

# NeuralEF: Deconstructing Kernels by Deep Neural Networks

Zhijie Deng

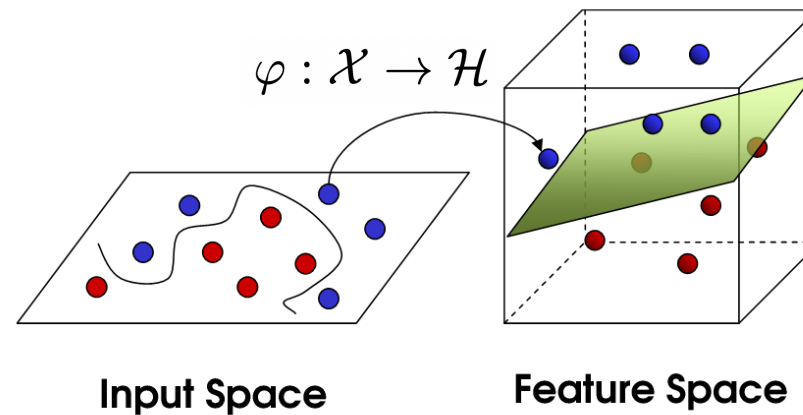
Tsinghua University

Contact: dengzhijiethu@gmail.com

*Joint work with: Jiaxin Shi and Jun Zhu*



# Kernel methods



- $\kappa(\mathbf{x}, \mathbf{x}') = \langle \varphi(\mathbf{x}), \varphi(\mathbf{x}') \rangle_{\mathcal{H}}$
- **Pros:** non-parametric flexibility & analytical inference
- **Cons:** limited scalability – at least  $O(N^2)$  complexity, typically  $O(N^3)$ ; inefficiency issue in the test phase

# Kernel approximation

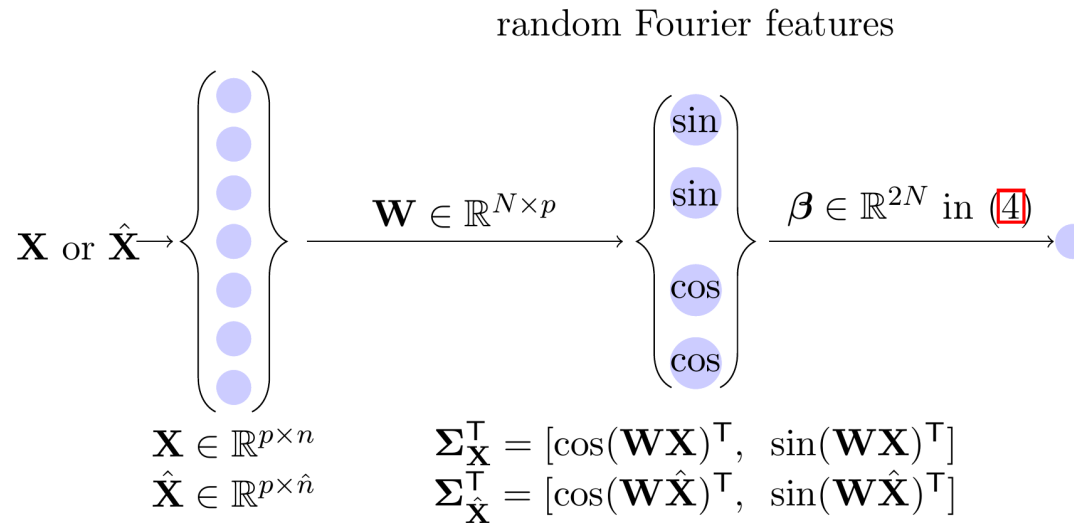
- Approximate the kernel with the inner product of **some explicit vector representations of the data**:

$$\kappa(\mathbf{x}, \mathbf{x}') \approx \nu(\mathbf{x})^\top \nu(\mathbf{x}') \quad \nu : \mathcal{X} \rightarrow \mathbb{R}^k$$

- A small  $k$  is desired for **scalability** while the approximation is **low-rank**
- Popular approaches:
  1. Random Fourier features [Rahimi & Recht, 2007; 2008]
  2. Nystrom method [Nystrom, 1930; Williams & Seeger, 2001]
  3. ....

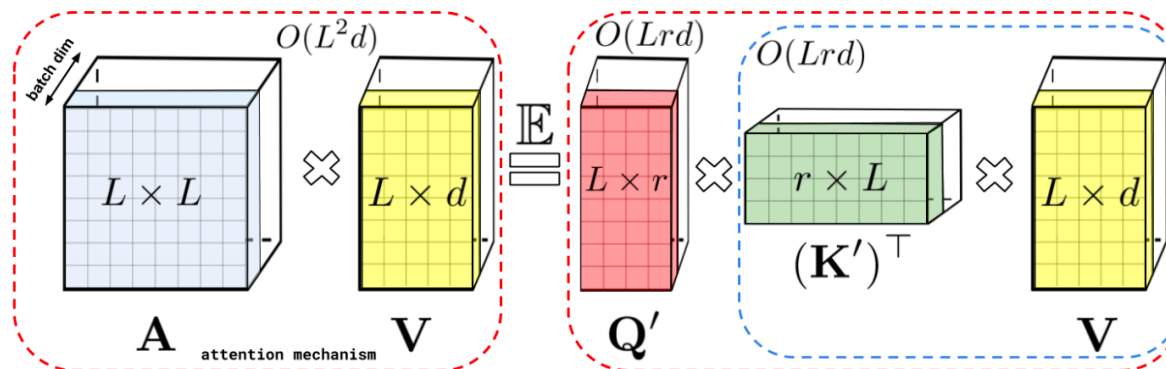
# Kernel approximation

## 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

- Mercer's theorem

$$\kappa(\mathbf{x}, \mathbf{x}') = \sum_{j \geq 1} \mu_j \psi_j(\mathbf{x}) \psi_j(\mathbf{x}')$$

where  $\psi_j$  denote the **eigenfunctions** of the kernel  $\kappa$  w.r.t. the probability measure  $q$ , and  $\mu_j \geq 0$  refer to the corresponding **eigenvalues**

