

# Promoting Ensemble Diversity with Interactive Bayesian Distributional Robustness for Fine-tuning Foundation Models

Ngoc-Quan Pham\*, Tuan Truong\*, Quyen Tran, Tan Nguyen, Dinh Phung, Trung Le



# How do we utilize large pre-trained models?

- Foundational models are increasingly demonstrating remarkable capabilities over a wide array of tasks.
- **Goal:** Effectively utilizing these pre-trained models for **downstream tasks**. However, adapting these models via **full fine-tuning** presents significant limitations: **high computational cost**, **overfitting**, **storage overhead**...



DALL-E



Claude



deepseek

# Parameter-efficient Fine-tuning (PEFT)

- Techniques to adapt large pre-trained models with minimal parameter updates.
- **Pros:** Reduce computational cost and memory usage while maintaining performance, and preserving pre-trained knowledge.
- **Cons:** PEFT methods can lead to overconfident predictions, especially when fine-tuned on small datasets.

# Motivations

Our method relies on:

- **Bayesian Inference**: enhances robustness, tackling uncertainty
- **Flat minimizers**: improve neural network generalization by helping models find broader local minima, making them more robust
- **Distributional Robustness**: a framework for learning under distributional uncertainty, which seeks the worst-case among a ball of “local distributions”.

# Motivations

- We also want to promote ensemble diversity
- Define the **approximate posterior distribution**  $Q^K$ , where samples are concatenated models  $\theta_{1:K}$
- We learn  $Q$  so that the sampled models reside in *low-loss, low-sharpness* regions while *maintaining ensemble diversity*.
- Leverage DRO to alleviate training instabilities

# Interactive Bayesian Distributional Robustness

- We propose a *distributional population loss* where  $l_{div}$  encourages the diversity among model particles

$$\mathcal{L}_{\mathcal{D}}(Q^K) = \mathbb{E}_{\theta \sim Q^K} \left[ \frac{1}{K} \sum_{i=1}^K \mathcal{L}_{\mathcal{D}}(\theta_i) + \alpha \mathbb{E}_{\mathcal{D}} \left[ l_{div}(\theta_{1:K}; x, y) \right] \right].$$

**Theorem 4.1.** *With the probability at least  $1 - \delta$  over the choice of  $\mathcal{S} \sim \mathcal{D}^N$ , we have*

$$\begin{aligned} \mathcal{L}_{\mathcal{D}}(Q^K) &\leq \min_{\lambda \geq 0} \left\{ \lambda \rho + \mathbb{E}_{\theta \sim Q^K} \left[ \max_{\theta'} \left\{ \mathcal{L}_{\mathcal{S}}(\theta') - \lambda c^K(\theta, \theta') \right\} \right] \right\} \\ &\quad + L \sqrt{\frac{K D_{KL}(Q, P) + \log \frac{1}{\delta}}{2N}} \end{aligned}$$

**Remark:** This framework operates on the joint distribution  $Q^K$  and incorporates the divergence loss  $l_{div}$ , enabling us to model interactions between the particle models  $\theta_{1:K}$

# Divergence Loss

- Let  $f_{-y}^i$  be the non-maximal prediction probabilities by eliminating the prediction probability of the ground-truth label  $y$
- We encourage the non-maximal predictions to diverge, while maximizing prediction probability of the ground-truth label.
- Motivated by the theory of Determinantal Point Processes, we define the ensemble diversity:

$$l_{div}(\theta_{1:K}; x, y) = \text{Vol}^2\left(\left[\tilde{f}_{-y}^i\right]_{i \in [C]}\right)$$

$$\text{where } \tilde{f}_{-y}^i = \frac{f_{-y}^i}{\|f_{-y}^i\|}, \left[\tilde{f}_{-y}^i\right]_{i \in [C]} \in \mathbb{R}^{(C-1) \times K}, [C] = \{1, \dots, C\}.$$

# Practical Method

- Define  $Q = \frac{1}{K} \sum_{i=1}^K \mathcal{N}(\mu_i, \sigma^2 \mathbb{I})$  and  $P = \mathcal{N}(\mathbf{0}, \mathbb{I})$ . With a few relaxations, the problem becomes

$$\min_{\mu_{1:K}, \sigma} \min_{\lambda \geq 0} \left\{ \lambda \rho + \mathbb{E}_{\theta_{1:K} \sim Q} \left[ \frac{1}{K} \sum_{i=1}^K \max_{\theta'_i} \tilde{\ell}(\theta'_i, \theta_i; x, y) \right] \right\} + \frac{\beta}{K} \left[ \sum_{i=1}^K \|\mu_i\|^2 + d(\sigma - \log \sigma) \right]$$

where  $\tilde{\ell}(\theta'_i, \theta_i; x, y) = l(\theta'_i; x, y) + \alpha l_{div}(\theta'_i, \theta_{-i}; x, y) - \lambda c(\theta_i, \theta'_i)$

