

# Traversing Between Modes in Function Space for Fast Ensembling

---

EungGu Yun<sup>1,2</sup> Hyungi Lee<sup>1</sup> Giung Nam<sup>1</sup> Juho Lee<sup>1,3</sup>

<sup>1</sup>Korea Advanced Institute of Science and Technology (KAIST)

<sup>2</sup>Saige Research

<sup>3</sup>AITRICS

- **Deep Ensemble (DE)** is a simple yet powerful way to improve the performance of deep neural networks.
- Using **mode connectivity**, one can efficiently collect ensemble parameters in low-loss subspaces.
- However, for inference, one should still execute multiple forward passes.
- We propose a novel framework “bridge network” to **reduce inference costs** using mode connectivity properties in function space.

**Deep Ensemble** [Lakshminarayanan et al., 2017] is a simple algorithm that ensembles multiple neural networks where each network is trained with different random seeds.

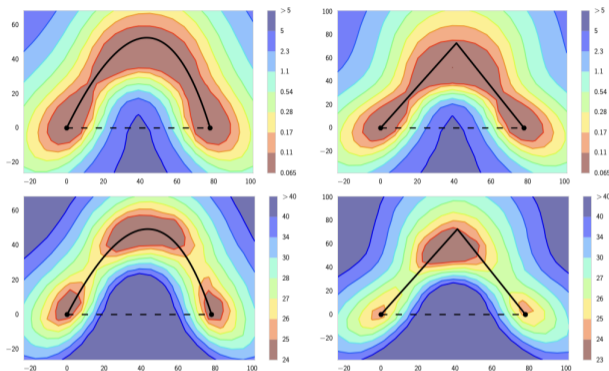
## **Benefits** of DE:

- Simple to implement.
- Improves both accuracy and uncertainty calibration.
- Easy to parallelize.

## **Drawbacks** of DE:

- Requires multiple training runs.
  - Requires multiple forward passes for inference.
- Computational cost increases with the number of ensembles.

# Mode Connectivity



**Figure 1:** Two modes in the loss surface and the connecting subspace. Left: *Bezier curve*, Right: *Polygonal chain*. (Figure from Garipov et al. [2018])

Garipov et al. [2018] and Draxler et al. [2018] showed that **modes** (local optima) in the loss surface of a deep neural network **are connected** by relatively simple low-dimensional subspaces where the loss in the subspace retains low values.

We focus on quadratic **Bezier curves** [Garipov et al., 2018]. Let  $\theta_i$  and  $\theta_j$  be two parameters of a neural network. The quadratic Bezier curve between them is defined as

$$\left\{ (1-r)^2\theta_i + 2r(1-r)\theta_{i,j}^{(\text{be})} + r^2\theta_j \mid r \in [0, 1] \right\}, \quad (1)$$

where  $\theta_{i,j}^{(\text{be})}$  is a *pin-point* parameter characterizing the curve. A low-loss subspace connecting  $(\theta_i, \theta_j)$  is found w.r.t.  $\theta_{i,j}^{(\text{be})}$  by minimizing

$$\int_0^1 \mathcal{L}(\theta_{i,j}^{(\text{be})}(r)) dr, \quad (2)$$

where  $\theta_{i,j}^{(\text{be})}(r)$  denotes the point at the position  $r$  of the curve,

$$\theta_{i,j}^{(\text{be})}(r) = (1-r)^2\theta_i + 2r(1-r)\theta_{i,j}^{(\text{be})} + r^2\theta_j, \quad (3)$$

and  $\mathcal{L} : \Theta \rightarrow \mathbb{R}$  is the loss function evaluating parameters.

Let  $\{\theta_1, \dots, \theta_m\}$  be a set of parameters independently trained as a deep ensemble. Then, for each pair  $(\theta_i, \theta_j)$ , we can construct a low-loss Bezier curve. For instance, choosing  $r = 0.5$ , we can collect  $\theta_{i,j}^{(\text{be})}(0.5)$  for all  $(i, j)$  pairs, and construct an ensembled predictor as

$$\frac{1}{m + \binom{m}{2}} \left( \sum_{i=1}^m f_{\theta_i}(\mathbf{x}) + \sum_{i < j} f_{\theta_{i,j}^{(\text{be})}(0.5)}(\mathbf{x}) \right). \quad (4)$$

While this strategy provide an effective way to increase the number of ensemble members, for inference, an additional  $\mathcal{O}(m^2)$  number of forward passes are required.

