

# Function Space Diversity for Uncertainty Prediction via Repulsive Last-Layer Ensembles

Anonymous authors

Paper under double-blind review

## Abstract

Bayesian inference in function space has gained attention due to its robustness against overparameterization in neural networks. However, approximating the infinite-dimensional function space introduces several challenges. In this work, we discuss function space inference via particle optimization and present practical modifications that improve uncertainty estimation and, most importantly, make it applicable for large and pretrained networks. First, we demonstrate that the input samples, where particle predictions are enforced to be diverse, are detrimental to the model performance. While diversity on training data itself can lead to underfitting, the use of label-destroying data augmentation, or unlabeled out-of-distribution data can improve prediction diversity and uncertainty estimates. Furthermore, we take advantage of the function space formulation, which imposes no restrictions on network parameterization other than sufficient flexibility. Instead of using full deep ensembles to represent particles, we propose a single multi-headed network that introduces a minimal increase in parameters and computation. This allows seamless integration to pretrained networks, where this repulsive last-layer ensemble can be used for uncertainty aware fine-tuning at minimal additional cost. We achieve competitive results in disentangling aleatoric and epistemic uncertainty, detecting out-of-distribution data, and providing calibrated uncertainty estimates under distribution shifts with minimal compute and memory.

## 1 Introduction

Deep learning is becoming ubiquitous in our lives, with applications ranging from medical diagnosis to autonomous driving. However, in safety-critical scenarios accurate predictions alone are not sufficient. In addition, models should provide well-calibrated uncertainty estimates to mitigate overconfidence and the potential risks associated with erroneous predictions. Uncertainty methods in deep learning typically distinguish between two types (Kendall & Gal, 2017; Hüllermeier & Waegeman, 2021): (a) Aleatoric uncertainty, which arises from inherent data ambiguity or noise, and (b) epistemic uncertainty, which corresponds to model uncertainty resulting from a lack of knowledge and observations.

Epistemic uncertainty is often estimated by deep ensembles (DEs) (Lakshminarayanan et al., 2017). For the estimate to be accurate, each ensemble member must entail a sufficiently different optimum of the posterior distribution. Particle-optimization variational inference (POVI) (Liu & Wang, 2016; Liu, 2017; Liu et al., 2019) achieves this diversity among ensemble members (i.e., the particles) by incorporating a repulsion term during parameter optimization. However, parameter diversity alone is insufficient. This is because neural networks with different parameters can still represent similar functions. It is thus advisable to avoid this issue by performing inference directly in the function space, enforcing diversity therein (Wang et al., 2019; D’Angelo & Fortuin, 2021). Yet, despite its theoretical appeal, function space POVI often performs worse than standard DEs in both accuracy and quality of uncertainty estimation (D’Angelo & Fortuin, 2021; Trinh et al., 2023; Yashima et al., 2022). In this work, we discuss the underlying reason for the performance gap; it is not because of the function space diversity but because of the challenges in accurately approximating the infinite-dimensional function space.

It remains practically infeasible to achieve function space diversity over the whole input domain (particularly for high dimensional input data). Therefore, good repulsion samples must not only be *diverse* but also capture the most *relevant parts* of the input domain. The training data itself is generally not rich enough and, as such, insufficient for accurate uncertainty estimation (D’Angelo & Fortuin, 2021; Trinh et al., 2023). We demonstrate how to improve this without sacrificing accuracy: the key is to utilize unlabeled out-of-distribution (OOD) data. If OOD data is unavailable, label-destroying data augmentation<sup>1</sup> can achieve improvements in certain tasks. Our evaluation confirms that well-chosen repulsion samples suppress reliance on spurious features, improve uncertainty estimation, and achieve reliable OOD detection.

Training and storing multiple ensemble members requires substantial computational resources, especially since the entire ensemble must be optimized jointly to maintain diversity. Function space POVI allows for flexible parameterizations. Drawing inspiration from ensemble distillation (Tran et al., 2020), we propose a multi-headed architecture (see Fig. 1). That is, we first learn a single base network that we subsequently equip with multiple heads – each representing one particle. While such last-layer ensembles (LL-Es) have been proposed in prior work, they are often not diverse enough; thus, we will reintroduce this diversity via repulsion either in parameter or in function space. This last layer repulsion has additional benefits. Specifically, the parameter space has lower dimension and suffers less from overparameterization while the function space allows for explicit prediction differences in relevant input-space regions.

Moreover, last-layer ensembles are especially well-suited for integration with pretrained networks. Pretraining is typically performed on large, diverse datasets to learn general and expressive features that benefit downstream tasks. This, however, implicitly prevents the use of deep ensembles that rely on random initializations to introduce diversity. By equipping a pretrained network with a last-layer ensemble, diversity can still be enforced therein.

From a practical standpoint, this raises important questions: Should one train multiple models from scratch to construct a DE, or is it preferable to use the rich feature set of a pretrained model and achieve diverse predictions through a repulsive last-layer ensemble in function-space (fs-RLL-E)? Setting efficiency consideration aside, can a pretrained model with repulsive heads match (or even surpass) a DE in terms of uncertainty estimation?

## Contributions

- We propose a parameter-efficient version of POVI via repulsive last-layer ensembles (Section 4.1). Diverse predictions are ensured by repulsion of the parameters or choosing an appropriate set of repulsion samples for the function space repulsion term (Section 4.2).
- We show that our method can be applied post hoc to pretrained networks. If the backbone avoids feature collapse, retraining only the repulsive last-layer ensemble is sufficient to obtain meaningful uncertainty estimates (Section 4.3).
- We evaluate the approach on synthetic and real-world tasks, including regression, classification, and distribution shift scenarios. The fs-RLL-E enables uncertainty decomposition, calibrated predictions, and improved OOD detection. On pretrained architectures, it matches or surpasses the uncertainty performance of full deep ensembles (Section 6).

## 2 Background

We consider supervised learning tasks. Let  $\mathcal{D} = \{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^N = (\mathbf{X}, \mathbf{Y})$  denote the training data set consisting of  $N$  i.i.d. data samples with inputs  $\mathbf{x}_i \in \mathcal{X}$  and targets  $\mathbf{y}_i \in \mathcal{Y}$ . We define a likelihood model  $p(\mathbf{y}|\mathbf{x}, \theta)$  with the mapping  $f(\cdot; \theta) : \mathcal{X} \rightarrow \mathbb{R}^K$  parameterized by a neural network (NN). Bayesian neural networks (BNNs) treat the parameters  $\theta$  as random variables rather than deterministic values. A prior  $p(\theta)$  is specified, and Bayes’ theorem yields the posterior

$$p(\theta|\mathcal{D}) \propto p(\theta)p(\mathbf{Y}|\mathbf{X}, \theta).$$

<sup>1</sup>Modification of input samples such that the original labels do not apply, e.g., shuffling of random image patches.

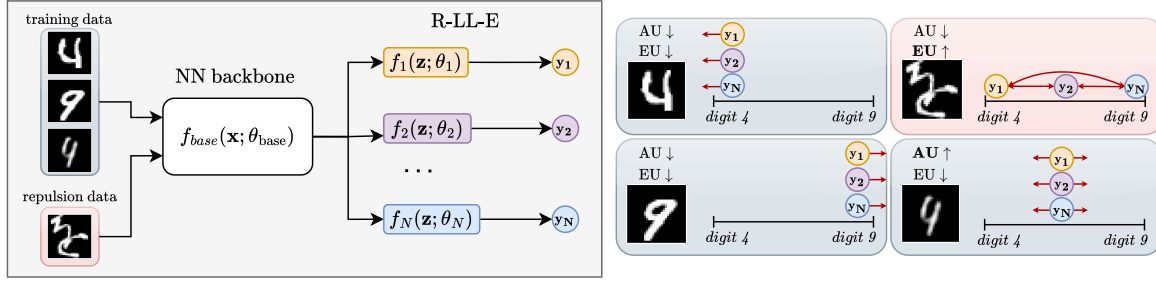


Figure 1: Repulsive last-layer ensemble in function-space (fs-RLL-E), with  $N$  particles. Colored dots correspond to particle predictions. Unlabeled OOD data is used as repulsion samples for the function space repulsion loss. Epistemic uncertainty (EU) is the lowest, when all particle predictions agree, and increases with the spread of the particles. Aleatoric uncertainty (AU) rises for ambiguous samples, e.g., the bottom-right digit belonging to both classes.

Predictions for a new input  $\mathbf{x}_*$  are obtained by marginalizing over  $\theta$

$$p(\mathbf{y}_*|\mathbf{x}_*, \mathcal{D}) = \int p(\mathbf{y}_*|\mathbf{x}_*, \theta) p(\theta|\mathcal{D}) d\theta.$$

This integral, however, is generally intractable and motivates approximate inference methods.

**Particle-optimization variational inference** Variational inference approximates the posterior  $p(\theta|\mathcal{D})$  by a simpler parametric distribution  $q(\theta)$ . POVI methods (Liu & Wang, 2016; Chen et al., 2018) aim to provide more flexibility by considering a non-parametric distribution, specified by a discrete set of particles  $\{\theta^{(i)}\}_{i=1}^n$  according to  $q(\theta) \approx \frac{1}{n} \sum_{i=1}^n \delta(\theta - \theta^{(i)})$ , where  $\delta(\cdot)$  is the Dirac function. The particles can then be optimized iteratively via

$$q(\theta) \approx \frac{1}{n} \sum_{i=1}^n \delta(\theta - \theta^{(i)}),$$

where  $\epsilon_l$  is the step size at time step  $l$ . By viewing the particle optimization as a gradient flow in Wasserstein space, D’Angelo & Fortuin (2021) derive the following update rule that decomposes into an attraction and repulsion term

$$\mathbf{v}(\theta^{(i)}) = \underbrace{\nabla_{\theta^{(i)}} \log p(\theta^{(i)}|\mathcal{D})}_{\text{ATTRACTION}} - \underbrace{\frac{\sum_{j=1}^n \nabla_{\theta^{(i)}} k(\theta^{(i)}, \theta^{(j)})}{\sum_{j=1}^n k(\theta^{(i)}, \theta^{(j)})}}_{\text{REPULSION}}, \quad (1)$$

where  $k(\cdot, \cdot)$  denotes a kernel function. The attraction term drives particles into high-density regions of the posterior distribution, while the repulsion term induces diversity by preventing particles from collapsing into the same optimum. For a single particle, this training procedure reduces to maximum a posteriori (MAP) training; for  $n \rightarrow \infty$  and a properly defined kernel, it converges to the true posterior distribution (D’Angelo & Fortuin, 2021).