- By the definition of eigenfunction, we have

$$\int \kappa(\mathbf{x}, \mathbf{x}') \psi_j(\mathbf{x}') q(\mathbf{x}') d\mathbf{x}' = \mu_j \psi_j(\mathbf{x}), \quad \forall j \geq 1$$

and

$$\int \psi_i(\mathbf{x}) \psi_j(\mathbf{x}) q(\mathbf{x}) d\mathbf{x} = \mathbb{1}[i = j], \quad \forall i, j \geq 1$$

# 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:

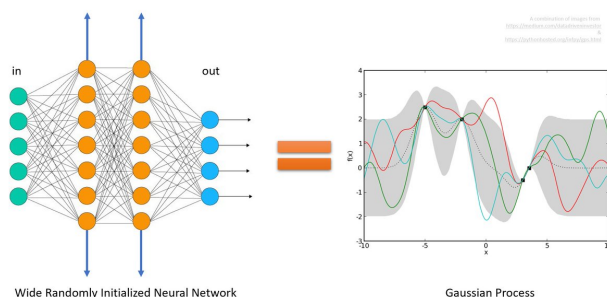
$$\hat{\psi}_j(\mathbf{x}) = \frac{1}{N \hat{\mu}_j} \sum_{n'=1}^N \kappa(\mathbf{x}, \mathbf{x}_{n'}) \hat{\psi}_j(\mathbf{x}_{n'}), \quad j = 1, \dots, k.$$

- Less scalable; the testing entails extensive computes

# The modern kernels

## Kernels meet NNs

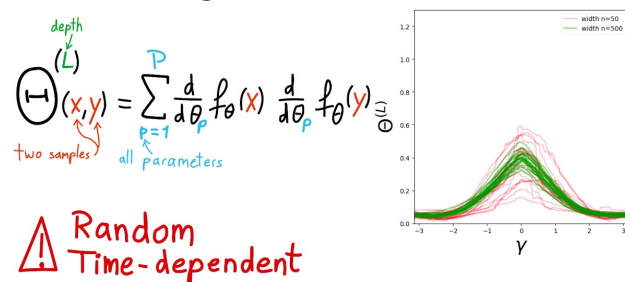
- Classic local kernels suffer from **curse of dimensionality** [Bengio et al., 2005]
- *Neural network Gaussian process (NNGP) kernels* [Neal, 1995; Lee et al., 2017; Khan et al., 2019]



$$\kappa_{\text{NN-GP}}(\mathbf{x}, \mathbf{x}') = \mathbb{E}_{\boldsymbol{\theta} \sim p(\boldsymbol{\theta})} g(\mathbf{x}, \boldsymbol{\theta}) g(\mathbf{x}', \boldsymbol{\theta})^\top$$

- *Neural tangent kernels (NTKs)* [Jacot et al., 2018]

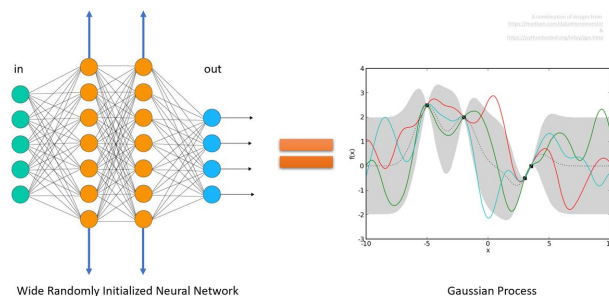
Neural Tangent Kernel



$$\kappa_{\text{NTK}}(\mathbf{x}, \mathbf{x}') = \partial_{\boldsymbol{\theta}} g(\mathbf{x}, \boldsymbol{\theta}) \partial_{\boldsymbol{\theta}} g(\mathbf{x}', \boldsymbol{\theta})^\top$$

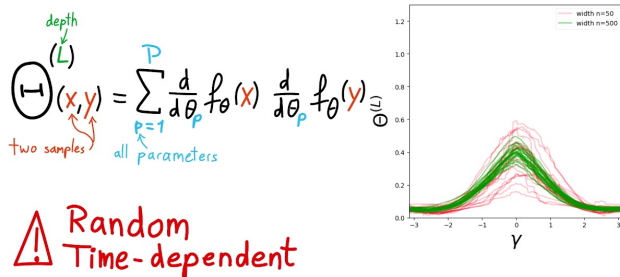
# The modern kernels

## Kernels meet NNs



$$\kappa_{\text{NN-GP}}(\mathbf{x}, \mathbf{x}') = \mathbb{E}_{\boldsymbol{\theta} \sim p(\boldsymbol{\theta})} g(\mathbf{x}, \boldsymbol{\theta}) g(\mathbf{x}', \boldsymbol{\theta})^\top$$

