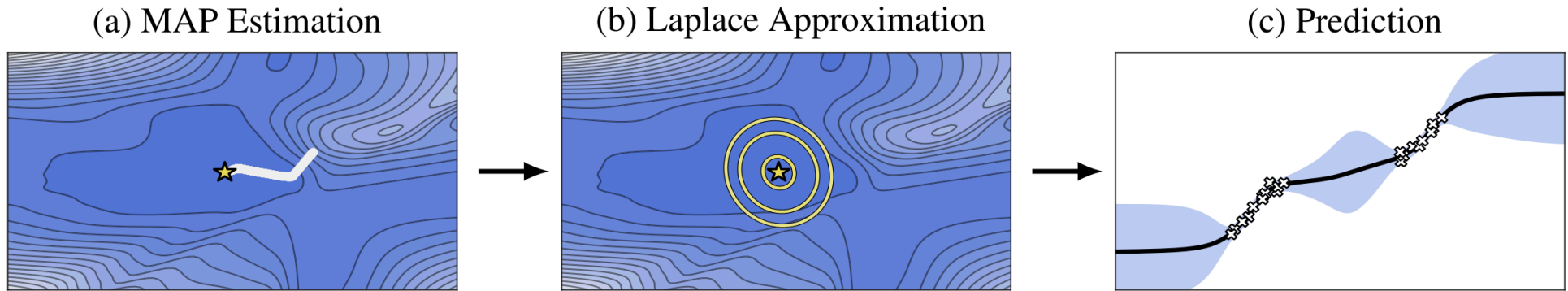


MCMC

Non-parametric and asymptotically exact yet typically with low convergence rate

Laplace approximation

Adapts a pre-trained DNN to BNN



[“Laplace Redux – Effortless Bayesian Deep Learning” by Daxberger et al.]

1. Find the maximum a posteriori (MAP) solution (the mode of Bayesian posterior)

$$\hat{\theta} = \arg \max_{\theta} \log p(\mathcal{D}|\theta) + \log p(\theta)$$

2. Construct a Gaussian approximation with the **Hessian** of log posterior

$$q(\theta) = \mathcal{N}(\theta; \hat{\theta}, \Sigma) \quad \Sigma^{-1} = -\nabla_{\theta\theta}^2 (\log p(\mathcal{D}|\theta) + \log p(\theta))|_{\theta=\hat{\theta}}$$

3. Make prediction by marginalization

The generalized Gauss-Newton (GGN) approximation and linearized LA (LLA)

Let $g_{\theta}(\mathbf{x})$ denote the model and set *an isotropic Gaussian prior*. GGN approx. sets

$$\Sigma^{-1} = \sum_i J_{\hat{\theta}}(\mathbf{x}_i)^{\top} \Lambda(\mathbf{x}_i, \mathbf{y}_i) J_{\hat{\theta}}(\mathbf{x}_i) + \mathbf{I}_P / \sigma_0^2,$$

where $J_{\hat{\theta}}(\mathbf{x}) \triangleq \nabla_{\theta} g_{\theta}(\mathbf{x})|_{\theta=\hat{\theta}}$ and $\Lambda(\mathbf{x}, \mathbf{y}) \triangleq -\nabla_{\mathbf{g}\mathbf{g}}^2 \log p(\mathbf{y}|\mathbf{g})|_{\mathbf{g}=g_{\hat{\theta}}(\mathbf{x})}$

- LLA applies LA to the first-order approximation of the NN of concern

$$g_{\theta}^{\text{lin}}(\mathbf{x}) = g_{\hat{\theta}}(\mathbf{x}) + J_{\hat{\theta}}(\mathbf{x})(\theta - \hat{\theta})$$

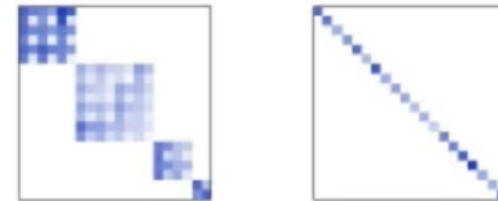
- LLA is **more sensible** than LA in the presence of GGN approximation [Immer et al., 22]
- LLA **can perform on par with or better than popular alternatives** on various uncertainty quantification (UQ) tasks
- The Laplace library [Daxberger et al., 21] further substantially advances LLA's applicability

Further approximations are required in practice!

- The GGN matrix of **size $P \times P$** is still unamenable in modern DL scenarios (P is the number of parameters)

- Further approximations **sparsifying the GGN** are introduced:

- **Diagonal and KFAC** [Martens & Grosse, 15] approximations



- Concern only the **last-layer inference**



- However, these strategies **sacrifice the fidelity of the learning outcomes** as the **approximation errors** in these cases can hardly be theoretically measured

Further approximations are required in practice!