**Why repulsion matters for a finite number of particles** The convergence guarantee to the true posterior distribution  $p(\theta|\mathcal{D})$  holds only in the limit of infinitely many particles. Although fascinating from a theoretical perspective, practical importance lies in the analysis of the behavior for a finite number of particles.

Typically, the disagreement between the predictions of ensemble members is used for uncertainty estimation. Following Depeweg et al. (2018), predictive uncertainty can be decomposed into aleatoric and epistemic components, represented as conditional entropy  $\mathbb{H}$  and mutual information  $\mathbb{I}$ :

$$\underbrace{\mathbb{H}[\mathbb{E}_{p(\theta|\mathcal{D})} p(\mathbf{y}|\mathbf{x}, \theta)]}_{\text{Total uncertainty}} = \underbrace{\mathbb{E}_{p(\theta|\mathcal{D})} \mathbb{H}[p(\mathbf{y}|\mathbf{x}, \theta)]}_{\text{Aleatoric}} + \underbrace{\mathbb{I}[\mathbf{y}; \theta|\mathbf{x}, \mathcal{D}]}_{\text{Epistemic}}.$$

Aleatoric uncertainty represents the variability in outcomes due to inherent randomness in the data, while epistemic uncertainty reflects model uncertainty due to limited data. Recent work has raised concerns about the validity of this decomposition (Wimmer et al., 2023). Still, it remains widely used in literature (Kirsch, 2024b) and serves as a practical measure for epistemic uncertainty.

For a finite particle approximation, the epistemic term reduces to a Monte Carlo estimate

$$\mathbb{I}[\mathbf{y}; \theta | \mathbf{x}, \mathcal{D}] \approx \frac{1}{n} \sum_{i=1}^n D_{\text{KL}} \left( p(\mathbf{y} | \mathbf{x}, \theta^{(i)}) \left\| \frac{1}{n} \sum_{j=1}^n p(\mathbf{y} | \mathbf{x}, \theta^{(j)}) \right\| \right). \quad (2)$$

If a test sample  $\mathbf{x}$  is explained by many disagreeing models  $p(\mathbf{y} | \mathbf{x}, \theta)$  under the posterior  $p(\theta | \mathcal{D})$ , epistemic uncertainty (EU) is high. Adding training data near  $\mathbf{x}$  reduces the space of plausible models and their disagreement.

Given practical constraints on the number of particles (typically five to ten), many posterior modes remain unexplored and the estimate of the epistemic uncertainty is shaped largely by a small number of posterior modes. This limitation stresses the need for guiding particles towards representative and *diverse posterior modes* to avoid underestimation of epistemic uncertainty. Deep ensembles can be viewed as an unregularized case of Eq. (1), lacking a repulsion term. Diversity stems from random weight initialization and often fails to capture truly distinct predictive functions, in some cases assigning lower EU to out-of-distribution inputs than a single model’s aleatoric uncertainty (Xia & Bouganis, 2022; Schweighofer et al., 2023b).

POVI methods address this via a repulsion kernel  $k(\theta^{(i)}, \theta^{(j)})$  to avoid mode collapse. In the infinite particle limit, this ensures convergence to the posterior distribution (D’Angelo & Fortuin, 2021; Wild et al., 2023). To improve finite-particle EU estimation, we propose two desiderata:

- D1** *The repulsion term should steer particles towards diverse posterior modes, which provide a useful approximation for the epistemic uncertainty in Eq. (2).*
- D2** *Particles should reach diverse posterior modes from the same initial parameters through the use of the repulsion term. This enables the fine-tuning of pretrained models to better approximate epistemic uncertainty.*

### 3 Where should we enforce diversity?

**Parameter space.** In BNNs, inference is often performed over parameters  $\theta$ , with repulsion via an  $\ell_2$ -based kernel (Wang et al., 2019; D’Angelo & Fortuin, 2021):

$$k(\theta^{(i)}, \theta^{(j)}) = \exp \left( - \frac{\|\theta^{(i)} - \theta^{(j)}\|_2^2}{\nu^2} \right),$$

where  $\nu$  is a hyperparameter, often set by the median heuristic. However, such regularization does not guarantee *diverse predictions* (D1). Overparameterization allows distinct  $\theta$  to yield identical outputs due to permutation and scaling symmetries (Pourzanjani et al., 2017). Moreover, in high-dimensional parameter space, Euclidean distances become less informative (Aggarwal et al., 2001), allowing particles to be far apart in  $\theta$  yet functionally equivalent.

**Function space.** To target predictive diversity, function space (fs) POVI (Wang et al., 2019; D’Angelo & Fortuin, 2021) operates directly on particle predictions  $f^{(i)}(\mathcal{X})$ , with the repulsion term

$$k(f^{(i)}, f^{(j)}) = \exp \left( - \frac{\|f^{(i)}(\mathcal{X}) - f^{(j)}(\mathcal{X})\|_2^2}{\nu^2} \right),$$

enforcing diverse outputs rather than parameters. Particles are updated via

$$f_{l+1}^{(i)}(\mathcal{X}) \leftarrow f_l^{(i)}(\mathcal{X}) + \epsilon_l v(f_l^{(i)}(\mathcal{X})),$$

which steers them toward distinct posterior modes with high KL in Eq. (2). However, to solve the problem we must rely on gradient based optimization procedures that in turn require a parameterized representation of the particles.

First, each particle  $f^{(i)}(\mathcal{X})$  is represented by a specific parameterization  $f^{(i)}(\mathcal{X}; \theta^{(i)})$ . The parameterization  $f^{(i)}(\mathcal{X}; \theta^{(i)})$  must be sufficiently flexible to effectively approximate the underlying function space (Wang et al., 2019).

Moreover, it remains prohibitive to evaluate  $f^{(i)}(\mathcal{X}; \theta^{(i)})$  across the entire input domain  $\mathcal{X}$ . Instead, prior work (Wang et al., 2019) adopted a mini-batch approximation, where the evaluation over the full set  $\mathcal{X}$  is replaced with  $B$  *repulsion samples* drawn from an arbitrary distribution  $\mathbf{x}_{rep} \sim \mu$  with support on  $\mathcal{X}^B$ . The variational distribution is shown to converge to the true posterior if the posterior can be determined by almost all  $B$ -dimensional marginals  $\{p(f(\mathbf{x})|\mathbf{X}, \mathbf{Y}) : \mathbf{x} \in \text{supp}(\mu)\}$  (Wang et al., 2019).

## 4 Improving function space approximations: Practical choices and implications

Function space inference mitigates issues of overparameterization and parameter identifiability (Kirsch, 2024a) by enforcing diversity directly on the predictions, reducing underestimation of uncertainty. However, empirical gains over unregularized DEs have often been limited, especially on large-scale image tasks (Wang et al., 2019; D’Angelo & Fortuin, 2021; Trinh et al., 2023). We argue that this is largely due to suboptimal practical choices when approximating the function space. Below, we revisit three key design decisions and provide improvements that satisfy our desiderata D1 and D2.

### 4.1 Choice of function parameterization

**Problem.** Prior fs-POVI work (Wang et al., 2019; D’Angelo & Fortuin, 2021) parameterizes each particle as a separate neural network. On large-scale tasks, training and storing multiple full networks is computationally expensive, and parallel training is required due to the repulsion term.

**Our approach.** We propose a *shared base network with multiple prediction heads*:

$$f^{(i)}(\mathbf{x}; \theta_{\text{base}}, \theta_{\text{head}}^{(i)}) = f_{\text{head}}^{(i)}(f_{\text{base}}(\mathbf{x}; \theta_{\text{base}}); \theta_{\text{head}}^{(i)}).$$

The base network  $f_{\text{base}}$  provides a shared latent representation; diversity is enforced across the heads  $f_{\text{head}}^{(i)}$ .

**Justification.** Multi-headed network architectures have been used successfully to distill DEs and replicate their functional behavior (Tran et al., 2020), demonstrating sufficient flexibility of single networks (Hinton et al., 2015). Performing particle optimization in function space mitigates the need for training a full DE prior to distillation. The use of a shared deterministic base network aligns with partially stochastic BNNs, where a subnetwork of the parameters is treated probabilistically. Most prominently, Bayesian last-layer networks are employed as practical means to reduce computational demands (Sharma et al., 2023).

### 4.2 Choice of repulsion samples

**Problem.** Evaluation of the function-space repulsion term requires selecting a set of *repulsion samples*  $\mathbf{x}_{rep} \in \mathcal{D}_{rep}$ . This choice impacts the valid input domain for uncertainty estimates of fs-POVI methods. Prior work proposed drawing repulsion samples from the kernel density estimate over the training data (Wang et al., 2019), or from the training data directly (D’Angelo & Fortuin, 2021), which ties the BNN approximation to the training distribution and thus does not guarantee reliable uncertainty estimates in OOD settings. In high-dimensional spaces, drawing random samples from the entire input domain is infeasible, requiring restriction to an informative subset that covers the domain of interest.

**Our approach.** We use *unlabeled OOD data* when available, or apply label-destroying augmentations to the training data. In image classification, this often includes natural images from varying distributions. We can thus exploit the abundance of available unlabeled image data. For example, using kMNIST as repulsion samples for models trained on MNIST, or Textures for models trained on CIFAR10/100, leverages natural variability across different image sets. If unlabeled OOD data is unavailable, repulsion samples can

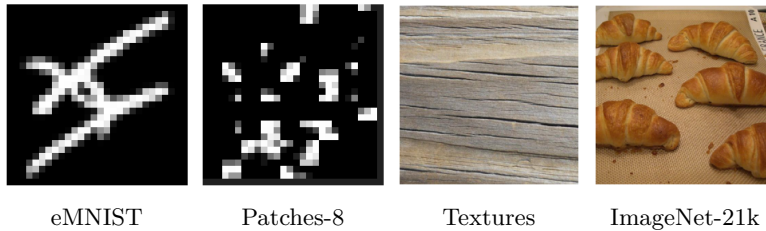


Figure 2: Examples of repulsion samples.

be generated from the training data by label-destroying data augmentation techniques. One such method is the random shuffling of image patches to destroy the shape information of objects that is crucial for human perception (see Fig. 2). The effectiveness of repulsion samples depends on their informativeness and domain coverage. If the samples are too close to the training distribution, they may fail to induce meaningful diversity and can degrade in-domain accuracy. If they are too far from the data manifold, such as random noise, they may have little or no effect on the model’s behavior. Selecting effective repulsion samples therefore requires some knowledge of what constitutes in-distribution and out-of-distribution inputs. This knowledge enables us to guide diversity in function space in a more targeted way, compensating for the lack of random weight initialization.