Neural Tangent Kernel

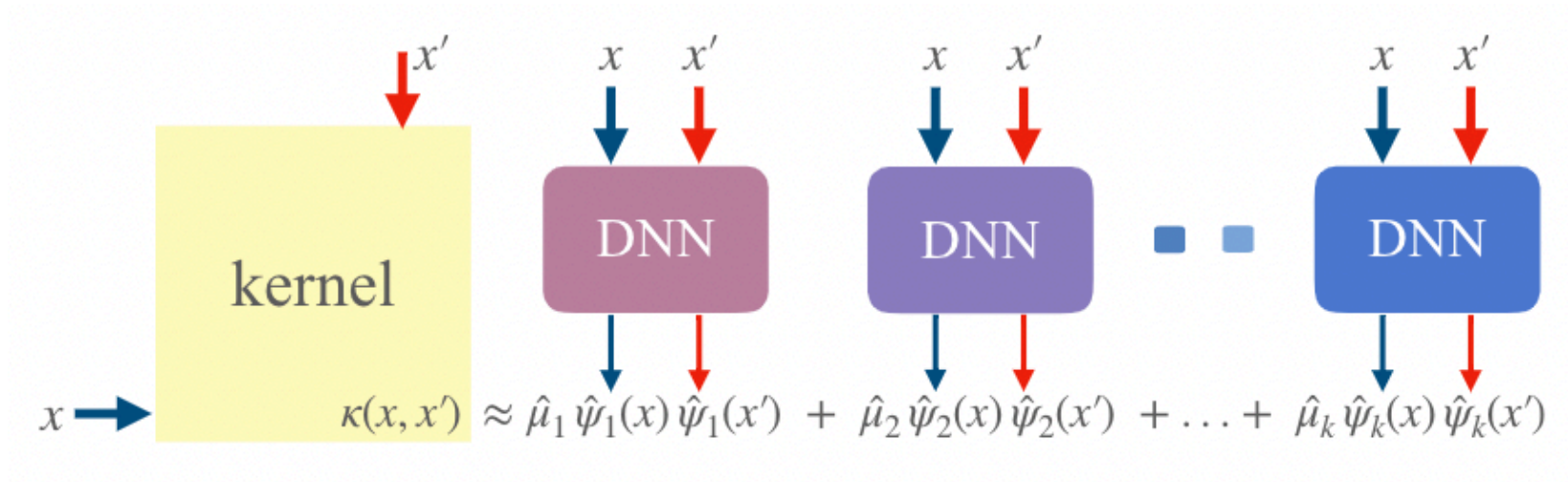


$$\kappa_{\text{NTK}}(\mathbf{x}, \mathbf{x}') = \partial_{\boldsymbol{\theta}} g(\mathbf{x}, \boldsymbol{\theta}) \partial_{\boldsymbol{\theta}} g(\mathbf{x}', \boldsymbol{\theta})^\top$$

- Nevertheless, writing down their **detailed mathematical formulae** is non-trivial [Arora et al., 2019] and evaluating them with recursion is both **time and memory consuming**.
- They have **poor compatibility** with standard kernel approximation methods.

NeuralEF: approximate the **eigenfunctions** of kernels by NNs

Our solution



## A closely related work

### Spectral Inference Networks (SpIN) [Pfau et al., 2018]

- Recover the top eigenfunctions with NNs due to their **universal approximation capability** and **parametric nature**
- Introduce a vector-valued NN function  $\Psi(\cdot, \mathbf{w}) : \mathcal{X} \rightarrow \mathbb{R}^k$  and solve:

$$\begin{aligned} \max_{\mathbf{w}} \text{Tr} \left( \iint \Psi(\mathbf{x}, \mathbf{w}) \Psi(\mathbf{x}', \mathbf{w})^\top \kappa(\mathbf{x}, \mathbf{x}') q(\mathbf{x}) q(\mathbf{x}') d\mathbf{x} d\mathbf{x}' \right) \\ \text{s.t.: } \int \Psi(\mathbf{x}, \mathbf{w}) \Psi(\mathbf{x}, \mathbf{w})^\top q(\mathbf{x}) d\mathbf{x} = \mathbf{I}_k, \end{aligned} \quad (9)$$

- However, this objective makes  $\Psi$  recover the **subspace spanned by the top-k eigenfunctions** rather than the **top-k eigenfunctions themselves** [Pfau et al., 2018].

## A closely related work

### Spectral Inference Networks (SpIN) [Pfau et al., 2018]

- Recover the top eigenfunctions with NNs due to their **universal approximation capability** and **parametric nature**
- Introduce a vector-valued NN function  $\Psi(\cdot, \mathbf{w}) : \mathcal{X} \rightarrow \mathbb{R}^k$  and solve:

$$\begin{aligned} \max_{\mathbf{w}} \text{Tr} \left( \iint \Psi(\mathbf{x}, \mathbf{w}) \Psi(\mathbf{x}', \mathbf{w})^\top \kappa(\mathbf{x}, \mathbf{x}') q(\mathbf{x}) q(\mathbf{x}') d\mathbf{x} d\mathbf{x}' \right) \\ \text{s.t.: } \int \Psi(\mathbf{x}, \mathbf{w}) \Psi(\mathbf{x}, \mathbf{w})^\top q(\mathbf{x}) d\mathbf{x} = \mathbf{I}_k, \end{aligned} \quad (9)$$

- To address this issue, SpIN relies on a gradient masking trick which involves a **Cholesky decomposition** per training iteration.
- SpIN also involves tracking the **exponential moving average (EMA)** of the **Jacobian** matrix to debias the stochastic optimization.



# Eigendecomposition as **asymmetric** maximization problems

## Our new results

Generalized  
Rayleigh quotient

Normalization  
constraint

**Theorem 1** (Proof in Appendix A.1). *The eigenpairs of the kernel  $\kappa(\mathbf{x}, \mathbf{x}')$  can be recovered by simultaneously solving the following series of constrained maximization problems:*

$$\max_{\hat{\psi}_j} R_{jj} \text{ s.t.: } C_j = 1, R_{1j} = 0, \dots, R_{(j-1)j} = 0, \forall j \geq 1, \quad (7)$$

where  $\hat{\psi}_j \in L^2(\mathcal{X}, q)$  represent the introduced approximate eigenfunctions, and

$$R_{ij} := \iint \hat{\psi}_i(\mathbf{x}) \kappa(\mathbf{x}, \mathbf{x}') \hat{\psi}_j(\mathbf{x}') q(\mathbf{x}') q(\mathbf{x}) d\mathbf{x}' d\mathbf{x}, \quad (8)$$

$$C_j := \int \hat{\psi}_j(\mathbf{x}) \hat{\psi}_j(\mathbf{x}) q(\mathbf{x}) d\mathbf{x}. \quad (9)$$

In particular,  $(R_{jj}, \hat{\psi}_j)$  will converge to the eigenpair associated with  $j$ -th largest eigenvalue of  $\kappa$ .

Orthogonality  
constraint





# Eigendecomposition as **asymmetric** maximization problems

## Proof scratch--the first problem

- The ground-truth eigenfunctions form **a set of orthonormal bases** of the  $L_2(X, q)$  space

- Represent the approximations in such a new axis system  $\hat{\psi}_1 = \sum_{i \geq 1} w_i \psi_i$

- The the maximization objective reduces to

$$R_{11} = \sum_{j \geq 1} \mu_j \langle \hat{\psi}_1, \psi_j \rangle^2 = \sum_{j \geq 1} \mu_j \langle \sum_{i \geq 1} w_i \psi_i, \psi_j \rangle^2 = \sum_{j \geq 1} \mu_j w_j^2.$$

- And the constraint reduces to

$$\langle \hat{\psi}_1, \hat{\psi}_1 \rangle = \langle \sum_{i \geq 1} w_i \psi_i, \sum_{j \geq 1} w_j \psi_j \rangle = \sum_{i, j \geq 1} w_i w_j \langle \psi_i, \psi_j \rangle = \sum_{j \geq 1} w_j^2 = 1$$