- However, these strategies **sacrifice the fidelity of the learning outcomes** as the **approximation errors** in these cases can hardly be theoretically measured
- Our solution: accElerated Linearized Laplace Approximation (ELLA), which **scales LLA up** to make probabilistic predictions **in a more assurable way**

The inherent connections between Neural Tangent Kernels (NTKs) and LLA

- Integrating $q(\theta) = \mathcal{N}(\theta; \hat{\theta}, \Sigma)$ with the linear model $g_{\theta}^{\text{lin}}(\mathbf{x})$ actually gives rise to a **function-space** approximate posterior

$$\mathcal{GP}(f|g_{\hat{\theta}}(\mathbf{x}), \kappa_{\text{LLA}}(\mathbf{x}, \mathbf{x}')) \quad \text{with} \quad \kappa_{\text{LLA}}(\mathbf{x}, \mathbf{x}') \triangleq J_{\hat{\theta}}(\mathbf{x})\Sigma J_{\hat{\theta}}(\mathbf{x}')^{\top}$$

- Doing some simple math, we have

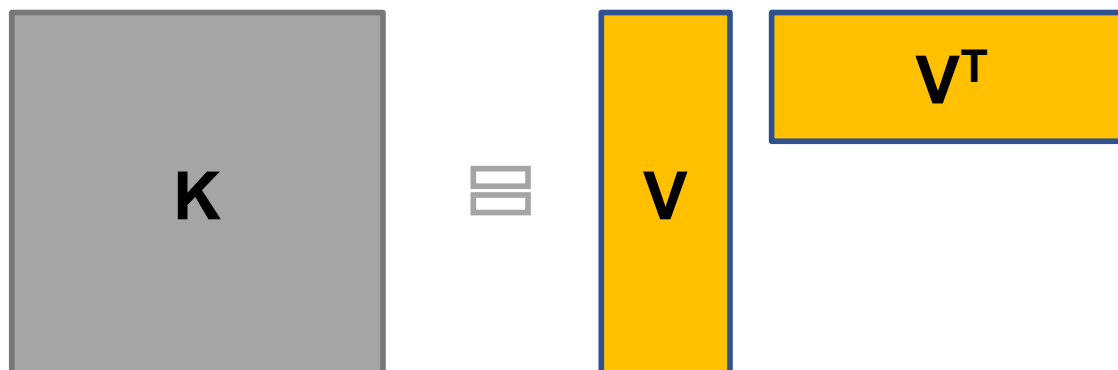
$$\kappa_{\text{LLA}}(\mathbf{x}, \mathbf{x}') = \sigma_0^2 \left(\kappa_{\text{NTK}}(\mathbf{x}, \mathbf{x}') - \kappa_{\text{NTK}}(\mathbf{x}, \mathbf{X}) [\Lambda_{\mathbf{X}, \mathbf{Y}}^{-1} / \sigma_0^2 + \kappa_{\text{NTK}}(\mathbf{X}, \mathbf{X})]^{-1} \kappa_{\text{NTK}}(\mathbf{X}, \mathbf{x}') \right)$$

$$\text{where } \kappa_{\text{NTK}}(\mathbf{x}, \mathbf{x}') \triangleq J_{\hat{\theta}}(\mathbf{x})J_{\hat{\theta}}(\mathbf{x}')^{\top}$$

- The **main challenge** then turns into the computation and inversion of the **gram matrix** $\kappa_{\text{NTK}}(\mathbf{X}, \mathbf{X})$ of size $NC \times NC$

Kernel approximation of NTK enables the acceleration of LLA

- If we can approximate $\kappa_{\text{NTK}}(\mathbf{x}, \mathbf{x}')$ with the inner product of some explicit K -dim representations of the data, i.e., $\kappa_{\text{NTK}}(\mathbf{x}, \mathbf{x}') \approx \varphi(\mathbf{x})\varphi(\mathbf{x}')^\top$,