**Justification.** Enforcing diversity directly on the training data has been shown to degrade performance by artificially inflating epistemic uncertainty at data points where independent training would yield confident predictions Abe et al. (2022); Jeffares et al. (2024). This approach often fails to detect OOD data, which may be characterized by spurious features present in the training data or features that are completely absent in the training set. Using unlabeled OOD data as repulsion samples provides an effective solution to these challenges. These samples may contain features that are spurious or absent in the training data, allowing the models to meaningfully enforce diversity and improve OOD detection capabilities. Similarly, label-destroying data augmentation mitigates robust features that are indicative of the class label. Compared to methods that rely on feature density to detect OOD data, repulsion samples offer the benefit of learning to ignore spurious features that may be present in the training data.

### 4.3 Post-hoc uncertainties for pretrained models

There are many scenarios where using a pretrained model is desirable. We may wish to equip an existing model with post-hoc uncertainty estimates without retraining from scratch, or finetune a base model pretrained on large datasets for a smaller target dataset. Pretraining yields strong, generalizable features, improves data efficiency, and speeds convergence. While some attribute these gains mainly to faster optimization (He et al., 2019), others find improved calibration even for single models (Hendrycks et al., 2019a). However, when the objective is to capture epistemic uncertainty, pretraining counteracts the benefit of deep ensembles that rely on random initializations with independent training to achieve diversity.

In such cases, multi-headed networks provide a principled way to estimate post-hoc uncertainties for the pretrained base. While LL-E approaches have been used before (Schweighofer et al., 2023b), they often lack diversity due to the shared base model. The repulsion term allows us to reintroduce this diversity and decouple representation learning (via the base model) from uncertainty estimation (via the diverse last-layer ensemble). If OOD-relevant features are mapped to ID features in feature space, known as feature collapse (van Amersfoort et al., 2020), the ensemble heads may not achieve the desired diverse predictions. Fortunately, many modern architectures include mechanisms that mitigate collapse.

**Spectral normalization and residual connections.** Distance-aware representations can be encouraged by imposing bi-Lipschitz constraints  $K_L d_I(\mathbf{x}_1, \mathbf{x}_2) \leq d_F(f_{\text{base}}(\mathbf{x}_1), f_{\text{base}}(\mathbf{x}_2)) \leq K_U d_I(\mathbf{x}_1, \mathbf{x}_2)$ , where  $d_I$  and  $d_F$  denote distances in input and feature space, and  $K_L$ ,  $K_U$  are lower and upper Lipschitz constants. These constraints bound how input distances translate into feature distances. Spectral normalization and residual connections are effective means to impose such bounds (Miyato et al., 2018). While most pretrained models are not trained with spectral normalization, residual architectures alone often preserve distance-awareness.

**Pretraining.** Large-scale pretraining, whether supervised, self-supervised, or contrastive (Krizhevsky et al., 2012), produces features that generalize well (Hendrycks et al., 2019a) and promote semantic separation in feature space. This mitigates collapse, but OOD inputs may still activate spurious features learned during pretraining, leading to overconfident errors. Repulsion in function space can counter this by encouraging prediction-head disagreement on such inputs.

## 5 Related Work

**(Repulsive) Deep Ensembles.** Deep ensembles (Lakshminarayanan et al., 2017) often outperform BNNs in accuracy, calibration, and OOD detection (Gustafsson et al., 2020; Ovadia et al., 2019). Repulsive variants promote diversity via kernelized losses (Wang et al., 2019; D’Angelo & Fortuin, 2021), feature/gradient differences (Yashima et al., 2022; Trinh et al., 2023), or disagreement on unlabeled data (Pagliardini et al., 2022). Weight-distribution entropy maximization (de Mathelin et al., 2025) also has improved epistemic uncertainty.

**Function-space inference.** A number of inference methods for BNNs consider the shift from inference in the space of network parameters to the function space Sun et al. (2019); Ma et al. (2019); Burt et al. (2020); Wang et al. (2019); Ma & Hernández-Lobato (2021); Rudner et al. (2022). This allows to specify meaningful prior distributions over the network parameters. Recent work proposes tractable VI via local linearization (Rudner et al., 2022; 2023). POVI methods approximate the posterior distribution using a set of discrete particles to capture its multimodal structure Wang et al. (2019); D’Angelo & Fortuin (2021).

**Auxiliary out-of-distribution data.** Function space inference methods enforce the function prior on a set of input points, in some work referred to as measurement (Sun et al., 2019; Wang et al., 2019; Ma & Hernández-Lobato, 2021) or context samples (Rudner et al., 2022; 2023). In low-dimensional problems, such samples can be obtained by drawing from a distribution with support over the domain of interest (Sun et al., 2019; Wang et al., 2019; Ma & Hernández-Lobato, 2021). For high-dimensional problems with structured data, such as natural images, samples from an OOD data set have shown improvements (Rudner et al., 2022; 2023). Related work on OOD detection methods for single networks has used auxiliary OOD datasets to maximize softmax entropy (Hendrycks et al., 2019b).

**Multi-headed architectures.** Multi-headed networks reduce memory requirements by sharing a backbone (Song & Chai, 2018; Sercu et al., 2016; Lee et al., 2015). Among the first, last-layer ensembles were analyzed by Lee et al. (2015) in terms of accuracy, parameters sharing, and diversity by random parameters versus bagging the data. They have been applied in RL (Osband et al., 2016), distillation (Zhu et al., 2018; Tran et al., 2020), and uncertainty estimation (Valdenegro-Toro, 2023), with diversity encouraged via decorrelation (Zhang et al., 2020; Lee et al., 2022).

**Distance-based uncertainty methods.** These relate epistemic uncertainty to distance from training support, typically in latent space (Charpentier et al., 2020; Postels et al., 2020; Mukhoti et al., 2023; Winkens et al., 2020; Liu et al., 2020; van Amersfoort et al., 2020; Tagasovska & Lopez-Paz, 2019). They require feature regularization (van Amersfoort et al., 2021) via gradient penalties (Gulrajani et al., 2017) or spectral normalization (Miyato et al., 2018). While strong for OOD detection, they often miscalibrate under shift (Postels et al., 2021); our approach benefits from similar regularization but enforces diversity directly in function space.

## 6 Experiments

We evaluate our repulsive last-layer ensemble applied to different base models across several benchmark experiments. Two cases are considered: i) training a base model with random initialization from scratch on the training data, and ii) using pretrained models that have been trained on ImageNet-1k. In the latter case, the pretrained model serves as basis for further optimization. Consequently, standard deep ensembles lose diversity from random weight initialization. In both cases, we aim to assess whether diversity lost due to the shared backbone can be meaningfully recovered by repulsion either in weight space or in function space and if our method can serve as a practical and lightweight uncertainty estimator.

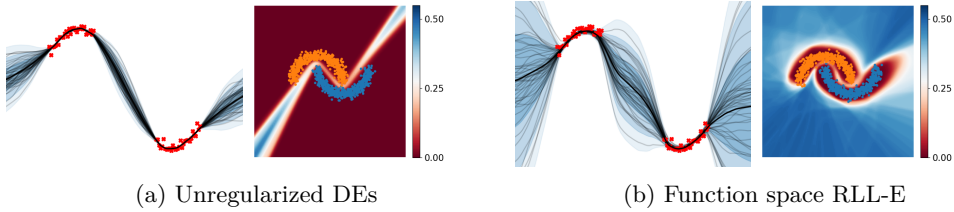


Figure 3: For regression, we show the prediction of individual particles/ensemble members, the mean and the standard deviation. For classification, we show the standard deviation of  $p(\mathbf{y}|\mathbf{x}, \theta)$ . DEs are highly confident in regions distant from training data (low standard deviation is colored in red), while fs-RLL-E predictions are enforced to be diverse outside of the training data domain.

Our experiments are structured as follows: We begin with synthetic toy examples to illustrate the behavior of our multi-head model and its function-space diversity (Section 6.1). We then evaluate the disentanglement of aleatoric and epistemic uncertainty on DirtyMNIST, a dataset with ambiguous labels (Section 6.2). We test calibration under distribution shifts using the CIFAR10-C and CIFAR100-C benchmarks (Section 6.3). Lastly, we evaluate epistemic uncertainty for transfer learning tasks with ImageNet-1k pretrained models (Section 6.4).

In the appendix, we use epistemic uncertainty for active learning on DirtyMNIST (Section B.1), and OOD detection using epistemic uncertainties where we compare base models that have been either (i) trained from scratch on the task or (ii) pretrained on ImageNet-1k (Section B.2).

**Models and baselines.** For from-scratch experiments we use ResNet18 and ResNet50 with spectral normalization; pretrained experiments additionally include ViT-B/16. Baselines comprise a single MAP-trained network, deep ensembles (DE-5), deterministic uncertainty (DDU), last-layer Laplace approximation (LL-Laplace), and non-repulsive last-layer ensembles (LL-E). Our variants apply repulsion in parameter space (RLL-E) or function space (fs-RLL-E), using identical backbones for fair comparison.

**Metrics.** Performance is measured via accuracy and negative log-likelihood (NLL). Calibration is quantified with the expected calibration error (ECE) (Naeini et al., 2015). For OOD detection, we report the area under the receiver operating characteristic curve (AUROC) and the area under the precision-recall curve (AUPR), using the specified epistemic uncertainty estimates to distinguish between ID and OOD samples. We use mutual information (Equation 2) for ensemble-based methods, maximum softmax probability (MSP)/maximum logit/entropy for the single base model, and GMM density for DDU.