- It is straight-forward to see the maxima

# Eigendecomposition as **asymmetric** maximization problems

## Proof scratch--the second problem

- Given  $\hat{\psi}_1 = \psi_1$

$$R_{12} = 0$$

$$\Rightarrow \iint \hat{\psi}_1(\mathbf{x}) \kappa(\mathbf{x}, \mathbf{x}') \hat{\psi}_2(\mathbf{x}') q(\mathbf{x}') q(\mathbf{x}) d\mathbf{x}' d\mathbf{x} = 0$$

$$\Rightarrow \int \hat{\psi}_2(\mathbf{x}') q(\mathbf{x}') \int \hat{\psi}_1(\mathbf{x}) \kappa(\mathbf{x}, \mathbf{x}') q(\mathbf{x}) d\mathbf{x} d\mathbf{x}' = 0$$

$$\Rightarrow \int \hat{\psi}_2(\mathbf{x}') q(\mathbf{x}') \int \psi_1(\mathbf{x}) \kappa(\mathbf{x}, \mathbf{x}') q(\mathbf{x}) d\mathbf{x} d\mathbf{x}' = 0$$

$$\Rightarrow \int \hat{\psi}_2(\mathbf{x}') q(\mathbf{x}') \mu_1 \psi_1(\mathbf{x}') d\mathbf{x}' = 0$$

$$\Rightarrow \langle \psi_1, \hat{\psi}_2 \rangle = 0.$$

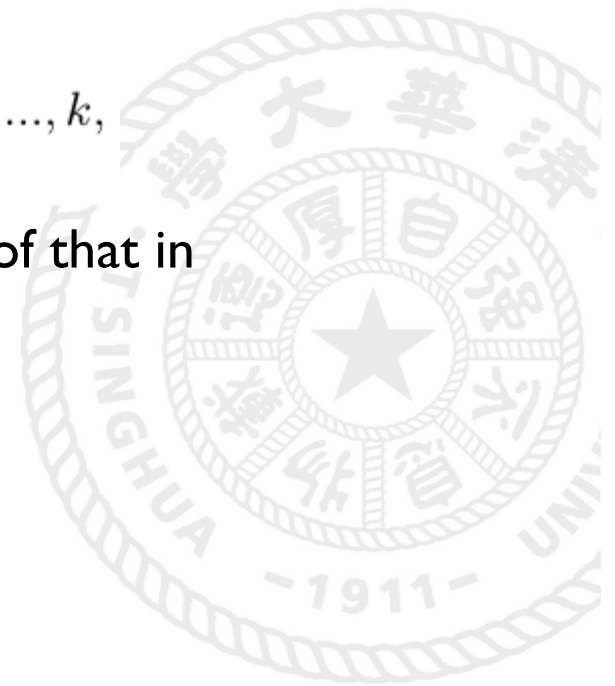
- $\hat{\psi}_2$  is constrained in the **orthogonal complement** of the subspace spanned by  $\psi_1$
- Then we can apply an analysis similar to that for the first problem
- Applying this procedure incrementally to the additional problems then finishes the proof

# Eigendecomposition as **asymmetric** maximization problems

- **Slack the constraints** on orthogonality as penalties and solve **the first k** optimization problems

$$\max_{\hat{\psi}_j} R_{jj} - \sum_{i=1}^{j-1} \frac{R_{ij}^2}{R_{ii}} \quad \text{s.t.: } C_j = 1, \text{ for } j = 1, \dots, k,$$

- Our objective forms a **function-space generalization** of that in EigenGame [Gemp et al., 2020]



# DNNs as eigenfunctions

Use an ensemble of  $k$  DNNs to approximate the top- $k$  eigenfunctions

- Mini-batch training -- by MC estimation:

$$\begin{aligned}\tilde{R}_{ij} &= \sum_{b=1}^B \sum_{b'=1}^B \frac{1}{B^2} \hat{\psi}_i(\mathbf{x}_b) \kappa(\mathbf{x}_b, \mathbf{x}_{b'}) \hat{\psi}_j(\mathbf{x}_{b'}) \\ &= \frac{1}{B^2} \hat{\psi}_i^{\mathbf{X}\top} \boldsymbol{\kappa}^{\mathbf{X},\mathbf{X}} \hat{\psi}_j^{\mathbf{X}},\end{aligned}$$

- L2 Batch normalization (L2BN) to absorb the normalization constraints:

$$h_b^{\text{out}} = \frac{h_b^{\text{in}}}{\sigma}, \text{ with } \sigma = \sqrt{\frac{1}{B} \sum_{b=1}^B h_b^{\text{in}2}}, \quad b = 1, \dots, B.$$

- The gradients:  $\nabla_{\mathbf{w}_j} \ell = -\frac{2}{B^2} \boldsymbol{\kappa}^{\mathbf{X},\mathbf{X}} \left( \hat{\psi}_j^{\mathbf{X}} - \sum_{i=1}^{j-1} \frac{\hat{\psi}_i^{\mathbf{X}\top} \boldsymbol{\kappa}^{\mathbf{X},\mathbf{X}} \hat{\psi}_j^{\mathbf{X}}}{\hat{\psi}_i^{\mathbf{X}\top} \boldsymbol{\kappa}^{\mathbf{X},\mathbf{X}} \hat{\psi}_i^{\mathbf{X}}} \hat{\psi}_i^{\mathbf{X}} \right) \cdot \partial_{\mathbf{w}_j} \hat{\psi}_j^{\mathbf{X}}.$

- Extension to *matrix-valued* kernels (e.g., NTKs):  
strategy 1: use multi-output DNNs  
strategy 2: make a factorization assumption