- We alternatively update  $\mu_{1:K}$  and  $\lambda$  with gradient descent, and update  $\theta'_i$  with a gradient ascent



# Interactive Bayesian Distributional Robustness

---

**Algorithm 1** Interactive Bayesian Distributional Robustness (IBDR)

---

**Input:** Initial particle means  $\mu_{1:K}$ ; ascend step size  $\alpha_1$ ; learning rates  $\alpha_\lambda, \alpha_\mu$

**Output:** Optimal particle means  $\mu_{1:K}$

**while** not converged **do**

    Sample batch  $\mathcal{B} = \{(x_1, y_1), \dots, (x_b, y_b)\}$

    Sample  $\epsilon_i \sim \mathbb{N}(0, \mathbb{I})$  and  $\theta_i \leftarrow \mu_i + \sigma \epsilon_i$

    Compute  $\theta'_i \leftarrow \theta_i + \alpha_1 \nabla_{\theta_i} \tilde{\ell}(\theta'_i, \theta_i; x, y)$

    Compute  $\lambda \leftarrow \lambda - \alpha_\lambda \nabla_\lambda \bar{\mathcal{L}}(\lambda, \theta'_i, \theta_i; x, y)$

    Compute  $\mu_i \leftarrow \lambda - \alpha_\mu \nabla_{\mu_i} \bar{\mathcal{L}}(\lambda, \theta'_i, \theta_i; x, y)$

**end while**

**return**  $\mu_{1:K}$

---

# Experiments

## Image Classification

Table 1. Top-1 Accuracy on VTAB-1K. The accuracies are reported with ViT-B/16 pre-trained on ImageNet-21K

	Natural							Specialized				Structured								
Method	CIFAR100	Caltech101	DTD	Flowers102	Pets	SVHN	Sun397	Camelyon	EuroSAT	Resisc45	Retinopathy	Clevr-Count	Clevr-Dist	DMLab	KITTI	dSpr-Loc	dSpr-Ori	sNORB-Azim	sNORB-Ele	AVG
FFT	68.9	87.7	64.3	97.2	86.9	87.4	38.8	79.7	<b>95.7</b>	84.2	73.9	56.3	58.6	41.7	65.5	57.5	46.7	25.7	29.1	62.3
LoRA	67.1	90.7	68.9	98.1	90.1	84.5	54.2	84.1	94.9	84.4	73.6	<b>82.9</b>	69.2	49.8	78.5	75.7	47.1	<b>31.0</b>	<b>44.0</b>	68.4
SAM	72.7	90.3	71.4	99.0	90.2	84.4	52.4	82.0	92.6	84.1	74.0	76.7	68.3	47.9	74.3	71.6	43.4	26.9	39.1	70.5
SA-BNN	65.1	91.5	71.0	98.9	89.4	89.3	55.2	<b>86.2</b>	94.5	86.4	75.2	61.4	63.2	40.0	71.3	64.5	34.5	27.2	31.2	68.2
SGLD	68.7	91.0	67.0	98.6	89.3	83.0	51.6	81.2	93.7	83.2	76.4	80.0	70.1	48.2	76.2	71.1	39.3	31.2	38.4	68.4
DeepEns	68.6	88.9	67.7	98.9	90.7	85.1	54.5	82.6	94.8	82.7	75.3	46.6	47.1	47.4	68.2	71.1	36.6	30.1	35.6	67.0
BayesTune	68.2	91.7	69.5	99.0	90.7	86.4	51.2	84.9	95.3	84.1	75.1	82.8	68.9	49.7	79.3	74.3	46.6	30.3	42.8	68.5
SVGD	71.3	90.2	71.0	98.7	90.2	84.3	52.7	83.4	93.2	86.7	75.1	75.8	70.7	49.6	79.9	69.1	41.2	30.6	33.1	70.9
IBDR	<b>73.0</b> (.11)	<b>92.1</b> (.31)	<b>71.7</b> (.12)	<b>99.3</b> (0.15)	<b>91.4</b> (0.16)	<b>91.3</b> (.36)	<b>56.7</b> (.18)	85.1 (.24)	95.0 (.44)	<b>87.3</b> (.14)	<b>76.5</b> (.12)	78.1 (.11)	<b>75.1</b> (.24)	<b>53.6</b> (.42)	<b>80.4</b> (.26)	<b>77.1</b> (.29)	<b>49.3</b> (.19)	28.9 (.13)	40.1 (.37)	<b>73.6</b>

# Experiments

## Image Classification

Table 2. Expected Calibration Errors (ECE) on VTAB-1K. The results are reported with ViT-B/16 pre-trained on ImageNet-21K