## 6.1 Synthetic data

On two toy examples, we illustrate the effectiveness of the multi-head architecture as a lightweight parameterization and inference in function space. We estimate the epistemic uncertainty for a one-dimensional regression and a two-dimensional classification problem using full DEs with 30 ensemble members and fs-RLL-E. A feed-forward neural network with 3 hidden layers and 128 neurons is used as the base model. The repulsive head consists of 30 particles with a linear layer. Results are shown in Figure 3. Deep ensemble predictions show low uncertainty far from the training data. By performing particle inference in function space, we can enforce diverse predictions outside the training distribution even with a simpler network structure.

## 6.2 Disentangling aleatoric and epistemic uncertainty

**Setup** We test uncertainty estimation using DirtyMNIST, a dataset containing 60k clean MNIST digits and 60k ambiguous digits with multiple valid labels (high aleatoric uncertainty). Our repulsive last-layer ensemble aims to model epistemic uncertainty by learning multiple hypotheses that agree on training data but disagree elsewhere. Epistemic uncertainty is generally difficult to evaluate empirically. A widely accepted assumption is that epistemic uncertainty should increase on OOD inputs (de Mathelin et al., 2025). Therefore, we use OOD detection as a proxy task. This involves computing uncertainty scores for ID and OOD samples, and



Table 1: **Uncertainty decomposition on DirtyMNIST.** Accuracy, NLL, and ECE are reported for the test split of DirtyMNIST. For distinguishing clean, ambiguous and OOD samples we report the AUROC scores. All results are averaged over 5 runs. Best results are **bold**, second-best are underlined.

Method	Uncertainty	DirtyMNIST			Clean MNIST vs OOD	Ambig. MNIST vs OOD
		Acc. [%] $\uparrow$	NLL $\downarrow$	ECE [%] $\downarrow$	AUROC /AUPR [%] $\uparrow$	AUROC /AUPR [%] $\uparrow$
MAP	MSP	80.78	0.5591	2.09	97.35 / 96.47	70.01 / 46.04
	Max Logit				98.58 / 98.68	73.75 / 52.88
	Entropy				97.57 / 96.82	73.76 / 52.08
DDU	GMM density	80.78	0.5591	2.09	<b>99.15</b> / <b>99.80</b>	<b>99.38</b> / <b>99.61</b>
DDU (diag)	GMM density				97.35 / 98.98	<u>96.47</u> / 96.25
LL-Laplace	MI				97.94 / 99.31	89.83 / 88.94
LL-E	MI				98.46 / <u>99.66</u>	92.01 / 94.79
RLL-E (ours)	MI	<u>83.51</u>	0.4839	1.37	<u>98.63</u> / 99.37	93.32 / 93.82
fs-RLL-E (ours)						
+ DIRTYMNIST	MI	<b>83.54</b>	<b>0.4834</b>	1.31	96.74 / 97.41	48.44 / 29.41
+ kMNIST	MI	82.97	0.4961	<b>0.92</b>	97.76 / 99.55	95.92 / 97.25
+ PATCHES-8	MI	83.04	0.4932	<u>1.20</u>	97.42 / 99.44	96.33 / <u>97.56</u>
DE-5	MI	83.32	0.5024	5.20	97.70 / 98.43	88.80 / 79.07

analyzing how well the two classes can be separated. As OOD datasets we use FashionMNIST (Xiao et al., 2017), EMNIST (Cohen et al., 2017), and Omniglot (Lake et al., 2015).

We evaluate: (i) classification and calibration on DirtyMNIST (Acc., NLL, ECE); (ii) OOD detection between clean MNIST and EMNIST, FashionMNIST, Omniglot, where aleatoric uncertainty often suffices; and (iii) OOD detection between ambiguous MNIST and the same datasets, requiring epistemic uncertainty.

All post-hoc methods (DDU, LL-Laplace, LL-E, RLL-E, fs-RLL-E) use the feature space of a ResNet18 MAP model. For LL-E variants, we reinitialize the last linear layer with 10 heads, freeze the backbone, and train with AdamW (LR=0.001). Repulsion is applied in parameter space (RLL-E) or function space (fs-RLL-E) with specified repulsion samples.

**Results** Table 1 summarizes the results across all tasks. First, we evaluate in-distribution performance on DirtyMNIST (Accuracy, NLL, ECE). The LL-E and our repulsive variants (RLL-E, fs-RLL-E) mostly improve upon the single-model baseline (MAP). However, the main advantage of our approach becomes evident in the OOD detection tasks. The AUROC and AUPR scores reflect how well the specified uncertainty estimator separates ID and OOD samples across all thresholds. Ideally, epistemic uncertainty should be small on ID samples and large on unseen OOD data.

On the clean-vs-OOD task, even using the aleatoric uncertainty (Entropy) of the single-model (MAP) as OOD detection scores performs reasonably well with an AUROC score of 97.57%. Nonetheless, our RLL-E extension improves detection further (98.63%) and even outperforms the full DE-5 (97.7%), despite requiring significantly fewer parameters and less training time. The more challenging ambiguous-vs-OOD task reveals the limitations of the single-model baseline. Here, the aleatoric uncertainty detector degrades to AUROC values of 73.76%. Our RLL-E still achieves 93.32%, outperforming the LL-Laplace and the full DE-5. Repulsion in function space using samples from the training distribution itself (DirtyMNIST) leads to degraded OOD detection. In contrast, using kMNIST or training data with label-destroying augmentations (such as Patches-8) introduces beneficial diversity. This results in improved separation of ambiguous and OOD samples, with AUROC values increasing to 95.92% and 96.33%, respectively.

Among all methods, only DDU performs better on OOD detection (99.38%). However, DDU follows a fundamentally different strategy by modeling the latent space using class-conditional Gaussian distributions. In high-dimensional feature spaces with many classes, the required covariance matrices become memory-intensive. We therefore also compare against a diagonalized version of DDU, which matches the performance of our fs-RLL-E with 96.47%. In addition, DDU performance may degrade when the latent features are not clustered such that they fit the classwise Gaussian assumption. We revisit this aspect in Section 6.4 when evaluating pretrained models as a backbone.

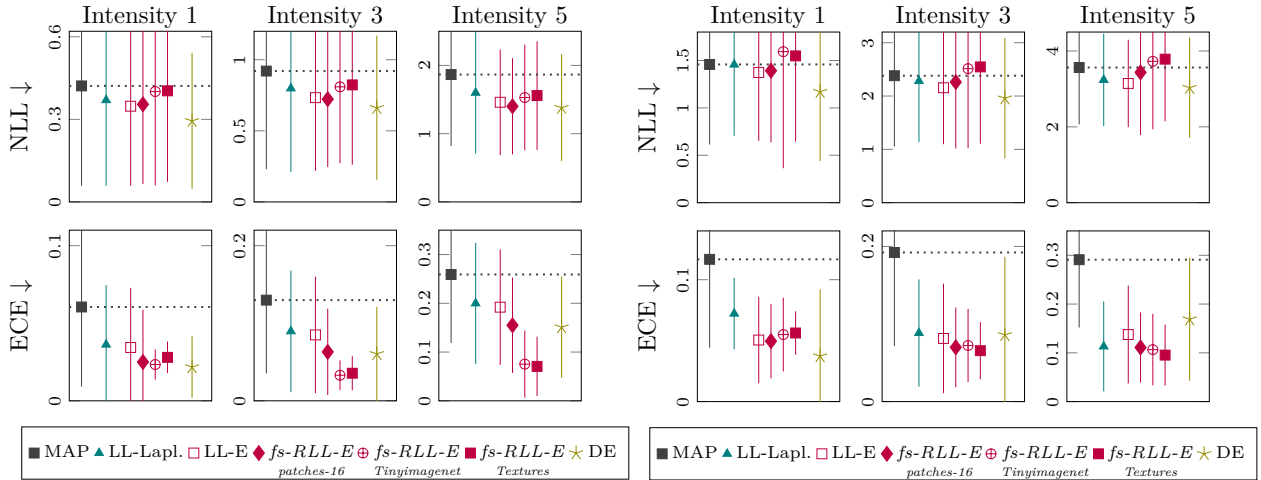


Figure 4: NLL and uncertainty calibration of the different methods on CIFAR10-C (left) and CIFAR100-C (right), for different levels of corruption intensity (columns), averaged over all corruption types. By retraining the last linear layer only, our method fs-RLL-E improves NLL, and achieves similar ECE scores as DEs.

### 6.3 Covariate shift calibration

**Setup** We evaluate how uncertainty estimators behave under covariate shift using the CIFAR10-C and CIFAR100-C benchmarks Hendrycks & Dietterich (2019). These datasets apply 19 types of common corruptions (e.g., blur, noise, compression), each at 5 severity levels, to simulate real-world data shifts. We use the corrupted test sets and report uncertainty calibration in terms of NLL and ECE, averaged over all corruption types.

**Results** Figure 4 summarizes the calibration performance under corruption. Retraining the last layer without regularization (LL-E) already improves the calibration of the single model (MAP) for both CIFAR10-C and CIFAR100-C. On CIFAR10-C, our function-space repulsion (fs-RLL-E) yields further improvements, outperforming LL-Laplace, and achieving calibration competitive with DEs, while requiring fewer parameters.

The improvement can be attributed to better feature selection in the presence of corruptions. If the final layer relies on spurious or fragile features, corrupted inputs may trigger high-confidence but incorrect predictions. By enforcing disagreement on repulsion samples, our method promotes diversity in predictive behavior and reduces over-reliance on non-robust features, resulting in better calibration.

### 6.4 Transfer learning with pretrained models

**Setup** In this task, we evaluate how well our method performs in transfer learning scenarios using high-resolution image classification tasks. Specifically, we consider the Food101 and Stanford Cars datasets, which are common benchmarks for fine-grained recognition. We use Resnet18, Resnet50, and ViT-B/16 models pretrained on ImageNet-1k. Note, that we do not train any model from scratch and thus do not compare to full deep ensembles. Instead we use LL-E, LL-Laplace and DDU as competitive uncertainty estimation methods. For fine-tuning, we jointly optimize the backbone model with the ensemble head by using separate learning rates of  $10^{-3}$  for the ensemble head and  $10^{-5}$  for the backbone model. For function-space repulsion, we sample mini-batches from Textures and ImageNet-21k as repulsion datasets. We test whether semantic and visual diversity leads to improved diversity. As OOD datasets we use SVHN (Netzer et al., 2011), Places365 (Zhou et al., 2017), ImageNet-O (Hendrycks et al., 2021), Country211 (Radford et al., 2021), and TinyImageNet (Le & Yang, 2015).