# NeuralEF

## The algorithm

---

**Algorithm 1** Find the top- $k$  eigenpairs of a kernel by NeuralEF

---

- 1: **Input:** Training data  $\mathbf{X}_{\text{tr}}$ , kernel  $\kappa$ , batch size  $B$ .
  - 2: Initialize NNs  $\hat{\psi}_j(\cdot) = \hat{\psi}(\cdot, \mathbf{w}_j)$  and scalars  $\hat{\mu}_j, j \in [k]$ ;
  - 3: Compute the kernel matrix  $\kappa^{\mathbf{X}_{\text{tr}}, \mathbf{X}_{\text{tr}}} = \kappa(\mathbf{X}_{\text{tr}}, \mathbf{X}_{\text{tr}})$ ;
  - 4: **for iteration do**
  - 5:   Draw a mini-batch  $\mathbf{X} \subseteq \mathbf{X}_{\text{tr}}$ ; retrieve  $\kappa^{\mathbf{X}, \mathbf{X}}$  from  $\kappa^{\mathbf{X}_{\text{tr}}, \mathbf{X}_{\text{tr}}}$ ;
  - 6:   Do forward propagation  $\hat{\psi}_j^{\mathbf{X}} = \hat{\psi}(\mathbf{X}, \mathbf{w}_j), j \in [k]$ ;
  - 7:    $\hat{\mu}_j \leftarrow \text{EMA}(\hat{\mu}_j, \frac{1}{B^2} \hat{\psi}_j^{\mathbf{X}^\top} \kappa^{\mathbf{X}, \mathbf{X}} \hat{\psi}_j^{\mathbf{X}}), j \in [k]$ ;
  - 8:   Compute  $\nabla_{\mathbf{w}_j} \ell, j \in [k]$  by Equation (18) and do SGD;
  - 9: **end for**
-

# Enable the learning of NN-GP kernels and NTKs

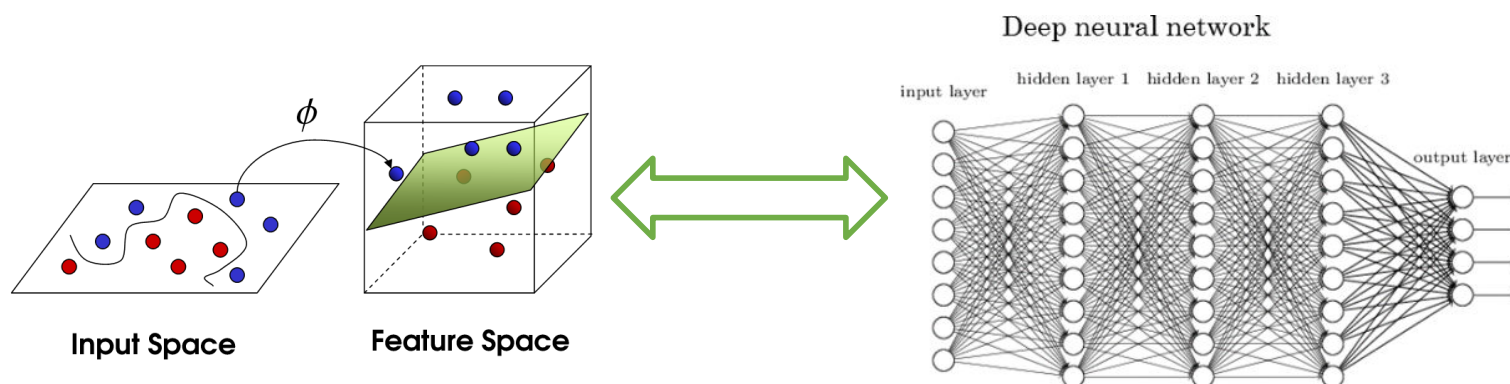
Based on thousands of random features

$$\begin{aligned}
 \begin{array}{c} \updownarrow n \\ \left[ \begin{array}{c} \mathbf{H}(\mathbf{w}_0) \end{array} \right] \\ \leftarrow n \end{array} &= \begin{array}{c} \left[ \begin{array}{c} \nabla_w \mathbf{y}(\mathbf{w}_0)^T \end{array} \right] \\ \leftarrow p \end{array} \begin{array}{c} \left[ \begin{array}{c} \nabla_w \mathbf{y}(\mathbf{w}_0) \end{array} \right] \\ \leftarrow n \end{array} \quad \begin{array}{c} \updownarrow p \\ \end{array} \\
 \text{NTK} & \\
 &= \begin{array}{c} \left[ \begin{array}{c} \text{---} \phi(\bar{\mathbf{x}}_1)^T \text{---} \\ \vdots \\ \text{---} \phi(\bar{\mathbf{x}}_n)^T \text{---} \end{array} \right] \\ \left[ \begin{array}{c} \phi(\bar{\mathbf{x}}_1) \dots \phi(\bar{\mathbf{x}}_n) \end{array} \right] \end{array}
 \end{aligned}$$

- Computing the training kernel matrices by **MC estimation** given a distribution  $p(v)$  satisfying  $\mathbb{E}_{p(v)}[vv^\top] = \mathbf{I}_{\dim(\theta)}$ , then

$$\begin{aligned}\kappa_{\text{NTK}}^{\mathbf{X}_{\text{tr}}, \mathbf{X}_{\text{tr}}} &= \mathbb{E}_{\mathbf{v} \sim p(\mathbf{v})} [\partial_{\theta} g(\mathbf{X}_{\text{tr}}, \theta) \mathbf{v}] [\partial_{\theta} g(\mathbf{X}_{\text{tr}}, \theta) \mathbf{v}]^{\top} \\ &\approx \frac{1}{S} \sum_{s=1}^S \left[ \frac{g(\mathbf{X}_{\text{tr}}, \theta + \epsilon \mathbf{v}_s) - g(\mathbf{X}_{\text{tr}}, \theta)}{\epsilon} \right] \left[ \frac{g(\mathbf{X}_{\text{tr}}, \theta + \epsilon \mathbf{v}_s) - g(\mathbf{X}_{\text{tr}}, \theta)}{\epsilon} \right]^{\top}\end{aligned}$$