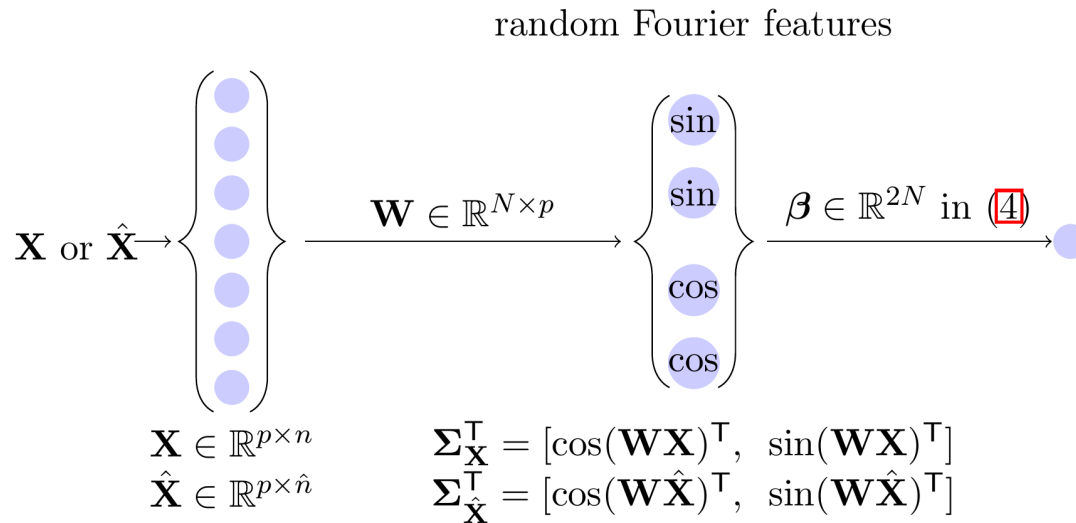


then

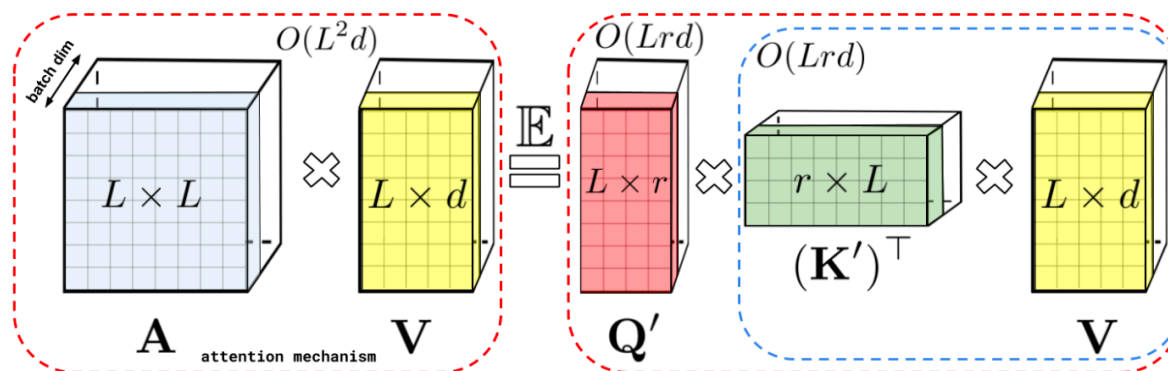
$$\begin{aligned} \kappa_{\text{LLA}}(\mathbf{x}, \mathbf{x}') &\approx \sigma_0^2 \left(\varphi(\mathbf{x})\varphi(\mathbf{x}')^\top - \varphi(\mathbf{x})\varphi_{\mathbf{X}}^\top \left[\Lambda_{\mathbf{X}, \mathbf{Y}}^{-1} / \sigma_0^2 + \varphi_{\mathbf{X}}\varphi_{\mathbf{X}}^\top \right]^{-1} \varphi_{\mathbf{X}}\varphi(\mathbf{x}')^\top \right) \\ &= \varphi(\mathbf{x}) \underbrace{\left[\sum_i \varphi(\mathbf{x}_i)^\top \Lambda(\mathbf{x}_i, \mathbf{y}_i) \varphi(\mathbf{x}_i) + \mathbf{I}_K / \sigma_0^2 \right]^{-1}}_{\mathbf{G}} \varphi(\mathbf{x}')^\top \triangleq \kappa_{\text{ELLA}}(\mathbf{x}, \mathbf{x}') \end{aligned}$$

Kernel approximation

Random features (RFs)



For shift-invariant kernels
[Rahimi et al., 2007]



Performer (RFs for $\exp(\mathbf{x}, \mathbf{x}')$)
[Choromanski et al., 2021]

Figure 1: Approximation of the regular attention mechanism \mathbf{AV} (before \mathbf{D}^{-1} -renormalization) via (random) feature maps. Dashed-blocks indicate order of computation with corresponding time complexities attached.

Kernel approximation

Nystrom method

- Given $\mathbf{X}_{\text{tr}} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ from q , perform MC integration:

$$\frac{1}{N} \sum_{n'=1}^N \kappa(\mathbf{x}, \mathbf{x}_{n'}) \psi_j(\mathbf{x}_{n'}) = \mu_j \psi_j(\mathbf{x}), \forall j \geq 1$$

- Eigendecompose $\kappa(\mathbf{X}_{\text{tr}}, \mathbf{X}_{\text{tr}})$ and get $\{(\hat{\mu}_j, [\hat{\psi}_j(\mathbf{x}_1), \dots, \hat{\psi}_j(\mathbf{x}_N)]^\top)\}_{j=1}^k$
- Kernelized solutions:

$$\hat{\psi}_j(\mathbf{x}) = \frac{1}{N \hat{\mu}_j} \sum_{n'=1}^N \kappa(\mathbf{x}, \mathbf{x}_{n'}) \hat{\psi}_j(\mathbf{x}_{n'}), \quad j = 1, \dots, k.$$

Adapt Nyström method to approximate the NTKs of multi-output NNs

- Rewrite $\kappa_{\text{NTK}}(\mathbf{x}, \mathbf{x}')$ as a scalar-valued kernel

$$\kappa_{\text{NTK}}((\mathbf{x}, i), (\mathbf{x}', i')) = J_{\hat{\theta}}(\mathbf{x}, i) J_{\hat{\theta}}(\mathbf{x}', i')^\top$$

By definition, the eigenfunctions can represent the spectral information of the kernel:

$$\int \kappa_{\text{NTK}}((\mathbf{x}, i), (\mathbf{x}', i')) \psi_k(\mathbf{x}', i') q(\mathbf{x}', i') = \mu_k \psi_k(\mathbf{x}, i), \forall k \geq 1,$$

while being orthonormal under q :

$$\int \psi_k(\mathbf{x}, i) \psi_{k'}(\mathbf{x}, i) q(\mathbf{x}, i) = \mathbb{1}[k = k'], \forall k, k' \geq 1.$$

By Monte Carlo estimation

$$\frac{1}{M} \sum_{m=1}^M \kappa_{\text{NTK}}((\mathbf{x}, i), (\mathbf{x}_m, i_m)) \psi_k(\mathbf{x}_m, i_m) = \mu_k \psi_k(\mathbf{x}, i), \forall k \in [K]$$

Adapt Nyström method to approximate the NTKs of multi-output NNs

- Applying this equation to the above samples gives rise to

$$\frac{1}{M} \sum_{m=1}^M \kappa_{\text{NTK}}((\mathbf{x}_{m'}, i_{m'}), (\mathbf{x}_m, i_m)) \psi_k(\mathbf{x}_m, i_m) = \mu_k \psi_k(\mathbf{x}_{m'}, i_{m'}), \forall k \in [K], m' \in [M]$$

then we arrive at

$$\frac{1}{M} \mathbf{K} \boldsymbol{\psi}_k \approx \mu_k \boldsymbol{\psi}_k, k \in [K]$$