Table 2: **In-distribution performance (Accuracy [%] / Negative Log Likelihood (NLL))**. We compare base models that are *fine-tuned from ImageNet-1k pretrained weights*.

Method	Food101			Stanford Cars		
	ResNet18 (Pretrained)	ResNet50 (Pretrained)	ViT-B/16 (Pretrained)	ResNet18 (Pretrained)	ResNet50 (Pretrained)	ViT-B/16 (Pretrained)
MAP	75.85 / 0.9357	82.44 / 0.7216	83.03 / 0.7887	64.68 / 1.3272	69.20 / 1.1758	80.62 / 0.7079
LL-Laplace	75.66 / 0.9083	82.06 / 0.7795	84.35 / 0.6400	61.01 / 1.7946	56.67 / 2.4190	74.67 / 1.6946
LL-E	75.81 / 0.9347	82.54 / 0.7218	83.31 / 0.7846	66.57 / 1.2438	69.76 / 1.1539	81.19 / 0.6909
RLL-E ( <i>ours</i> )	75.87 / 0.9292	82.67 / 0.7156	82.79 / 0.7987	66.60 / 1.2426	69.96 / 1.1496	81.12 / 0.6932
fs-RLL-E ( <i>ours</i> )						
+ TEXTURES	75.89 / 0.9342	82.58 / 0.7227	82.81 / 0.8031	66.54 / 1.2426	70.05 / 1.1538	81.20 / 0.6916
+ IMAGENET-21K	73.68 / 1.0386	81.26 / 0.6802	83.07 / 0.6436	66.52 / 1.2443	69.97 / 1.1503	81.17 / 0.6913

Table 3: **Semantic shift detection (AUROC [%] / AUPR [%])**. We compare base models that are *fine-tuned from ImageNet-1k pretrained weights*. All results are averaged over 5 runs. Best results for each base model are **bold**, second-best are underlined.

Method	Uncertainty	Food101			Stanford Cars		
		ResNet18 (Pretrained)	ResNet50 (Pretrained)	ViT-B/16 (Pretrained)	ResNet18 (Pretrained)	ResNet50 (Pretrained)	ViT-B/16 (Pretrained)
MAP	MSP	93.55 / 83.70	90.04 / 76.93	94.74 / 86.25	86.99 / 84.77	92.14 / 92.64	98.90 / 99.10
	Max Logit	71.38 / 45.76	38.84 / 31.20	99.11 / 96.82	67.46 / 65.73	10.82 / 44.88	99.56 / 99.67
	Entropy	95.53 / 88.34	92.71 / 83.13	96.96 / 92.06	90.31 / 88.34	95.37 / 95.77	99.63 / 99.77
LL-Laplace	MI	89.08 / 77.56	<u>99.09</u> / <b>96.79</b>	<b>99.41</b> / <b>98.02</b>	95.81 / 97.16	95.27 / 93.72	99.42 / 99.31
DDU (diag)	GMM density	50.27 / 35.30	86.82 / 76.99	98.68 / 96.52	50.27 / 35.3	67.01 / 71.96	99.90 / 99.95
LL-E	MI	87.11 / 70.54	96.19 / 90.96	98.36 / 95.37	93.95 / 95.58	96.37 / 95.74	99.67 / 99.81
RLL-E ( <i>ours</i> )	MI	90.54 / 74.80	97.33 / 91.25	98.66 / 95.78	94.56 / 95.96	98.51 / 98.0	99.79 / 99.88
fs-RLL-E ( <i>ours</i> )							
+ TEXTURES	MI	95.70 / 90.03	97.93 / 94.43	98.83 / 96.76	99.23 / 99.48	98.87 / 98.98	99.87 / 99.93
+ IMAGENET-21K	MI	<b>98.25</b> / <b>93.61</b>	<b>99.33</b> / <u>96.40</u>	<u>99.40</u> / <u>97.98</u>	<b>99.68</b> / <b>99.56</b>	<b>99.63</b> / <b>99.28</b>	<b>99.95</b> / <b>99.96</b>

**Results** Table 2 and 3 summarize the ID performance and OOD detection, respectively. Again, we do not expect significant improvements on the ID data in terms of accuracy and NLL scores with the repulsion terms. Our main objective is to estimate epistemic uncertainty, and thus learn ensemble heads which agree on the training data but disagree for distribution shifts. In Table 2 we want to emphasize that the repulsion terms do not harm accuracy and calibration for in-distribution data.

Clear benefits emerge in the OOD detection tasks (Table 3). While we do not evaluate on the Textures dataset, it provides valuable repulsion samples that improve separation of OOD data, particularly for ResNet18. For example, on Food101, fs-RLL-E improves AUROC from 87.11% (LL-E) to 95.7%, and AUPR from 70.54% to 90.03%. On Stanford Cars, the improvement is even stronger: fs-RLL-E achieves 99.23 / 99.48, compared to 93.95 / 95.58 for LL-E. Changing the repulsion samples to the more diverse ImageNet-21k dataset further increases AUROC across all datasets and base models.

For better-performing backbones, ResNet50 and ViT-B/16, the effect of repulsion is still positive but less pronounced. On Stanford Cars, the single ViT-B/16 model (MAP) already achieves strong OOD scores (e.g., 99.63 / 99.77 for Entropy). Nonetheless, fs-RLL-E with repulsion on samples from the Textures dataset improves this further to 99.87 / 99.93, and similar gains are seen for ResNet50 (from 95.37 / 95.77 with Entropy to 98.87 / 98.98 with fs-RLL-E). For ViT-B/16 on Food101, the AUROC increases from 96.96% (Entropy) to 98.83% with fs-RLL-E.

LL-Laplace performs best on Food101 with ResNet50 and ViT-B/16, outperforming other methods, but underperforms on Stanford Cars. DDU (with diagonal covariance) shows mixed performance: it excels on ViT-B/16 (e.g., 99.90% / 99.95% on Stanford Cars), but performs poorly on convolutional models like ResNet18 and ResNet50 (e.g., 50.27% / 35.30% on Food101). We attribute this to (i) its reliance on class-

conditional Gaussian assumptions, which may not hold for CNNs, and (ii) the limitations of the diagonal covariance approximation.

Overall, we observe that LL-Es mostly improves over the single-model baseline, with RLL-E and fs-RLL-E offering further benefits. While repulsion using the Textures dataset already improves performance, using a more diverse dataset such as ImageNet-21k for repulsion samples further improves AUROC.

## 7 Conclusion

We have shown that particle optimization in function space is not limited to DE architectures. A significant number of parameters can be saved by exploring different network architectures to parameterize the function space. We proposed a hybrid approach using a multi-headed network where the shared base model extracts the features for the repulsive ensemble head. This offers a principled way to equip pretrained networks with post-hoc uncertainty estimates and to incorporate prior functional knowledge into the training procedure. We highlighted the inherent limitations of enforcing diversity on training data alone. By utilizing augmented training data, or unlabeled OOD data, we achieve improvements on OOD detection without harming classification accuracy. We show that fine-tuning pretrained models outperforms deep ensembles trained from scratch in both in-distribution and OOD settings. Post-hoc uncertainty extensions, including our fs-RLL-E, further improve epistemic uncertainty estimation with minimal computational overhead.

### Future work

The function space formulation of the inference problem requires a selection of repulsion samples. If those repulsion samples do not cover the domain of interest during deployment of the model, the function space repulsion term may fail to improve the quality of epistemic uncertainty estimates and their effectiveness in downstream tasks such as OOD detection. As future work, it is of interest to give more rigorous statements about the selection of repulsion samples and the implications for uncertainty estimates. We plan to further explore data augmentation schemes for the generation of repulsion samples for different tasks and their respective limitations. Complementary to this, it would be valuable to investigate how different uncertainty decomposition schemes (Wimmer et al., 2023; Schweighofer et al., 2023a) influence epistemic uncertainty estimates. Additionally, our evaluation is currently limited to classification tasks. In future work, we want to extend our approach to real-world regression tasks.

## References

- Taiga Abe, E Kelly Buchanan, Geoff Pleiss, and John Patrick Cunningham. The best deep ensembles sacrifice predictive diversity. 2022.
- Charu C Aggarwal, Alexander Hinneburg, and Daniel A Keim. On the surprising behavior of distance metrics in high dimensional space. In *Database Theory—ICDT 2001: 8th International Conference London, UK, January 4–6, 2001 Proceedings* 8, pp. 420–434. Springer, 2001.
- Zeyuan Allen-Zhu and Yuanzhi Li. Towards understanding ensemble, knowledge distillation and self-distillation in deep learning. *arXiv preprint arXiv:2012.09816*, 2020.
- Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101—mining discriminative components with random forests. In *European Conference on Computer Vision*, pp. 446–461. Springer, 2014.
- David R Burt, Sebastian W Ober, Adrià Garriga-Alonso, and Mark van der Wilk. Understanding variational inference in function-space. In *Third Symposium on Advances in Approximate Bayesian Inference*, 2020.
- Bertrand Charpentier, Daniel Zügner, and Stephan Günnemann. Posterior network: Uncertainty estimation without ood samples via density-based pseudo-counts. In *Advances in Neural Information Processing Systems*, volume 33, pp. 1356–1367, 2020.
- Changyou Chen, Ruiyi Zhang, Wenlin Wang, Bai Li, and Liqun Chen. A unified particle-optimization framework for scalable bayesian sampling. *arXiv preprint arXiv:1805.11659*, 2018.