# The impact of NeuralEF

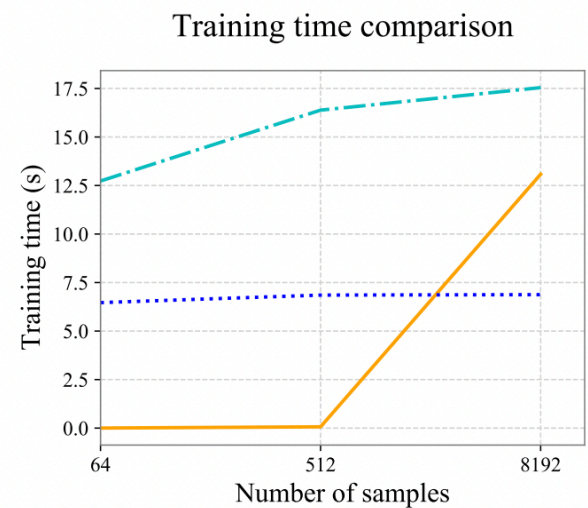
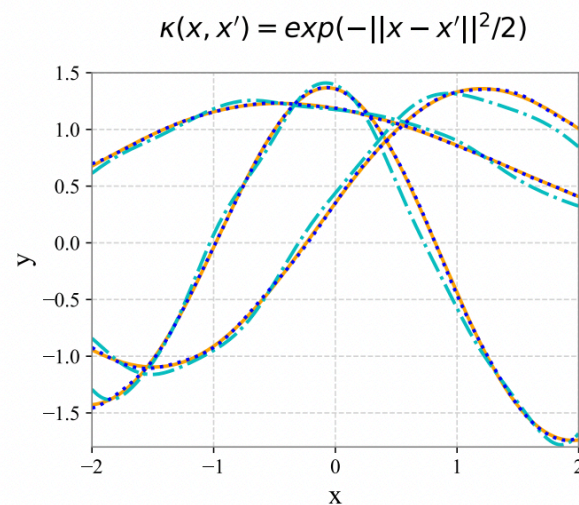
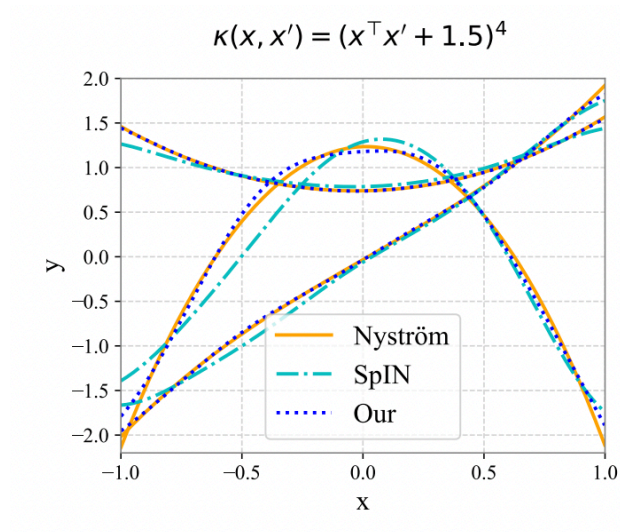


- NeuralEF approximate NTKs and NN-GP kernels with **less NN forward passes** than RFs
- It gives rise to an **unsupervised representation learning** paradigm, where the pairwise similarity captured by kernels is embedded into NNs
- It *relates two fields of research*



# The applications

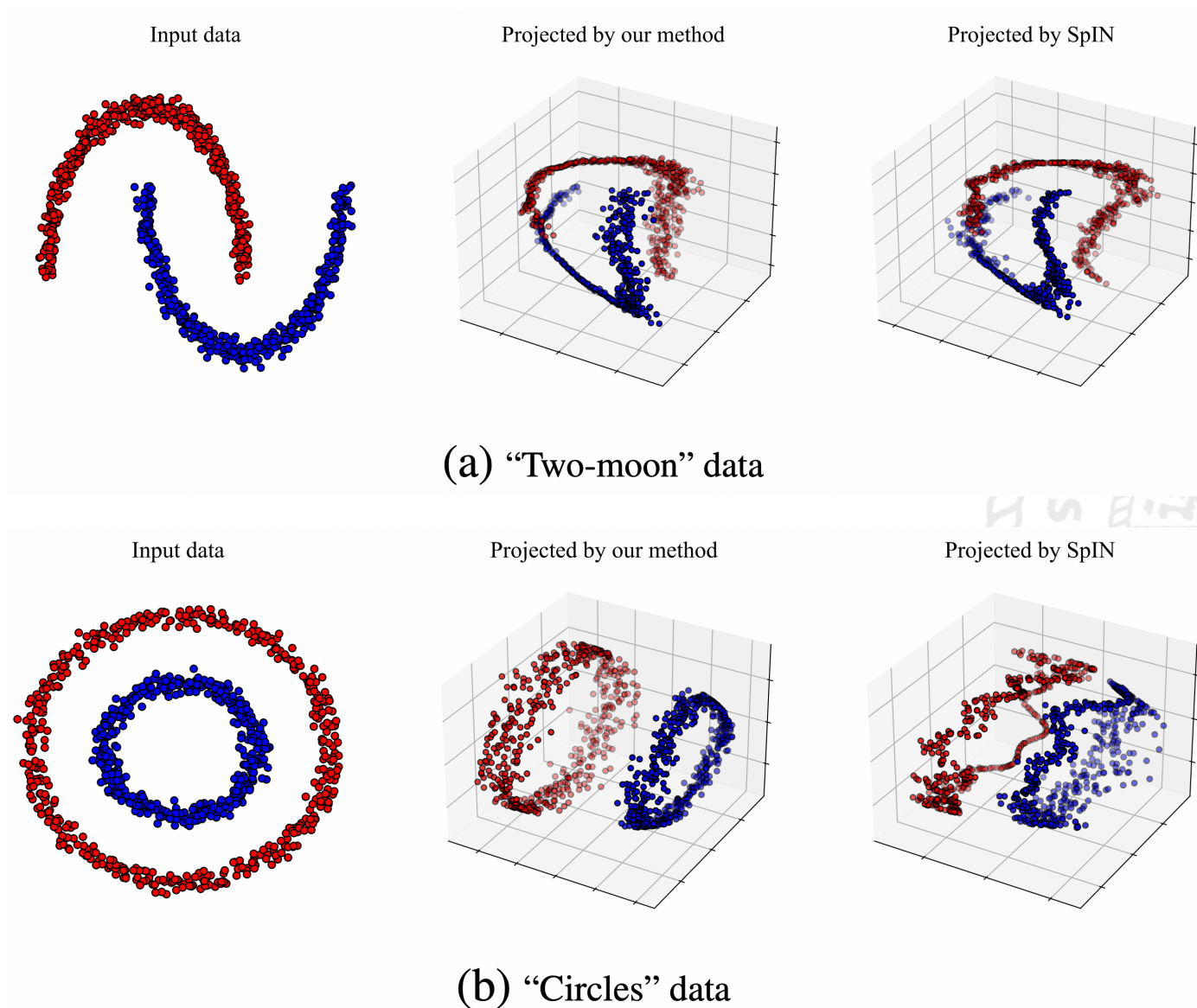
Find the eigenfunctions of classic kernels