- We compute the top-K eigenvalues $\lambda_1, \dots, \lambda_K$ of matrix \mathbf{K} and record the corresponding orthonormal eigenvectors $\mathbf{u}_1, \dots, \mathbf{u}_K$. Then it is easy to get the Nyström approximation of the top-K eigenfunctions:

$$\hat{\psi}_k(\mathbf{x}, i) = \frac{\sqrt{M}}{\lambda_k} \sum_{m=1}^M \mathbf{u}_k^{(m)} \kappa_{\text{NTK}}((\mathbf{x}, i), (\mathbf{x}_m, i_m)) = \frac{\sqrt{M}}{\lambda_k} J_{\hat{\theta}}(\mathbf{x}, i) \mathbf{J}_{\hat{\theta}, \tilde{\mathbf{X}}}^{\top} \mathbf{u}_k$$

- Finally, we have

$$\boldsymbol{\varphi}(\mathbf{x}) = [J_{\hat{\theta}}(\mathbf{x}) \mathbf{v}_1, \dots, J_{\hat{\theta}}(\mathbf{x}) \mathbf{v}_K] \text{ with } \mathbf{v}_k = \mathbf{J}_{\hat{\theta}, \tilde{\mathbf{X}}}^{\top} \mathbf{u}_k / \sqrt{\lambda_k}$$

Implementation

- The estimation of φ on a data point \mathbf{x} degenerates as K JVPs, which can be accomplished by invoking forward mode automatic differentiation (fwAD) for K times:

$$\begin{pmatrix} \mathbf{x}, \hat{\boldsymbol{\theta}} \\ \mathbf{v}_k \end{pmatrix} \xrightarrow{\text{forward}} \begin{pmatrix} g_{\hat{\boldsymbol{\theta}}}(\mathbf{x}) \\ J_{\hat{\boldsymbol{\theta}}}(\mathbf{x})\mathbf{v}_k \end{pmatrix}, k \in [K]$$

where the model output $g_{\hat{\boldsymbol{\theta}}}(\mathbf{x})$ and the JVP $J_{\hat{\boldsymbol{\theta}}}(\mathbf{x})\mathbf{v}_k$ are simultaneously computed in one single forward pass

- Prevalent DL libraries like PyTorch and Jax have already been armed with the capability for fwAD

Algorithms

Algorithm 1: Build the LLA posterior.

```
#  $g_{\hat{\theta}}$ : NN pre-trained by MAP;  $(\mathbf{X}, \mathbf{Y})$ :  
# training set;  $C$ : number of classes  
#  $M, K, \sigma_0^2$ : hyper-parameters  
def estimate_G( $\varphi, \mathbf{X}, \mathbf{Y}, K, \sigma_0^2$ ):  
     $\mathbf{G} = \text{zeros}(K, K)$   
    for  $(\mathbf{x}, \mathbf{y})$  in  $(\mathbf{X}, \mathbf{Y})$ :  
         $\mathbf{g}_x, \varphi_x = \varphi(\mathbf{x})$   
         $\Lambda_{x,y} = \text{hessian}(\text{nll}(\mathbf{g}_x, \mathbf{y}), \mathbf{g}_x)$   
         $\mathbf{G} += \varphi_x^\top \Lambda_{x,y} \varphi_x$   
    return  $\mathbf{G} + \text{eye}(K) / \sigma_0^2$   
def _q_f( $\varphi, \mathbf{G}^{-1}, \mathbf{x}$ )  
     $\mathbf{g}_x, \varphi_x = \varphi(\mathbf{x})$   
     $\kappa_{x,x} = \varphi_x \mathbf{G}^{-1} \varphi_x^\top$   
    return  $\mathcal{N}(\mathbf{g}_x, \kappa_{x,x})$   
 $\varphi = \text{build}_\varphi(g_{\hat{\theta}}, \mathbf{X}, C, M, K)$   
 $\mathbf{G}^{-1} = \text{inv}(\text{estimate}_\mathbf{G}(\varphi, \mathbf{X}, \mathbf{Y}, K, \sigma_0^2))$   
 $\text{q\_f} = \text{partial}(\_q\_f, \varphi, \mathbf{G}^{-1})$ 
```

Algorithm 2: Build φ .

```
def build_φ( $g_{\hat{\theta}}, \mathbf{X}, C, M, K$ ):  
    def _φ( $g_{\hat{\theta}}, C, \{\mathbf{v}_k\}_{k=1}^K, \mathbf{x}$ ):  
         $\varphi_x = \text{zeros}(C, K)$   
        for  $k$  in range( $K$ ):  
            with  $\text{fwAD.enable}()$ :  
                 $\mathbf{g}_x, \mathbf{jvp} = g_{(\hat{\theta}, \mathbf{v}_k)}(\mathbf{x})$   
                 $\varphi_x[:, k] = \mathbf{jvp}$   
        return  $\mathbf{g}_x, \varphi_x$   
     $\mathbf{J}_{\hat{\theta}, \tilde{\mathbf{x}}} = \text{zeros}(M, \text{dim}(\hat{\theta}))$   
    for  $m$  in range( $M$ ):  
         $\mathbf{x}_m = \text{uniform\_sample}(\mathbf{X})$   
         $i_m = \text{uniform\_sample}([C])$   
         $\mathbf{J}_{\hat{\theta}, \tilde{\mathbf{x}}}[m] = \text{grad}(g_{\hat{\theta}}(\mathbf{x}_m)[i_m], \hat{\theta})$   
     $\{\lambda_k, \mathbf{u}_k\}_{k=1}^K = \text{eig}(\mathbf{J}_{\hat{\theta}, \tilde{\mathbf{x}}} \mathbf{J}_{\hat{\theta}, \tilde{\mathbf{x}}}^\top, \text{top} = K)$   
    for  $k$  in range( $K$ ):  
         $\mathbf{v}_k = \mathbf{J}_{\hat{\theta}, \tilde{\mathbf{x}}}^\top \mathbf{u}_k / \sqrt{\lambda_k}$   
    return  $\text{partial}(\_φ, g_{\hat{\theta}}, C, \{\mathbf{v}_k\}_{k=1}^K)$ 
```