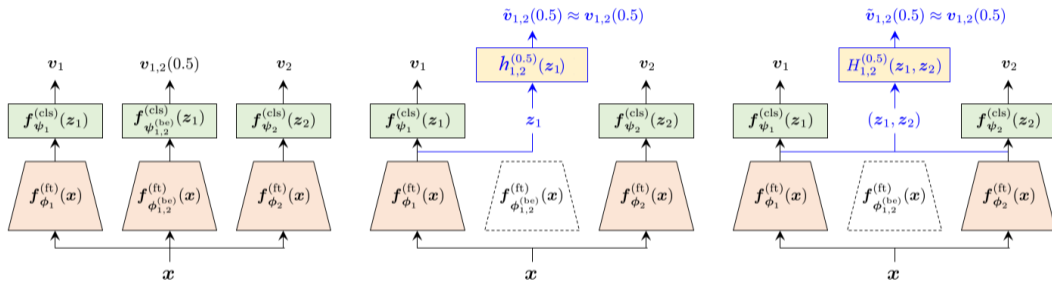
*Reduce the inference cost using mode connectivity properties in function space.*

**How:** Directly approximate the outputs evaluated at the subspace with a small auxiliary network, which is called “**bridge network**”.

**Assumption:** if two modes are connected by a simple subspace, we can predict the outputs corresponding to the parameters on the subspace using *only the outputs computed from the modes*.

*The bridge network lets us travel between modes in the function space.*

# Method Overview



**Figure 2:** Ensembles with a Bezier curve (**left**), a type I bridge network (**center**), and a type II bridge network (**right**).



**Assumption (revisited):** if two modes are connected by a simple low-loss subspace (Bezier curve), then we can predict the outputs corresponding to the parameters on the subspace using only the information obtained from the modes.

If such mapping exists, we may learn them via a **lightweight neural network**.

- **Features:**  $\mathbf{z}_i := f_{\phi_i}^{(\text{ft})}(\mathbf{x})$
- **Output:**  $\mathbf{v}_i := f_{\theta_i}(\mathbf{x}) = f_{\psi_i}^{(\text{cls})}(\mathbf{z}_i)$
- **Features (Bezier curve):**  $\mathbf{z}_{i,j}(r) := f_{\phi_{i,j}^{(\text{be})}(r)}(\mathbf{x})$
- **Output (Bezier curve):**  $\mathbf{v}_{i,j}(r) := f_{\theta_{i,j}^{(\text{be})}(r)}(\mathbf{x}) = f_{\psi_{i,j}^{(\text{be})}(r)}^{(\text{cls})}(\mathbf{z}_{i,j}(r))$

We **reuse** features  $\mathbf{z}_i$  to predict  $\mathbf{v}_{i,j}(r)$  with a lightweight neural network, which lets us **directly move from  $\mathbf{v}_i$  to  $\mathbf{v}_{i,j}(r)$  in the function space**.

A bridge network is usually constructed with a Convolutional Neural Network (CNN) whose inference cost is much lower than that of  $f_{\theta_i}$ .

## Type I Bridge Networks

A type I bridge network  $h_{i,j}^{(r)}$  takes a feature  $\mathbf{z}_i$  from only one mode, and predicts

$$\mathbf{v}_{i,j}(r) \approx \tilde{\mathbf{v}}_{i,j}(r) = h_{i,j}^{(r)}(\mathbf{z}_i). \quad (5)$$