# The applications

## Process MLP-GP kernels



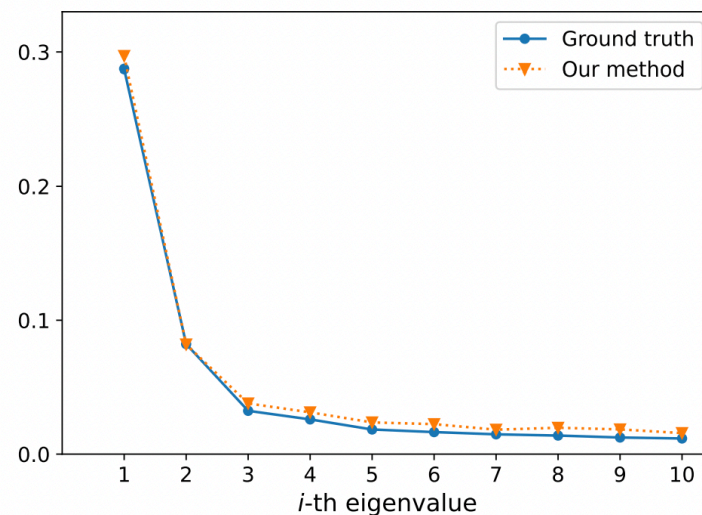
# The applications

Process CNN-GP kernels

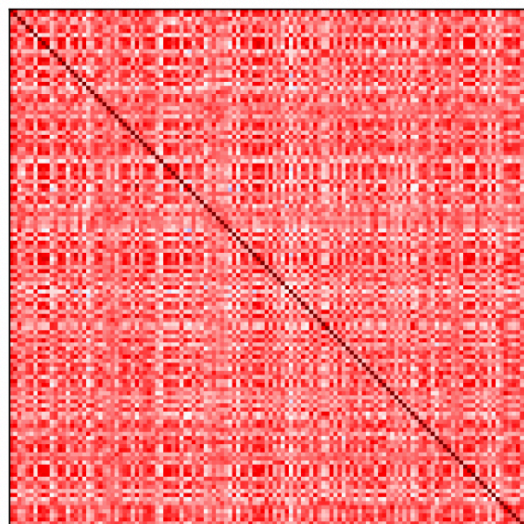
Method	LR test accuracy
<i>Our method (CNN-GP kernel)</i>	<b>84.98%</b>
<i>Nyström (CNN-GP kernel)</i>	N/A
<i>Nyström (polynomial kernel)</i>	78.00%
<i>Nyström (RBF kernel)</i>	77.55%

# The applications

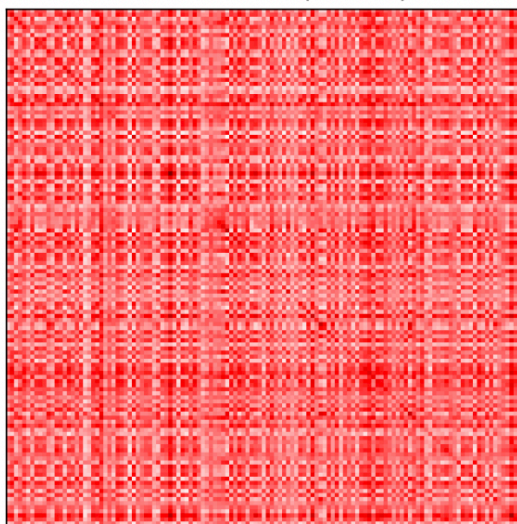
Find the eigenfunctions of NTK which itself is hard to compute



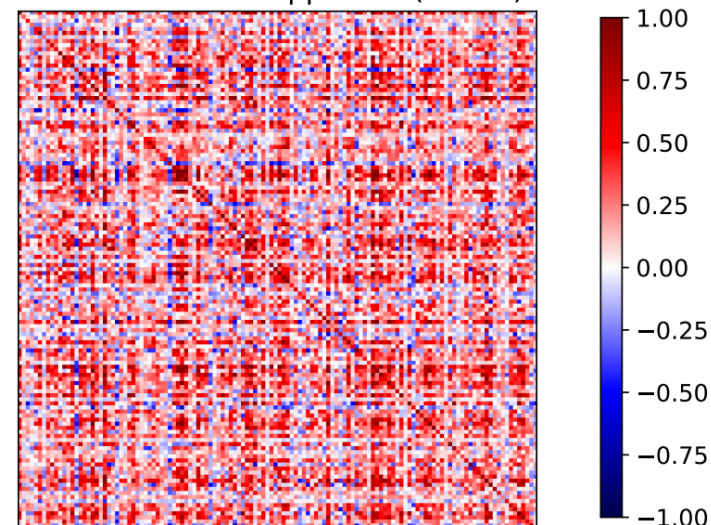
Ground truth



Our method ( $k = 10$ )



Random feature approach ( $S = 10$ )





# The applications

Improve linearised Laplace approximation with NeuralEF

$$\mathcal{GP}(f|g(\mathbf{x}, \boldsymbol{\theta}_{\text{MAP}}), \partial_{\boldsymbol{\theta}} g(\mathbf{x}, \boldsymbol{\theta}_{\text{MAP}}) \Sigma \partial_{\boldsymbol{\theta}} g(\mathbf{x}', \boldsymbol{\theta}_{\text{MAP}})^{\top})$$

Inverse of the Gauss-Newton of size # of params  $\times$  # of params

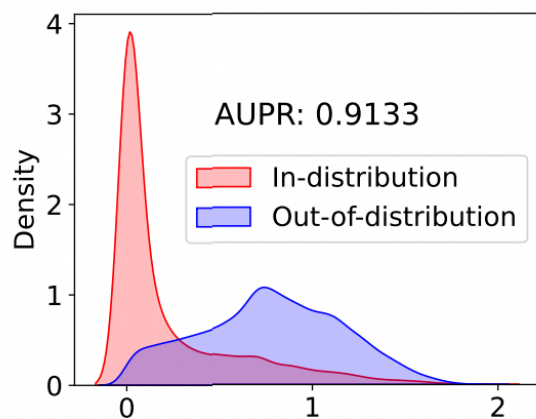
$$\Sigma^{-1} = \sum_i \partial_{\boldsymbol{\theta}} g(\mathbf{x}_i, \boldsymbol{\theta}_{\text{MAP}})^{\top} \Lambda_i \partial_{\boldsymbol{\theta}} g(\mathbf{x}_i, \boldsymbol{\theta}_{\text{MAP}}) + 1/\sigma_0^2 \mathbf{I}_{\dim(\boldsymbol{\theta})}$$