Theoretical Analysis

κ_{ELLA} can be reformulated as follows (in the seek of theoretical analysis)

$$\kappa_{\text{ELLA}}(\mathbf{x}, \mathbf{x}') = J_{\hat{\theta}}(\mathbf{x}) \underbrace{\mathbf{J}_{\hat{\theta}, \tilde{\mathbf{x}}}^{\top} \left[\mathbf{J}_{\hat{\theta}, \tilde{\mathbf{x}}} \mathbf{J}_{\hat{\theta}, \mathbf{x}}^{\top} \Lambda_{\mathbf{x}, \mathbf{y}} \mathbf{J}_{\hat{\theta}, \mathbf{x}} \mathbf{J}_{\hat{\theta}, \tilde{\mathbf{x}}}^{\top} + \mathbf{J}_{\hat{\theta}, \tilde{\mathbf{x}}} \mathbf{J}_{\hat{\theta}, \tilde{\mathbf{x}}}^{\top} / \sigma_0^2 \right]^{-1} \mathbf{J}_{\hat{\theta}, \tilde{\mathbf{x}}}^{\top}}_{\Sigma'} J_{\hat{\theta}}(\mathbf{x})^{\top}$$

Theorem 1 (Proof in Appendix A.4). *Let c_{Λ} be a finite constant associated with Λ , and \mathcal{E}' the error of Nyström approximation $\|\mathbf{J}_{\hat{\theta}, \mathbf{x}} \mathbf{J}_{\hat{\theta}, \tilde{\mathbf{x}}}^{\top} (\mathbf{J}_{\hat{\theta}, \tilde{\mathbf{x}}} \mathbf{J}_{\hat{\theta}, \tilde{\mathbf{x}}}^{\top})^{-1} \mathbf{J}_{\hat{\theta}, \tilde{\mathbf{x}}} \mathbf{J}_{\hat{\theta}, \mathbf{x}}^{\top} - \mathbf{J}_{\hat{\theta}, \mathbf{x}} \mathbf{J}_{\hat{\theta}, \mathbf{x}}^{\top}\|$. It holds that*

$$\mathcal{E} \leq \sigma_0^4 c_{\Lambda} \mathcal{E}' + \sigma_0^2.$$

Theorem 2 (Error bound of Nyström approximation). *With probability at least $1 - \delta$, it holds that:*

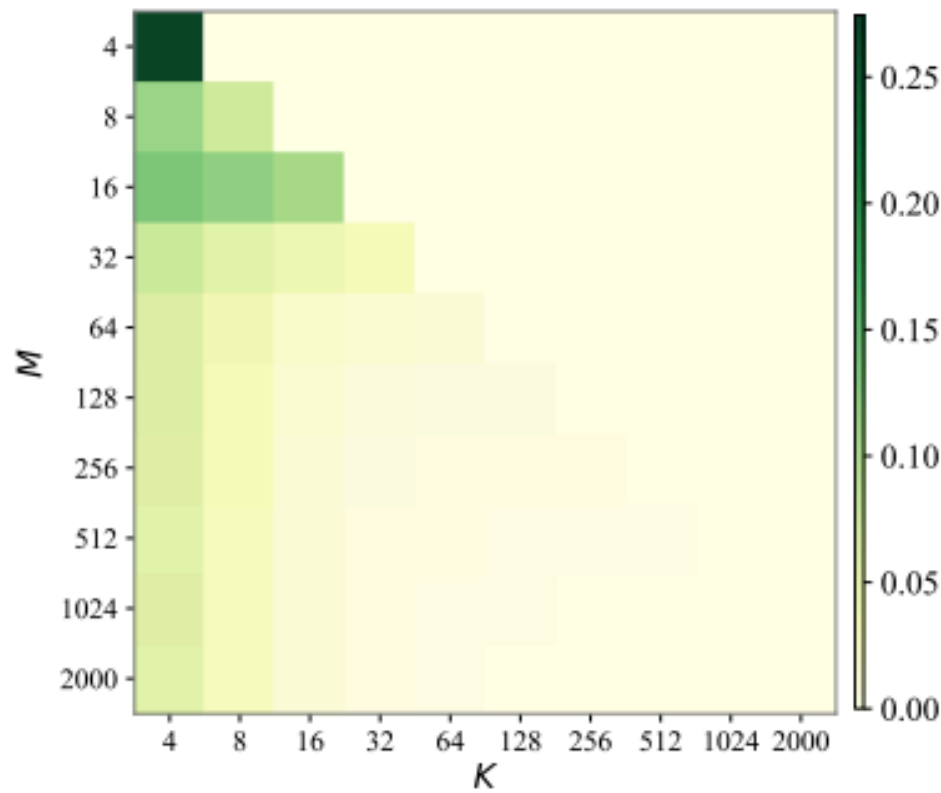
$$\mathcal{E}' \leq \tilde{\lambda}_{M+1} + \frac{NC}{\sqrt{M}} c_{\kappa} \left(2 + \log \frac{1}{\delta}\right).$$