An ensembled prediction with the type I bridge network is then constructed as

$$\frac{1}{2} \left( \mathbf{v}_i + h_{i,j}^{(r)}(\mathbf{z}_i) \right), \quad (6)$$

whose inference cost is not much higher than that of  $\mathbf{v}_i$ . One can also connect  $\theta_i$  with multiple modes  $\{\theta_{j_1}, \dots, \theta_{j_k}\}$ , learn bridge networks between  $(i, j_1), \dots, (i, j_k)$ , and construct an ensemble

$$\frac{1}{1+k} \left( \mathbf{v}_i + \sum_{j=1}^k h_{i,j_k}^{(r)}(\mathbf{z}_i) \right). \quad (7)$$

Still, since the costs for  $h_{i,j_k}^{(r)}$ s are far lower than  $\mathbf{v}_i$ , the inference cost does not significantly increase.

A type II bridge network  $H_{i,j}^{(r)}$  between  $(\theta_i, \theta_j)$  takes two features  $(\mathbf{z}_i, \mathbf{z}_j)$ , and predicts

$$\mathbf{v}_{i,j}(r) \approx \tilde{\mathbf{v}}_{i,j}(r) = H_{i,j}^{(r)}(\mathbf{z}_i, \mathbf{z}_j). \quad (8)$$

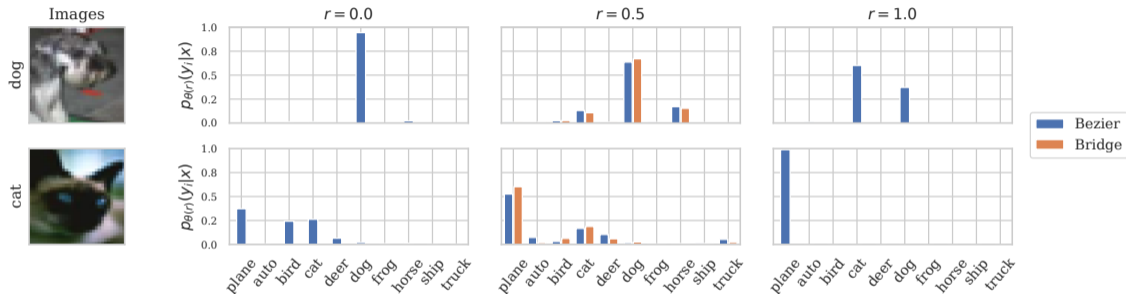
An ensemble prediction with the type II bridge network is then constructed as

$$\frac{1}{3} \left( \mathbf{v}_i + \mathbf{v}_j + H_{i,j}^{(r)}(\mathbf{z}_i, \mathbf{z}_j) \right), \quad (9)$$

where we construct an ensemble of three models with effectively two forward passes (for  $\mathbf{v}_i$  and  $\mathbf{v}_j$ ). Similar to the type I bridge networks, we may construct multiple bridges between a single curves and use them together for an ensemble

$$\frac{1}{k + \binom{k}{2}} \left( \sum_{i=1}^k \mathbf{v}_i + \sum_{i < j \leq k} H_{i,j}^{(r)}(\mathbf{z}_i, \mathbf{z}_j) \right). \quad (10)$$

# Correspondence



**Figure 3:** Bar plots in the third column show the class probability outputs of the bridge network (**orange**) and the base model with the Bezier parameters (**blue**) for given images displayed in the first column. We also depict the predicted outputs from the base model with  $\theta_1$  and  $\theta_2$  in the second and fourth columns, respectively.

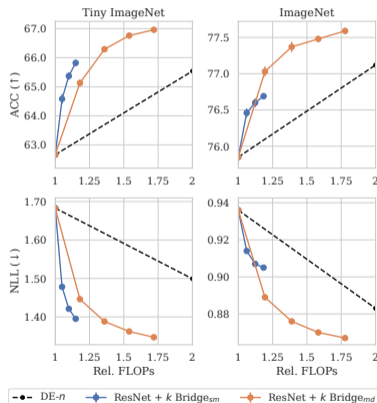
# Classification Performance (Type I)

**Table 1:** Performance improvement of the ensemble by adding type I bridges to the single base ResNet model on Tiny ImageNet and ImageNet datasets. Bridge<sub>sm</sub> and Bridge<sub>md</sub> denote the small and the medium versions of the bridge network based on their FLOPs.