- M. Cimpoi, S. Maji, I. Kokkinos, S. Mohamed, and A. Vedaldi. Describing textures in the wild. In *Proceedings of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2014.
- Tarin Clanuwat, Mikel Bober-Irizar, Asanobu Kitamoto, Alex Lamb, Kazuaki Yamamoto, and David Ha. Deep learning for classical japanese literature. *arXiv preprint arXiv:1812.01718*, 2018.
- Gregory Cohen, Saeed Afshar, Jonathan Tapson, and Andre Van Schaik. Emnist: Extending mnist to handwritten letters. In *2017 international joint conference on neural networks (IJCNN)*, pp. 2921–2926. IEEE, 2017.
- Francesco D’Angelo and Vincent Fortuin. Repulsive deep ensembles are Bayesian. In *Advances in Neural Information Processing Systems*, volume 34, pp. 3451–3465, 2021.
- Antoine de Mathelin, François Deheeger, Mathilde Mougéot, and Nicolas Vayatis. Deep out-of-distribution uncertainty quantification via weight entropy maximization. *Journal of Machine Learning Research*, 26(4):1–68, 2025.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 248–255. Ieee, 2009.
- Stefan Depeweg, Jose-Miguel Hernandez-Lobato, Finale Doshi-Velez, and Steffen Udluft. Decomposition of uncertainty in Bayesian deep learning for efficient and risk-sensitive learning. In *International Conference on Machine Learning*, pp. 1184–1193. PMLR, 2018.
- Edrina Gashi, Jiankang Deng, and Ismail Elezi. Deep active learning: A reality check. *arXiv preprint arXiv:2403.14800*, 2024.
- Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville. Improved training of wasserstein GANs. In *Advances in Neural Information Processing Systems*, volume 30, 2017.
- Fredrik K Gustafsson, Martin Danelljan, and Thomas B Schon. Evaluating scalable bayesian deep learning methods for robust computer vision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition workshops*, pp. 318–319, 2020.
- Kaiming He, Ross Girshick, and Piotr Dollár. Rethinking imagenet pre-training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 4918–4927, 2019.
- Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. In *International Conference on Learning Representations*, 2019.
- Dan Hendrycks, Kimin Lee, and Mantas Mazeika. Using pre-training can improve model robustness and uncertainty. In *International Conference on Machine Learning*, pp. 2712–2721. PMLR, 2019a.
- Dan Hendrycks, Mantas Mazeika, and Thomas Dietterich. Deep anomaly detection with outlier exposure. In *International Conference on Learning Representations*, 2019b.
- Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adversarial examples. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15262–15271, 2021.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- Eyke Hüllermeier and Willem Waegeman. Aleatoric and epistemic uncertainty in machine learning: An introduction to concepts and methods. *Machine Learning*, 110:457–506, 2021.
- Alan Jeffares, Tennison Liu, Jonathan Crabbé, and Mihaela van der Schaar. Joint training of deep ensembles fails due to learner collusion. In *Advances in Neural Information Processing Systems*, volume 36, 2024.

- Alex Kendall and Yarin Gal. What uncertainties do we need in Bayesian deep learning for computer vision? In *Advances in Neural Information Processing Systems*, volume 30, 2017.
- Andreas Kirsch. Bridging the data processing inequality and function-space variational inference. In *The Third Blogpost Track at ICLR 2024*, 2024a.
- Andreas Kirsch. (implicit) ensembles of ensembles: Epistemic uncertainty collapse in large models. *arXiv preprint arXiv:2409.02628*, 2024b.
- Jonathan Krause, Jia Deng, Michael Stark, and Li Fei-Fei. Collecting a large-scale dataset of fine-grained cars. 2013.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems*, volume 25, 2012.
- Brenden M Lake, Ruslan Salakhutdinov, and Joshua B Tenenbaum. Human-level concept learning through probabilistic program induction. *Science*, 350(6266):1332–1338, 2015.
- Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In *Advances in Neural Information Processing Systems*, volume 30, 2017.
- Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. *CS 231N*, 7(7):3, 2015.
- Stefan Lee, Senthil Purushwalkam, Michael Cogswell, David Crandall, and Dhruv Batra. Why m heads are better than one: Training a diverse ensemble of deep networks. *arXiv preprint arXiv:1511.06314*, 2015.
- Yoonho Lee, Huaxiu Yao, and Chelsea Finn. Diversify and disambiguate: Out-of-distribution robustness via disagreement. In *International Conference on Learning Representations*, 2022.
- Chang Liu, Jingwei Zhuo, Pengyu Cheng, Ruiyi Zhang, and Jun Zhu. Understanding and accelerating particle-based variational inference. In *International Conference on Machine Learning*, pp. 4082–4092. PMLR, 2019.
- Jeremiah Liu, Zi Lin, Shreyas Padhy, Dustin Tran, Tania Bedrax Weiss, and Balaji Lakshminarayanan. Simple and principled uncertainty estimation with deterministic deep learning via distance awareness. In *Advances in Neural Information Processing Systems*, volume 33, pp. 7498–7512, 2020.
- Qiang Liu. Stein variational gradient descent as gradient flow. In *Advances in Neural Information Processing Systems*, volume 30, 2017.
- Qiang Liu and Dilin Wang. Stein variational gradient descent: A general purpose Bayesian inference algorithm. In *Advances in Neural Information Processing Systems*, volume 29, 2016.
- Chao Ma and José Miguel Hernández-Lobato. Functional variational inference based on stochastic process generators. In *Advances in Neural Information Processing Systems*, volume 34, pp. 21795–21807, 2021.
- Chao Ma, Yingzhen Li, and José Miguel Hernández-Lobato. Variational implicit processes. In *International Conference on Machine Learning*, pp. 4222–4233. PMLR, 2019.
- Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative adversarial networks. *arXiv preprint arXiv:1802.05957*, 2018.
- Jishnu Mukhoti, Andreas Kirsch, Joost van Amersfoort, Philip HS Torr, and Yarin Gal. Deep deterministic uncertainty: A new simple baseline. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 24384–24394, 2023.
- Mahdi Pakdaman Naeini, Gregory Cooper, and Milos Hauskrecht. Obtaining well calibrated probabilities using Bayesian binning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 29, 2015.

- Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Baolin Wu, Andrew Y Ng, et al. Reading digits in natural images with unsupervised feature learning. In *NIPS workshop on deep learning and unsupervised feature learning*, volume 2011, pp. 7. Granada, 2011.
- Ian Osband, Charles Blundell, Alexander Pritzel, and Benjamin Van Roy. Deep exploration via bootstrapped dqn. In *Advances in Neural Information Processing Systems*, volume 29, 2016.
- Yaniv Ovadia, Emily Fertig, Jie Ren, Zachary Nado, David Sculley, Sebastian Nowozin, Joshua Dillon, Balaji Lakshminarayanan, and Jasper Snoek. Can you trust your model’s uncertainty? evaluating predictive uncertainty under dataset shift. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- Matteo Pagliardini, Martin Jaggi, François Fleuret, and Sai Praneeth Karimireddy. Agree to disagree: Diversity through disagreement for better transferability. *arXiv preprint arXiv:2202.04414*, 2022.
- Janis Postels, Hermann Blum, Cesar Cadena, Roland Siegwart, Luc Van Gool, and Federico Tombari. Quantifying aleatoric and epistemic uncertainty using density estimation in latent space. *arXiv preprint arXiv:2012.03082*, 2020.
- Janis Postels, Mattia Segu, Tao Sun, Luca Sieber, Luc Van Gool, Fisher Yu, and Federico Tombari. On the practicality of deterministic epistemic uncertainty. *arXiv preprint arXiv:2107.00649*, 2021.
- Arya A Pourzanjani, Richard M Jiang, and Linda R Petzold. Improving the identifiability of neural networks for bayesian inference. In *NIPS workshop on bayesian deep learning*, volume 4, pp. 31, 2017.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pp. 8748–8763. PMLR, 2021.
- Tim GJ Rudner, Zonghao Chen, Yee Whye Teh, and Yarin Gal. Tractable function-space variational inference in Bayesian neural networks. In *Advances in Neural Information Processing Systems*, volume 35, pp. 22686–22698, 2022.
- Tim GJ Rudner, Sanyam Kapoor, Shikai Qiu, and Andrew Gordon Wilson. Function-space regularization in neural networks: A probabilistic perspective. In *International Conference on Machine Learning*, pp. 29275–29290. PMLR, 2023.
- Kajetan Schweighofer, Lukas Aichberger, Mykyta Ielanskyi, and Sepp Hochreiter. Introducing an improved information-theoretic measure of predictive uncertainty. *arXiv preprint arXiv:2311.08309*, 2023a.
- Kajetan Schweighofer, Lukas Aichberger, Mykyta Ielanskyi, Günter Klambauer, and Sepp Hochreiter. Quantification of uncertainty with adversarial models. In *Advances in Neural Information Processing Systems*, volume 36, pp. 19446–19484, 2023b.
- Tom Sercu, Christian Puhersch, Brian Kingsbury, and Yann LeCun. Very deep multilingual convolutional neural networks for lvcsr. In *2016 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pp. 4955–4959. IEEE, 2016.
- Mrinank Sharma, Sebastian Farquhar, Eric Nalisnick, and Tom Rainforth. Do Bayesian neural networks need to be fully stochastic? In *International Conference on Artificial Intelligence and Statistics*, pp. 7694–7722. PMLR, 2023.
- Guocong Song and Wei Chai. Collaborative learning for deep neural networks. In *Advances in Neural Information Processing Systems*, volume 31, 2018.
- Shengyang Sun, Guodong Zhang, Jiaxin Shi, and Roger Grosse. Functional variational Bayesian neural networks. *arXiv preprint arXiv:1903.05779*, 2019.
- Natasa Tagasovska and David Lopez-Paz. Single-model uncertainties for deep learning. In *Advances in Neural Information Processing Systems*, volume 32, 2019.