Corollary 1. *With probability at least $1 - \delta$, the following bound exists:*

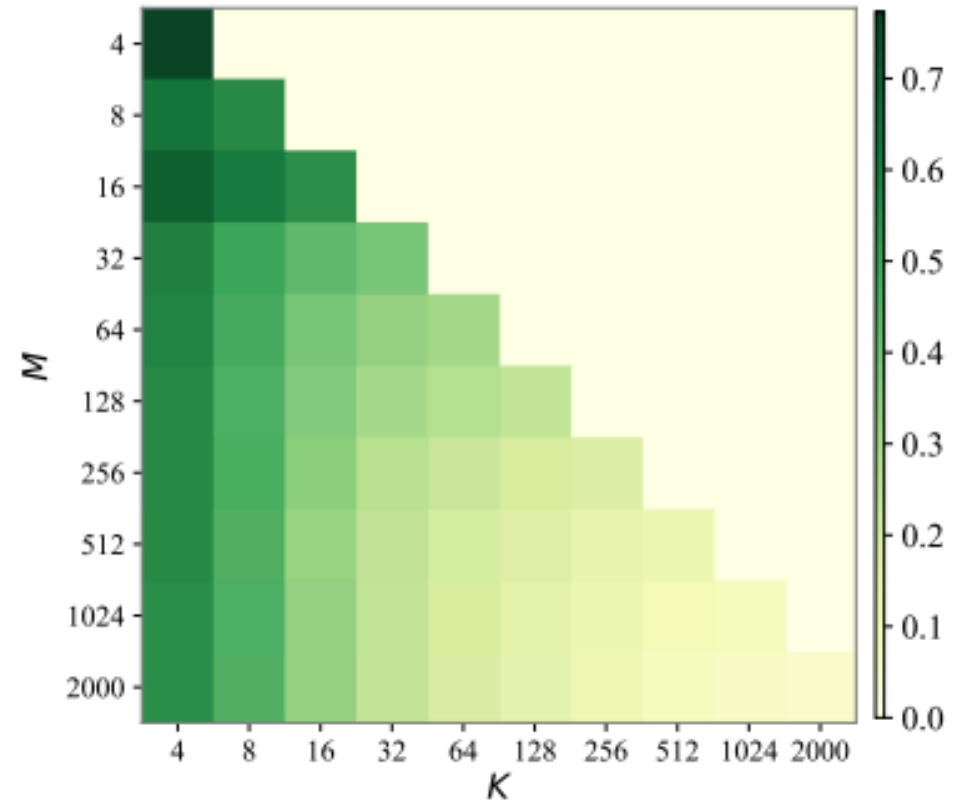
$$\mathcal{E} \leq \sigma_0^4 c_{\Lambda} \left(\tilde{\lambda}_{M+1} + \frac{NC}{\sqrt{M}} c_{\kappa} \left(2 + \log \frac{1}{\delta}\right) \right) + \sigma_0^2.$$

- As desired, the upper bound of approximation error decreases along with the growing of the number of MC samples in Nystrom approximation

How the approximation errors vary w.r.t. K and M? (MNIST, CNNs)



(a)



(b)

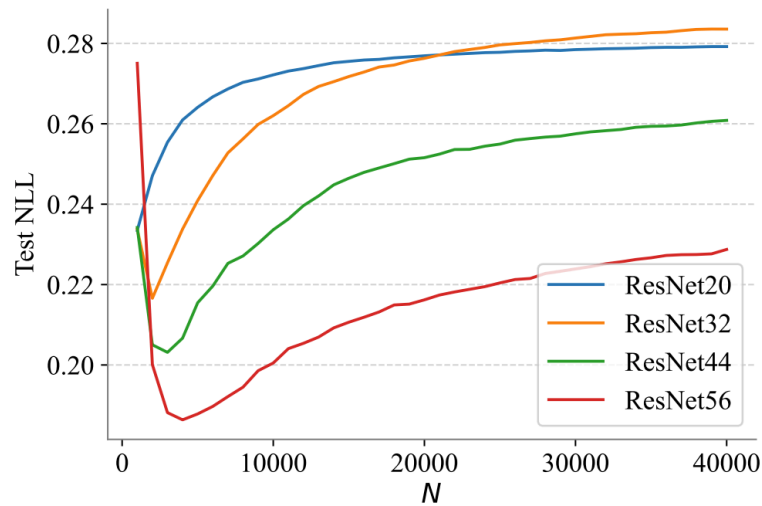
$$\epsilon_{\text{Nystrom}} \triangleq \|\mathbf{J}_{\hat{\theta}, \mathbf{X}} \mathbf{J}_{\hat{\theta}, \tilde{\mathbf{X}}}^{\top} (\mathbf{J}_{\hat{\theta}, \tilde{\mathbf{X}}} \mathbf{J}_{\hat{\theta}, \tilde{\mathbf{X}}}^{\top})^{-1} \mathbf{J}_{\hat{\theta}, \tilde{\mathbf{X}}} \mathbf{J}_{\hat{\theta}, \mathbf{X}}^{\top} - \mathbf{J}_{\hat{\theta}, \mathbf{X}} \mathbf{J}_{\hat{\theta}, \mathbf{X}}^{\top}\| / \|\mathbf{J}_{\hat{\theta}, \mathbf{X}} \mathbf{J}_{\hat{\theta}, \mathbf{X}}^{\top}\|$$

$$\epsilon_{\text{ELLA}} \triangleq \frac{1}{|\mathbf{X}_{\text{val}}|} \sum_{\mathbf{x} \in \mathbf{X}_{\text{val}}} \|\kappa_{\text{ELLA}}(\mathbf{x}, \mathbf{x}) - \kappa_{\text{LLA}}(\mathbf{x}, \mathbf{x})\| / \|\kappa_{\text{LLA}}(\mathbf{x}, \mathbf{x})\|$$