Tiny ImageNet						
Model	FLOPs ( $\downarrow$ )	#Params ( $\downarrow$ )	ACC ( $\uparrow$ )	NLL ( $\downarrow$ )	ECE ( $\downarrow$ )	DEE ( $\uparrow$ )
ResNet (DE-1)	$\times 1.000$	$\times 1.000$	$62.66 \pm 0.23$	$1.683 \pm 0.009$	$0.050 \pm 0.004$	1.000
+ 1 Bridge <sub>sm</sub>	$\times 1.050$	$\times 1.057$	$64.58 \pm 0.17$	$1.478 \pm 0.006$	$0.025 \pm 0.002$	$2.280 \pm 0.086$
+ 2 Bridge <sub>sm</sub>	$\times 1.099$	$\times 1.114$	$65.37 \pm 0.13$	$1.421 \pm 0.004$	$0.018 \pm 0.002$	$3.087 \pm 0.118$
+ 3 Bridge <sub>sm</sub>	$\times 1.149$	$\times 1.171$	<b><math>65.82 \pm 0.10</math></b>	<b><math>1.395 \pm 0.003</math></b>	<b><math>0.015 \pm 0.001</math></b>	<b><math>3.680 \pm 0.133</math></b>
+ 1 Bridge <sub>md</sub>	$\times 1.180$	$\times 1.206$	$65.13 \pm 0.12$	$1.446 \pm 0.002$	$0.034 \pm 0.002$	$2.709 \pm 0.049$
+ 2 Bridge <sub>md</sub>	$\times 1.359$	$\times 1.412$	$66.29 \pm 0.06$	$1.388 \pm 0.004$	$0.025 \pm 0.001$	$3.845 \pm 0.171$
+ 3 Bridge <sub>md</sub>	$\times 1.539$	$\times 1.618$	<b><math>66.76 \pm 0.09</math></b>	<b><math>1.362 \pm 0.003</math></b>	<b><math>0.023 \pm 0.001</math></b>	<b><math>4.708 \pm 0.209</math></b>
DE-2	$\times 2.000$	$\times 2.000$	$65.54 \pm 0.25$	$1.499 \pm 0.007$	$0.029 \pm 0.003$	2.000

ImageNet						
Model	FLOPs ( $\downarrow$ )	#Params ( $\downarrow$ )	ACC ( $\uparrow$ )	NLL ( $\downarrow$ )	ECE ( $\downarrow$ )	DEE ( $\uparrow$ )
ResNet (DE-1)	$\times 1.000$	$\times 1.000$	$75.85 \pm 0.06$	$0.936 \pm 0.003$	$0.019 \pm 0.001$	1.000
+ 1 Bridge <sub>sm</sub>	$\times 1.061$	$\times 1.071$	$76.46 \pm 0.06$	$0.914 \pm 0.000$	$0.012 \pm 0.001$	$1.418 \pm 0.034$
+ 2 Bridge <sub>sm</sub>	$\times 1.123$	$\times 1.141$	$76.60 \pm 0.06$	$0.907 \pm 0.000$	$0.012 \pm 0.001$	$1.537 \pm 0.026$
+ 3 Bridge <sub>sm</sub>	$\times 1.184$	$\times 1.212$	<b><math>76.69 \pm 0.04</math></b>	<b><math>0.905 \pm 0.000</math></b>	<b><math>0.011 \pm 0.001</math></b>	<b><math>1.584 \pm 0.021</math></b>
+ 1 Bridge <sub>md</sub>	$\times 1.194$	$\times 1.222$	$77.03 \pm 0.07$	$0.889 \pm 0.001$	<b><math>0.013 \pm 0.000</math></b>	$1.881 \pm 0.022$
+ 2 Bridge <sub>md</sub>	$\times 1.389$	$\times 1.444$	$77.37 \pm 0.07$	$0.876 \pm 0.001$	<b><math>0.013 \pm 0.001</math></b>	$2.341 \pm 0.076$
+ 3 Bridge <sub>md</sub>	$\times 1.583$	$\times 1.665$	<b><math>77.48 \pm 0.03</math></b>	<b><math>0.870 \pm 0.000</math></b>	<b><math>0.013 \pm 0.000</math></b>	<b><math>2.618 \pm 0.062</math></b>
DE-2	$\times 2.000$	$\times 2.000$	$77.12 \pm 0.04$	$0.883 \pm 0.001$	$0.012 \pm 0.001$	2.000