- Linh Tran, Bastiaan S Veeling, Kevin Roth, Jakub Swiatkowski, Joshua V Dillon, Jasper Snoek, Stephan Mandt, Tim Salimans, Sebastian Nowozin, and Rodolphe Jenatton. Hydra: Preserving ensemble diversity for model distillation. *arXiv preprint arXiv:2001.04694*, 2020.
- Trung Q Trinh, Markus Heinonen, Luigi Acerbi, and Samuel Kaski. Input gradient diversity for neural network ensembles. *CoRR*, 2023.
- Matias Valdenegro-Toro. Sub-ensembles for fast uncertainty estimation in neural networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 4119–4127, 2023.
- Joost van Amersfoort, Lewis Smith, Yee Whye Teh, and Yarin Gal. Uncertainty estimation using a single deep deterministic neural network. In *International Conference on Machine Learning*, pp. 9690–9700. PMLR, 2020.
- Joost van Amersfoort, Lewis Smith, Andrew Jesson, Oscar Key, and Yarin Gal. On feature collapse and deep kernel learning for single forward pass uncertainty. *arXiv preprint arXiv:2102.11409*, 2021.
- Ziyu Wang, Tongzheng Ren, Jun Zhu, and Bo Zhang. Function space particle optimization for Bayesian neural networks. *arXiv preprint arXiv:1902.09754*, 2019.
- Veit David Wild, Sahra Ghalebikesabi, Dino Sejdinovic, and Jeremias Knoblauch. A rigorous link between deep ensembles and (variational) bayesian methods. In *Advances in Neural Information Processing Systems*, 2023.
- Lisa Wimmer, Yusuf Sale, Paul Hofman, Bernd Bischl, and Eyke Hüllermeier. Quantifying aleatoric and epistemic uncertainty in machine learning: Are conditional entropy and mutual information appropriate measures? In *Uncertainty in Artificial Intelligence*, pp. 2282–2292. PMLR, 2023.
- Jim Winkens, Rudy Bunel, Abhijit Guha Roy, Robert Stanforth, Vivek Natarajan, Joseph R Ledsam, Patricia MacWilliams, Pushmeet Kohli, Alan Karthikesalingam, Simon Kohl, et al. Contrastive training for improved out-of-distribution detection. *arXiv preprint arXiv:2007.05566*, 2020.
- Guoxuan Xia and Christos-Savvas Bouganis. On the usefulness of deep ensemble diversity for out-of-distribution detection. *arXiv preprint arXiv:2207.07517*, 2022.
- Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- Shingo Yashima, Teppei Suzuki, Kohta Ishikawa, Ikuro Sato, and Rei Kawakami. Feature space particle inference for neural network ensembles. In *International Conference on Machine Learning*, pp. 25452–25468. PMLR, 2022.
- Shaofeng Zhang, Meng Liu, and Junchi Yan. The diversified ensemble neural network. In *Advances in Neural Information Processing Systems*, volume 33, pp. 16001–16011, 2020.
- Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017.
- Xiatian Zhu, Shaogang Gong, et al. Knowledge distillation by on-the-fly native ensemble. In *Advances in Neural Information Processing Systems*, volume 31, 2018.



## A Experimental details

For particle-based inference in function space (Wang et al., 2019; D’Angelo & Fortuin, 2021), we relied on the implementation available at [https://github.com/ratschlab/repulsive\\_ensembles](https://github.com/ratschlab/repulsive_ensembles), and for DDU (Mukhoti et al., 2023) at <https://github.com/omegafragger/DDU>. Table 4 summarizes relevant hyperparameters for training the base networks and our repulsive last-layer ensemble (RLL-E). For networks with spectral normalization we follow the implementation of (Mukhoti et al., 2023). Online spectral normalization with a one step power iteration is applied to convolutional weights, and exact spectral normalization is applied to 1x1 convolutional layers (Mukhoti et al., 2023). For the optimization of the ensemble head, we use a batch size of 32 for the training data and 16 for the repulsion samples for all image classification tasks.

Table 4: Implementation details and hyperparameter for the different experiments.

TASK	ARCHITECTURE	HYPERPARAMETER		VALUE
IMAGE CLASSIFICATION BASE NETWORK	RESNET-18 RESNET50	EPOCHS		50 (DirtyMNIST) 300 (CIFAR10/100)
		OPTIMIZER		SGD
		LEARNING RATE		0.1 0.01 – epoch 25 (DirtyMNIST), epoch 150 (CIFAR10/100) 0.001 – epoch 40 (DirtyMNIST), epoch 250 (CIFAR10/100)
		MOMENTUM		0.9
ACTIVE LEARNING BASE NETWORK	RESNET-18	EPOCHS		20
		OPTIMIZER		Adam
		LEARNING RATE		0.001

### Computational overhead

In all experiments, we use a pretrained base network as the feature extractor for the uncertainty estimation. For the image classification experiments, we do not modify the base network and add an ensemble head consisting of linear layers only. Thus, the number of trainable parameters is determined by the dimension of the feature space of the base network  $d$ , the number of classes  $K$ , and the number of particles  $n$ , i.e.,  $(d \times K + K) \times n$ . The feature space dimension for various base networks is shown in Table 5.

### Datasets

*DirtyMNIST*: This dataset (Mukhoti et al., 2023) consists of 60,000 clean MNIST digits with unique class labels and 60,000 synthetically generated ambiguous digits with multiple labels. All images are grayscale and of resolution  $28 \times 28$ .

*CIFAR10 / CIFAR100*: The original datasets (Krizhevsky et al., 2009) contain 50,000 training and 10,000 test images. CIFAR-10 has 10 classes; CIFAR-100 has 100 fine-grained classes. All images are  $32 \times 32$  RGB.

*CIFAR10-C / CIFAR100-C*: These benchmarks (Hendrycks & Dietterich, 2019) contain corrupted versions of the CIFAR-10 and CIFAR-100 test sets. Each consists of 10,000 test images per corruption type across 19 corruption types and 5 severity levels. Image resolution is  $32 \times 32$  RGB.

Table 5: Feature space dimension of different base network architectures.

RESNET-18	$d = 512$
WIDERESNET-28-10	$d = 640$
RESNET-50	$d = 2048$
ViT-B/16	$d = 768$

*Fine-tuning datasets:* For high-resolution evaluation, we use Stanford Cars (Krause et al., 2013) and Food101 (Bossard et al., 2014). Stanford Cars contains 16,185 images (8,144 train / 8,041 test) and 196 classes. Food101 consists of 101,000 images and 101 classes (750 train / 250 test per class). All inputs are center-cropped and resized to  $224 \times 224$ .

*OOD benchmarks:* For OOD detection on DirtyMNIST, we use FashionMNIST (Xiao et al., 2017), EMNIST (Cohen et al., 2017), and Omniglot (Lake et al., 2015) as OOD datasets. For the remaining datasets, we use Places365 (Zhou et al., 2017), SVHN (Netzer et al., 2011), ImageNet-O (Hendrycks et al., 2021), Country211 (Radford et al., 2021), and TinyImageNet (Le & Yang, 2015).

*Repulsion samples:* For the repulsion samples, we use kMNIST (Clanuwat et al., 2018), Textures (Cimpoi et al., 2014), and ImageNet-21k (Deng et al., 2009).

## B Additional Results

### B.1 Uncertainty decomposition for active learning

We evaluate the performance of our fs-RLL-E uncertainty estimates on an active learning task proposed in (Mukhoti et al., 2023). Given a small number of initial training samples and a large pool of unlabeled data, the goal is to iteratively select the most informative samples to improve model performance.

We initialize training with 20 labeled samples and use a pool composed of clean and ambiguous MNIST samples in a 1:60 ratio. In each acquisition step, the 5 samples with the highest uncertainty are selected and added to the training set. Since ambiguous samples are inherently noisy, disentangling aleatoric and epistemic uncertainty is essential for avoiding uninformative acquisitions. After each acquisition, the network is retrained from scratch on the extended dataset.

As shown in Fig. 5, acquisition based on predictive entropy, both for single models and deep ensembles, consistently performs worse. These methods tend to select ambiguous samples with high aleatoric uncertainty, which are of limited value for improving model performance. In contrast, acquisition strategies based on epistemic uncertainty – such as DEs, DDU, LL-E, and fs-RLL-E – are more effective in avoiding these uninformative samples.

However, despite this qualitative advantage, all epistemic acquisition strategies perform comparably to the random baseline. This finding is consistent with recent results by Gashi et al. (Gashi et al., 2024), who show that many active learning methods fail to outperform simple baselines across a range of experimental settings.

At the same time, our results highlight the limitations of predictive entropy as an acquisition function in the presence of significant aleatoric uncertainty. While epistemic uncertainty-based methods may not yield substantial improvements in final accuracy under the current setup, they do provide more principled guidance in avoiding harmful acquisitions.

### B.2 Full ensembling versus pretraining with post-hoc uncertainty

**Setup** We evaluate epistemic uncertainty estimates on OOD detection tasks using CIFAR-10 and CIFAR-100 as training datasets. We train a ResNet50 from scratch with spectral normalization to mitigate feature collapse. For the pretrained models, we use a ResNet50 and ViT-B/16 pretrained on ImageNet-1k, using the weight checkpoint provided by Torchvision.

For the base model trained from scratch, we reinitialize the final linear layer and retrain an ensemble of 10 linear heads, keeping the backbone frozen. We use the AdamW optimizer with a learning rate of  $10^{-4}$ . For the ImageNet-1k pretrained backbones, we follow the same protocol but fine-tune the heads and backbone jointly. We use a learning rate of  $10^{-3}$  for the linear ensemble heads and  $10^{-5}$  for the backbone. To obtain the single-model (MAP) baseline for the ImageNet-1k pretrained model, we use one linear layer and the same training procedure as for the last-layer ensembles. For function-space repulsion, we use mini batches from the Textures dataset as repulsion samples. As OOD datasets we use SVHN (Netzer et al., 2011),

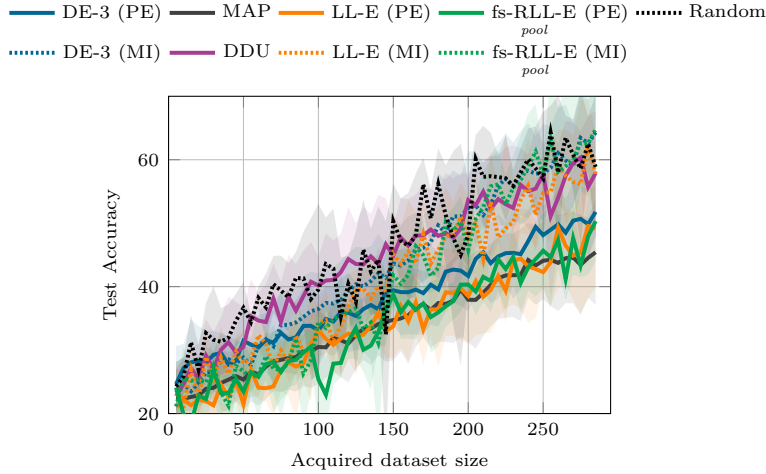


Figure 5: Test accuracy of the model as a function of the data samples that are acquired using the different uncertainty estimates. Predictive entropy (PE) combines aleatoric and epistemic uncertainty. Using the mutual information (MI) of the LL-E and fs-RLL-E prediction outperforms softmax entropy of the single network and performs on par with the other uncertainty baselines. The results are averaged over 5 runs.