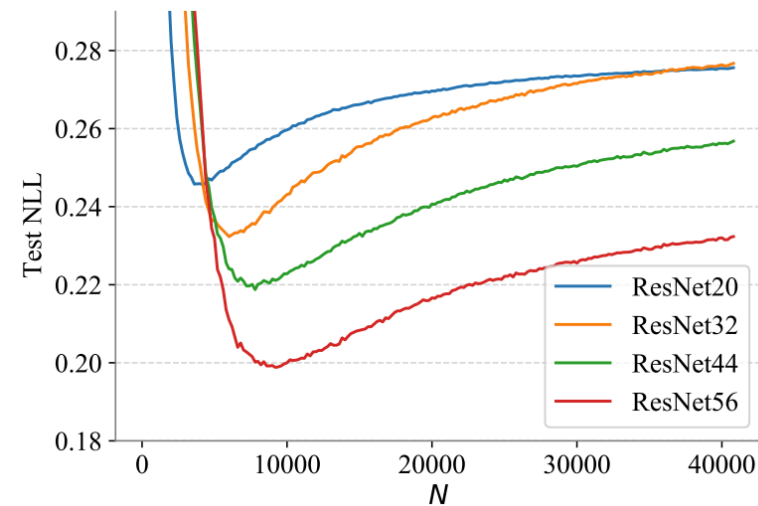
The overfitting issue of ELLA and general LLA (Cifar-10, ResNets)

$$\Sigma^{-1} = \sum_i J_{\hat{\theta}}(\mathbf{x}_i)^\top \Lambda(\mathbf{x}_i, \mathbf{y}_i) J_{\hat{\theta}}(\mathbf{x}_i) + \mathbf{I}_P / \sigma_0^2$$

With more training data involved, the covariance in LA, LLA, and ELLA shrinks and **the uncertainty dissipates**.



ELLA



Last-layer LLA

Solution: early stopping

Illustrative Regression

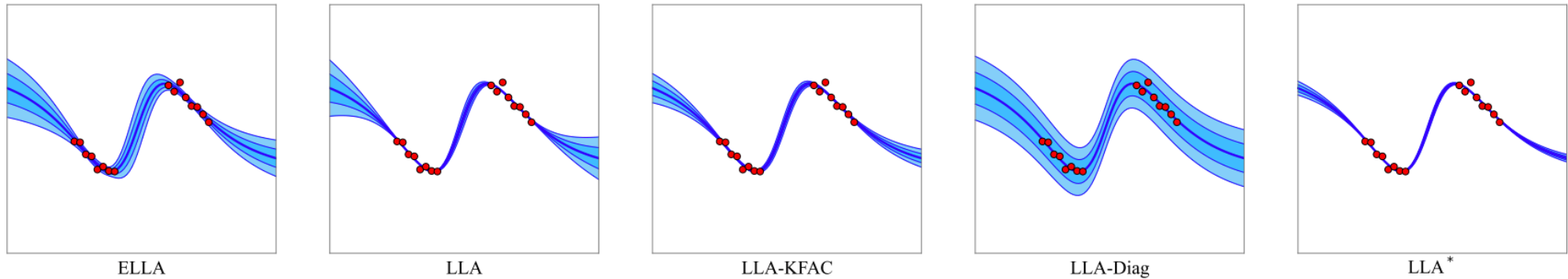


Figure 1: 1-D regression on $y = \sin 2x + \epsilon$, $\epsilon \sim \mathcal{N}(0, 0.2)$. Red dots, central blue curves, and shaded regions refer to the training data, mean predictions, and uncertainty respectively. The model is a pretrained multilayer perceptron (MLP) with 3 hidden layers. As shown, the predictive uncertainty of ELLA is on par with or better than the competitors such as LLA with KFAC approximation (LLA-KFAC), LLA with diagonal approximation (LLA-Diag), and last-layer LLA (LLA*).

Cifar-10 classification

Table 1: Comparison on test accuracy (%) \uparrow , NLL \downarrow , and ECE \downarrow on CIFAR-10. We report the average results over 5 random runs. As the accuracy values of most methods are close, we do not highlight the best.

Method	ResNet-20			ResNet-32			ResNet-44			ResNet-56		
	Acc.	NLL	ECE	Acc.	NLL	ECE	Acc.	NLL	ECE	Acc.	NLL	ECE
<i>ELLA</i>	92.5	0.233	0.009	93.5	0.215	0.008	93.9	0.204	0.007	94.4	0.187	0.007
<i>MAP</i>	92.6	0.282	0.039	93.5	0.292	0.041	94.0	0.275	0.039	94.4	0.252	0.037
<i>MFVI-BF</i>	92.7	0.231	0.016	93.5	0.222	0.020	93.9	0.206	0.018	94.4	0.188	0.016
<i>LLA*</i>	92.6	0.269	0.034	93.5	0.259	0.033	94.0	0.237	0.028	94.4	0.213	0.022
<i>LLA*-KFAC</i>	92.6	0.271	0.035	93.5	0.260	0.033	94.0	0.232	0.028	94.4	0.202	0.024
<i>LLA-Diag</i>	92.2	0.728	0.404	92.7	0.755	0.430	92.8	0.778	0.445	92.9	0.843	0.480
<i>LLA-KFAC</i>	92.0	0.852	0.467	91.8	1.027	0.547	91.4	1.091	0.566	89.8	1.174	0.579