	Natural							Specialized				Structured								
Method	CIFAR100	Caltech101	DTD	Flowers102	Pets	SVHN	Sun397	Camelyon	EuroSAT	Resisc45	Retinopathy	Clevr-Count	Clevr-Dist	DMLab	KITTI	dSpr-Loc	dSpr-Ori	sNORB-Azim	sNORB-Ele	AVG
FFT	0.29	0.23	0.20	0.13	0.27	0.19	0.45	0.21	0.13	0.18	0.17	0.41	0.44	0.42	0.22	0.14	0.23	0.24	0.40	0.26
LoRA	0.38	0.19	0.18	0.05	0.09	0.10	0.14	<b>0.11</b>	0.09	0.12	0.11	0.12	0.19	0.34	0.18	0.14	0.21	0.18	0.31	0.17
SAM	0.21	0.25	0.20	0.11	0.12	0.15	0.14	0.17	0.16	0.14	<b>0.09</b>	0.12	0.17	0.24	0.16	0.21	<b>0.19</b>	0.13	0.16	0.16
SA-BNN	0.22	0.08	0.19	0.15	0.12	0.12	0.24	0.13	<b>0.06</b>	0.12	0.18	0.14	<b>0.21</b>	0.22	0.24	0.25	0.41	0.46	0.34	0.20
SGLD	0.26	0.20	<b>0.17</b>	0.05	0.18	0.14	0.23	0.18	0.09	0.12	0.32	0.26	0.29	<b>0.21</b>	0.26	0.42	0.39	<b>0.11</b>	0.24	0.22
DeepEns	0.24	0.12	0.22	0.04	0.10	0.13	0.23	0.16	0.07	0.15	0.21	0.31	0.32	0.36	0.13	0.32	0.31	0.16	0.29	0.20
BayesTune	0.32	0.93	0.20	0.03	0.85	0.12	0.22	0.13	0.07	0.13	0.22	<b>0.12</b>	0.23	0.30	0.24	0.28	0.28	0.31	0.26	0.23
SVGD	0.20	0.13	0.19	0.04	0.16	0.09	0.20	0.15	0.11	0.13	0.12	0.17	0.21	0.30	0.18	0.21	0.25	0.14	0.26	0.18
IBDR	<b>0.16</b> (.03)	<b>0.08</b> (.02)	0.19 (.02)	<b>0.02</b> (.01)	<b>0.07</b> (.01)	<b>0.07</b> (.01)	<b>0.13</b> (.02)	0.12 (.03)	0.06 (.02)	<b>0.11</b> (.02)	0.11 (.01)	0.13 (.01)	0.24 (.02)	0.30 (.03)	<b>0.12</b> (.01)	<b>0.11</b> (.01)	0.30 (.05)	0.30 (.04)	<b>0.16</b> (.02)	<b>0.14</b>

# Experiments

## Commonsense Reasoning

Table 3. Accuracy/ECE on six common-sense reasoning datasets

Metric		Datasets						
Type	Method	WG-S	ARC-C	ARC-E	WG-M	OBQA	BoolQ	AVG
ACC ( $\uparrow$ )	MLE	68.99	69.10	85.65	74.53	81.52	86.53	77.72
	MAP	68.62	67.59	86.55	75.61	81.38	86.50	77.71
	MCD	69.26	68.43	86.07	76.18	81.49	87.15	78.10
	ENS	69.57	66.20	84.40	75.32	81.38	87.09	77.33
	BBB	67.54	68.11	85.63	73.41	81.72	<b>87.19</b>	77.27
	LAP	69.20	66.78	80.05	75.55	82.12	86.95	76.78
	BLoB	70.89	<b>70.83</b>	86.68	74.55	82.73	86.80	78.75
	IBDR	<b>72.51</b>	70.56	<b>86.95</b>	<b>76.46</b>	<b>84.60</b>	86.89	<b>79.66</b>
ECE ( $\downarrow$ )	MLE	29.83	29.00	13.12	20.62	12.55	3.18	18.05
	MAP	29.76	29.42	12.07	23.07	13.26	3.16	18.46
	MCD	28.06	27.73	12.31	18.27	15.12	3.49	17.50
	ENS	28.52	29.16	12.57	20.86	15.34	9.61	19.34
	BBB	21.93	25.84	12.42	15.89	11.23	3.76	15.18
	LAP	<b>4.15</b>	<b>16.25</b>	33.29	<b>7.40</b>	8.70	<b>1.30</b>	<b>11.85</b>
	BLoB	20.62	20.61	<b>9.43</b>	11.23	8.36	2.46	12.12
	IBDR	24.17	21.20	9.71	11.19	<b>5.82</b>	1.54	12.27

# Conclusion

- We introduce a novel Bayesian framework that explicitly models the interaction between particles
- We propose Interactive Bayesian Distributional Robustness, which simultaneously enhances ensemble diversity, generalization ability, and distributional robustness

**Thank you**