**Figure 4:** The cost-performance plots of type I bridges compared to DE on Tiny ImageNet and ImageNet datasets. On the basis of DE (black dashed line), the upper left is preferable in ACC, and the lower left is preferable in NLL.

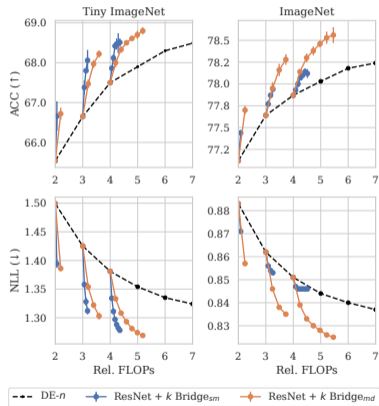
# Classification Performance (Type II)

**Table 2:** Performance improvement of the ensemble by adding type II bridges as members to existing DE ensembles on Tiny ImageNet and ImageNet datasets. Bridge<sub>sm</sub> and Bridge<sub>md</sub> denote the small and the medium versions of the bridge network based on their FLOPs.

Tiny ImageNet						
Model	FLOPs ( $\downarrow$ )	#Params ( $\downarrow$ )	ACC ( $\uparrow$ )	NLL ( $\downarrow$ )	ECE ( $\downarrow$ )	DEE ( $\uparrow$ )
DE-4	$\times 4.000$	$\times 4.000$	$67.50 \pm 0.11$	$1.381 \pm 0.004$	$0.018 \pm 0.001$	4.000
+ 1 Bridge <sub>sm</sub>	$\times 4.058$	$\times 4.067$	$67.86 \pm 0.05$	$1.334 \pm 0.003$	$0.017 \pm 0.002$	$6.051 \pm 0.181$
+ 2 Bridge <sub>sm</sub>	$\times 4.117$	$\times 4.135$	$68.12 \pm 0.09$	$1.311 \pm 0.005$	$0.015 \pm 0.001$	$8.174 \pm 0.465$
+ 4 Bridge <sub>sm</sub>	$\times 4.234$	$\times 4.269$	$68.47 \pm 0.14$	$1.288 \pm 0.004$	$0.015 \pm 0.001$	$10.340 \pm 0.773$
+ 6 Bridge <sub>sm</sub>	$\times 4.351$	$\times 4.404$	<b><math>68.51 \pm 0.10</math></b>	<b><math>1.278 \pm 0.003</math></b>	<b><math>0.014 \pm 0.001</math></b>	<b><math>11.268 \pm 0.871</math></b>
+ 1 Bridge <sub>md</sub>	$\times 4.198$	$\times 4.226$	$68.00 \pm 0.11$	$1.333 \pm 0.003$	<b><math>0.019 \pm 0.001</math></b>	$6.183 \pm 0.120$
+ 2 Bridge <sub>md</sub>	$\times 4.395$	$\times 4.453$	$68.33 \pm 0.08$	$1.308 \pm 0.003$	<b><math>0.019 \pm 0.001</math></b>	$8.489 \pm 0.481$
+ 4 Bridge <sub>md</sub>	$\times 4.791$	$\times 4.906$	$68.61 \pm 0.05$	$1.281 \pm 0.004$	$0.021 \pm 0.003$	$10.897 \pm 0.800$
+ 6 Bridge <sub>md</sub>	$\times 5.186$	$\times 5.359$	<b><math>68.80 \pm 0.09</math></b>	<b><math>1.269 \pm 0.003</math></b>	$0.021 \pm 0.001$	<b><math>12.110 \pm 1.083</math></b>
DE-5	$\times 5.000$	$\times 5.000$	$67.90 \pm 0.14$	$1.354 \pm 0.003$	$0.019 \pm 0.001$	5.000