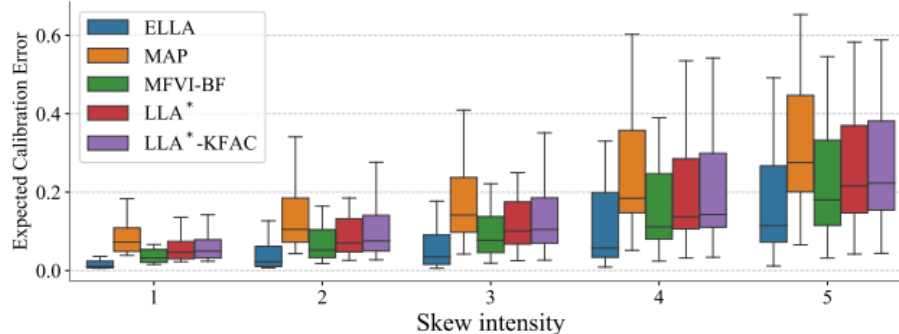
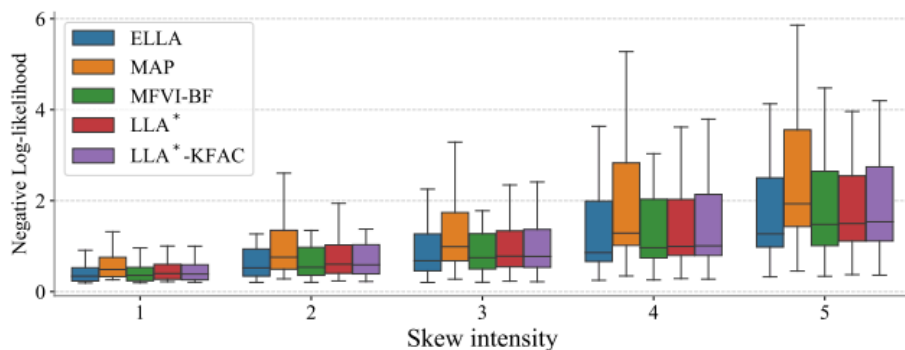


Figure 3: NLL (Left) and ECE (Right) on CIFAR-10 corruptions for models trained with ResNet-56 architecture. Each box corresponds to a summary of the results across 19 types of skew.

ImageNet classification

Table 2: Comparison on test accuracy (%) \uparrow , NLL \downarrow , and ECE \downarrow on ImageNet. We report the average results over 3 random runs.

Method	ResNet-18			ResNet-34			ResNet-50		
	Acc.	NLL	ECE	Acc.	NLL	ECE	Acc.	NLL	ECE
<i>ELLA</i>	69.8	1.243	0.015	73.3	1.072	0.018	76.2	0.948	0.018
<i>MAP</i>	69.8	1.247	0.026	73.3	1.081	0.035	76.2	0.962	0.037
<i>MFVI-BF</i>	70.3	1.218	0.042	73.7	1.043	0.033	76.1	0.945	0.030

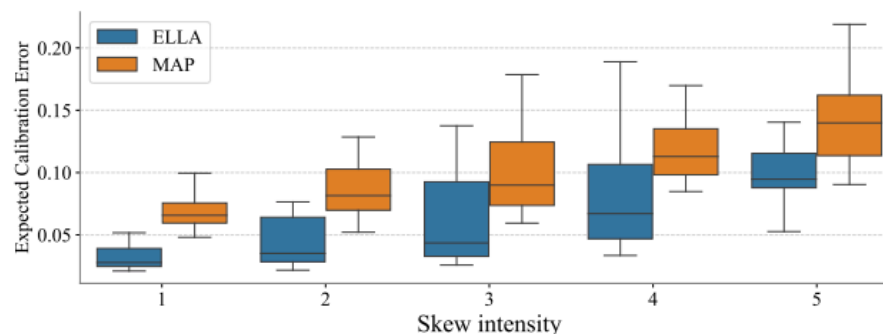
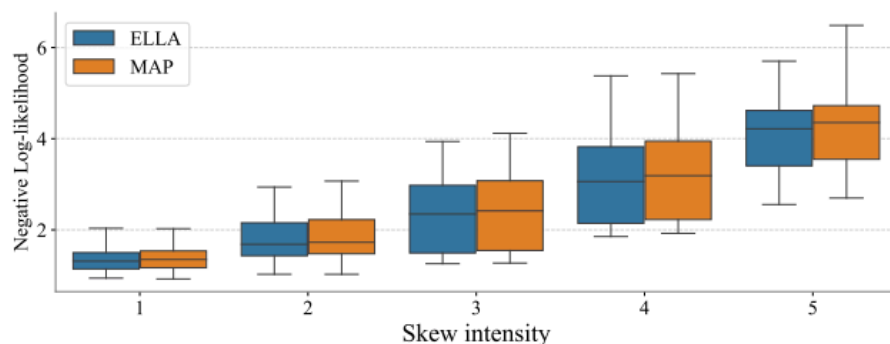


Figure 5: NLL (Left) and ECE (Right) on ImageNet corruptions for models trained with ViT-B architecture. Each box corresponds to a summary of the results across 19 types of skew.

Results on ViT

Method	Acc.	NLL	ECE
<i>ELLA</i>	81.6	0.695	0.022
<i>MAP</i>	81.5	0.700	0.039



Thanks!