By Woodbury matrix identity and Mercer's theorem

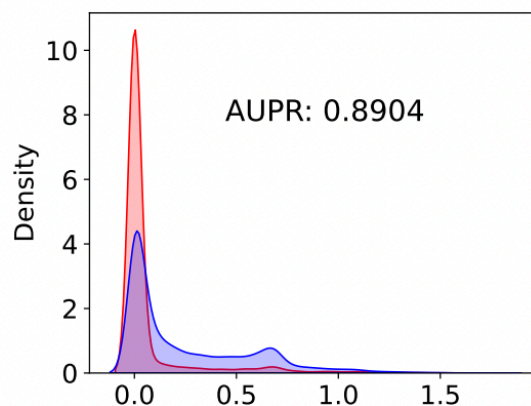
$$\mathcal{GP}\left(f|g(\mathbf{x}, \boldsymbol{\theta}_{\text{MAP}}), \Psi(\mathbf{x}) \left[ \sum_i \Psi(\mathbf{x}_i)^{\top} \Lambda_i \Psi(\mathbf{x}_i) + \frac{1}{\sigma_0^2} \mathbf{I}_k \right]^{-1} \Psi(\mathbf{x}')^{\top}\right),$$

Size:  $k \times k$

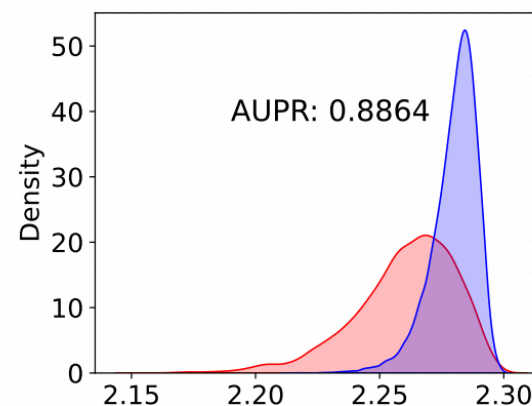
The concatenation of the  $k$  eigenfunctions for the NTK



(a) Ours



(b) MAP



(c) KFAC LLA

# The applications

## Bayesian deep learning by modeling SGD trajectory

$$\kappa_{\text{SGD}}(\mathbf{x}, \mathbf{x}') = \frac{1}{M} \sum_{i=1}^M (g(\mathbf{x}, \boldsymbol{\theta}_i) - \bar{g}(\mathbf{x})) (g(\mathbf{x}', \boldsymbol{\theta}_i) - \bar{g}(\mathbf{x}'))^\top$$

$$p(\mathbf{x}_{\text{new}}) = \int \mathcal{GP}(f | \bar{g}(\mathbf{x}), \kappa_{\text{SGD}}(\mathbf{x}, \mathbf{x}')) p(\mathbf{x}_{\text{new}} | f) df$$

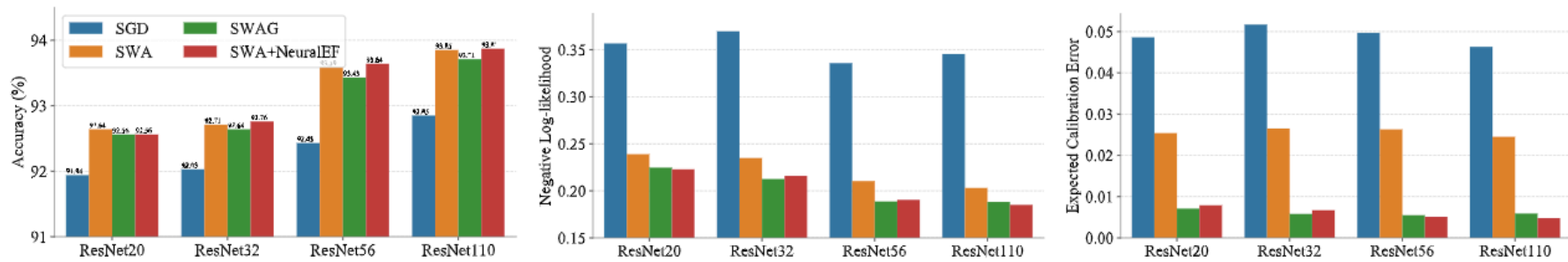
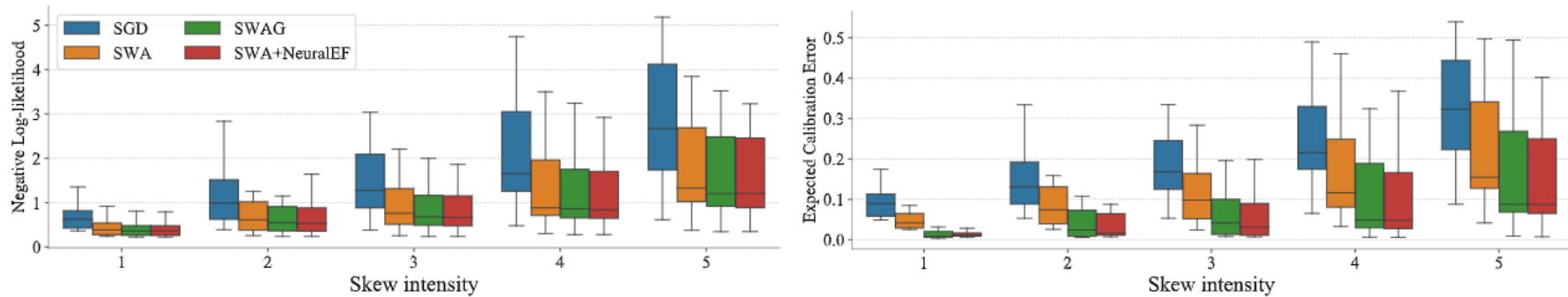


Figure 6: Test accuracy  $\uparrow$ , NLL  $\downarrow$ , and ECE  $\downarrow$  comparisons among models on CIFAR-10.



Thanks!