ImageNet						
Model	FLOPs ( $\downarrow$ )	#Params ( $\downarrow$ )	ACC ( $\uparrow$ )	NLL ( $\downarrow$ )	ECE ( $\downarrow$ )	DEE ( $\uparrow$ )
DE-4	$\times 4.000$	$\times 4.000$	$77.87 \pm 0.04$	$0.851 \pm 0.001$	$0.012 \pm 0.001$	4.000
+ 1 Bridge <sub>sm</sub>	$\times 4.086$	$\times 4.088$	$77.93 \pm 0.02$	$0.847 \pm 0.000$	$0.012 \pm 0.001$	$4.580 \pm 0.052$
+ 2 Bridge <sub>sm</sub>	$\times 4.172$	$\times 4.176$	$78.00 \pm 0.04$	<b><math>0.846 \pm 0.000</math></b>	<b><math>0.011 \pm 0.000</math></b>	$4.739 \pm 0.052$
+ 4 Bridge <sub>sm</sub>	$\times 4.343$	$\times 4.351$	$78.10 \pm 0.03$	<b><math>0.846 \pm 0.000</math></b>	<b><math>0.011 \pm 0.001</math></b>	<b><math>4.768 \pm 0.041</math></b>
+ 6 Bridge <sub>sm</sub>	$\times 4.515$	$\times 4.527$	<b><math>78.12 \pm 0.05</math></b>	<b><math>0.846 \pm 0.001</math></b>	<b><math>0.011 \pm 0.001</math></b>	$4.659 \pm 0.037$
+ 1 Bridge <sub>md</sub>	$\times 4.243$	$\times 4.256$	$78.14 \pm 0.03$	$0.839 \pm 0.000$	<b><math>0.011 \pm 0.001</math></b>	$6.123 \pm 0.121$
+ 2 Bridge <sub>md</sub>	$\times 4.487$	$\times 4.512$	$78.30 \pm 0.05$	$0.833 \pm 0.000$	$0.012 \pm 0.001$	$8.068 \pm 0.144$
+ 4 Bridge <sub>md</sub>	$\times 4.973$	$\times 5.024$	$78.46 \pm 0.04$	$0.828 \pm 0.000$	$0.012 \pm 0.000$	$9.951 \pm 0.163$
+ 6 Bridge <sub>md</sub>	$\times 5.460$	$\times 5.536$	<b><math>78.56 \pm 0.09</math></b>	<b><math>0.825 \pm 0.000</math></b>	$0.012 \pm 0.001$	<b><math>10.760 \pm 0.202</math></b>
DE-5	$\times 5.000$	$\times 5.000$	$78.03 \pm 0.03$	$0.844 \pm 0.001$	$0.012 \pm 0.001$	5.000



**Figure 5:** The cost-performance plots of type II bridges compared to DE on Tiny ImageNet and ImageNet datasets. On the basis of DE (black dashed line), the upper left is preferable in ACC, and the lower left is preferable in NLL.

We proposed a novel framework for efficient ensembling that reduces inference costs of ensembles with a lightweight network called **bridge networks**.

Through empirical validation, we show that

1. **Bridge networks can approximate outputs** of connecting subspaces quite accurately with minimal computation cost.
2. DES augmented with bridge networks can significantly **reduce inference costs** without big sacrifice in performance.

## References

- F. Draxler, K. Veschgini, M. Salmhofer, and F. Hamprecht. Essentially no barriers in neural network energy landscape. In *Proceedings of The 35th International Conference on Machine Learning (ICML)*, 2018.
- T. Garipov, P. Izmailov, D. Podoprikin, D. Vetrov, and A. G. Wilson. Loss surfaces, mode connectivity, and fast ensembling of DNNs. In *Advances in Neural Information Processing Systems 31 (NeurIPS)*, 2018.
- B. Lakshminarayanan, A. Pritzel, and C. Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In *Advances in Neural Information Processing Systems 30 (NIPS)*, 2017.