Places365 (Zhou et al., 2017), ImageNet-O (Hendrycks et al., 2021), Country211 (Radford et al., 2021), and TinyImageNet (Le & Yang, 2015).

To ensure a fair comparison between models trained from scratch and ImageNet-pretrained models, we matched the ID/OOD preprocessing steps for each scenario: the low-resolution models were trained and evaluated at native CIFAR-10/100 resolution ( $32 \times 32$ ), with OOD data accordingly downsampled. Conversely, pretrained models fine-tuned on CIFAR-10/100 at  $224 \times 224$  were evaluated on OOD samples first downsampled to CIFAR resolution and then re-upscaled to avoid access to more detailed OOD data.

It has previously been shown that pretrained models can improve uncertainty estimates of a single model (Hendrycks et al., 2019a). However, as demonstrated in the DirtyMNIST experiment (see Section 6.2), a single model is not sufficient to disentangle aleatoric and epistemic uncertainty. While this may not be evident for high-performing models, it leads to performance degradation when in-distribution data exhibits increased aleatoric uncertainty due to label ambiguity or a lack of prediction accuracy.

When using a pretrained model as a fixed checkpoint, deep ensembles lose the diversity from random weight initialization. Moreover, full deep ensembles are resource-heavy (parameters/time), making post-hoc, and last-layer approaches attractive for memory-efficient uncertainty estimation. In this experiment, we aim to answer the following questions:

- (i) Is a last-layer ensemble sufficiently diverse for epistemic uncertainty estimation? Can repulsion in parameter or function space recover missing diversity?
- (ii) Given the choice, should we train a full deep ensemble from scratch, or use a pretrained model extended with post-hoc uncertainty methods?

**Results** Tables 7 and 8 summarize the ID performance and OOD detection under semantic shift, respectively. For the ID performance, we report accuracy and NLL scores. Since our main objective is to estimate epistemic uncertainty, we do not expect substantial improvements in accuracy or calibration. We use OOD detection as a proxy task for evaluating epistemic uncertainty, following the assumption that uncertainty should increase for unfamiliar inputs (de Mathelin et al., 2025).

We begin by evaluating a ResNet50 model trained from scratch with random initialization. Across all post-hoc uncertainty methods (LL-Laplace, LL-E, RLL-E, and fs-RLL-E), ID performance remains comparable to the single-model (MAP) baseline. In contrast, the deep ensemble (DE-5) improves accuracy by 1 to 3%.

Table 6: **Model configurations and parameter counts.** We report the number of trainable parameters for each base model and the increase introduced by a LL-E with 10 particles. DE-5 replicates the entire base model. LL-E shares the backbone and only replicates the final classifier layer. Reported values reflect configurations used in Tables 7 and 8 for CIFAR-10 and CIFAR-100.

Base model	Input Res.	# Parameters	
		10 classes	100 classes
ResNet50	32×32	20.74M	20.93M
× 5 (DE-5)		103.72M	104.64M
+ (fs)-RLL-E		+ 0.20M	+ 2.05M
ResNet50	224×224	23.53M	23.71M
+ (fs)-RLL-E		+ 0.20M	+ 2.05M
ViT-B/16	224×224	85.81M	85.88M
+ (fs)-RLL-E		+ 77k	+ 0.77M

This performance gain is consistent with prior work suggesting that ensembles benefit from learning more diverse features through different initializations (Allen-Zhu & Li, 2020).

Next, we consider a ResNet50 model pretrained on ImageNet-1k. Fine-tuning this model on CIFAR-10 and CIFAR-100 yields strong accuracy and calibration, matching DE-5 on CIFAR100 and outperforming it on CIFAR10. Using a more powerful backbone such as ViT-B/16 further improves ID metrics and results in the best overall performance. Table 6 reports the parameter count for each model configuration.

We now turn to OOD detection based on epistemic uncertainty. For models trained from scratch, DE-5 performs best on CIFAR10 with an AUROC of 91.15%, followed by DDU with 90.90%. In this setting, our methods do not outperform these baselines. However, on CIFAR100, function-space repulsion improves AUROC to 79.83% for LL-E to 84.05% for fs-RLL-E, outperforming the other methods.

The benefit of pretraining becomes more apparent in OOD detection. On ResNet50 and ViT-B/16, even the entropy-based uncertainty of a single model (MAP) outperforms DE-5. For example, MAP with softmax entropy on pretrained ResNet50 achieves an AUROC and AUPR of 94.19 and 93.36% on CIFAR10. Adding a LL-E improves these scores to 94.85 and 94.36%, while fs-RLL-E reaches 94.95 and 95.10%. On CIFAR100, the AUROC increases from 82.29% with MAP entropy to 84.62% with fs-RLL-E on ResNet50, and from 86.22% to 87.87% on ViT-B/16.

DDU, which was the best-performing OOD detector on DirtyMNIST, performs less effectively on this benchmark. The increased dimensionality of the feature space and the larger number of classes in CIFAR100 limit the applicability of the full GMM model. The diagonal covariance approximation and the assumption of class-conditional Gaussian features may be too restrictive. Similarly, the max logit baseline performs well on ViT-B/16 but degrades notably on ResNet-based models. LL-Laplace, while less effective on DirtyMNIST, achieves the best results in this setting.

In summary, using pretrained models not only improves ID performance but also enhances epistemic uncertainty estimation. A single pretrained model achieves similar or better results than DE-5, with significantly lower computational cost. Repulsion in function space with fs-RLL-E provides improvements over the unregularized LL-E. LL-Laplace also achieves strong results and remains a competitive baseline.

We note that CIFAR10 and CIFAR100 consist of low-resolution images, which may limit clear separability between ID and OOD samples. In Section 6.4, we considered fine-grained image classification tasks with high-resolution images. There, we demonstrated the applicability of our fs-RLL-E, especially using a more diverse set of repulsion samples.

Table 7: **In-distribution performance (Accuracy [%] / Negative Log Likelihood (NLL))**. We compare base models that are either trained from *random initialization* or *fine-tuned from ImageNet-1k pretrained weights*.

Method	CIFAR10			CIFAR100		
	ResNet50 (Random)	ResNet50 (Pretrained)	ViT-B/16 (Pretrained)	ResNet50 (Random)	ResNet50 (Pretrained)	ViT-B/16 (Pretrained)
MAP	94.93 / 0.1992	96.39 / 0.1171	98.29 / 0.0757	79.74 / 0.8087	81.86 / 0.6852	88.16 / 0.5604
LL-Laplace	94.93 / 0.1886	96.24 / 0.1180	98.27 / 0.0641	79.67 / 0.7873	80.83 / 0.7501	86.94 / 0.6191
LL-E	94.67 / 0.2005	96.49 / 0.1138	98.02 / 0.0797	79.00 / 0.8133	81.99 / 0.6722	88.17 / 0.5671
RLL-E ( <i>ours</i> )	94.66 / 0.2005	96.44 / 0.1151	98.26 / 0.0705	79.00 / 0.8133	81.98 / 0.6768	88.28 / 0.5574
fs-RLL-E ( <i>ours</i> )						
+ TEXTURES	94.77 / 0.1923	96.38 / 0.1214	98.30 / 0.0732	78.69 / 0.8227	81.90 / 0.6809	87.09 / 0.5870
DE-5	96.10 / 0.1247	— / —	— / —	82.88 / 0.6270	— / —	— / —

Table 8: **Semantic shift detection (AUROC [%] / AUPR [%])**. We compare base models that are either trained from *random initialization* or *fine-tuned from ImageNet-1k pretrained weights*. All results are averaged over 5 runs. Best results for each base model are **bold**, second-best are underlined.

Method	Uncertainty	CIFAR10			CIFAR100		
		ResNet50 (Random)	ResNet50 (Pretrained)	ViT-B/16 (Pretrained)	ResNet50 (Random)	ResNet50 (Pretrained)	ViT-B/16 (Pretrained)
MAP	MSP	88.71 / 86.12	93.45 / 91.63	97.07 / 96.91	79.54 / 76.74	80.33 / 77.83	85.09 / 83.23
	Max Logit	88.97 / 87.93	93.08 / 90.07	<u>97.92 / 98.34</u>	80.13 / 77.23	75.28 / 71.68	<u>89.32 / 89.24</u>
	Entropy	89.04 / 87.24	94.19 / 93.36	<u>97.3 / 97.49</u>	80.29 / 77.36	82.29 / 80.25	<u>86.22 / 85.7</u>
LL-Laplace	MI	89.22 / 87.01	<b>95.97 / 96.11</b>	<b>97.99 / 98.48</b>	81.27 / 78.43	<b>86.35 / 85.32</b>	<b>89.91 / 89.16</b>
DDU (diag)	GMM density	<u>90.90 / 87.54</u>	87.4 / 88.54	94.90 / 94.41	80.26 / 76.92	72.09 / 71.44	87.19 / 85.01
LL-E	MI	<u>88.91 / 86.9</u>	94.85 / 94.36	97.43 / 97.92	79.83 / 77.37	83.19 / 81.85	87.45 / 86.88
RLL-E ( <i>ours</i> )	MI	89.15 / 87.76	94.02 / 92.48	97.69 / 98.15	80.45 / 78.24	83.57 / 80.7	87.92 / 87.75
fs-RLL-E ( <i>ours</i> )							
+ TEXTURES	MI	88.36 / 86.52	<u>94.95 / 95.10</u>	97.78 / 98.33	<b>84.05 / 81.0</b>	<u>84.62 / 83.81</u>	87.87 / 88.15
DE-5	MI	<b>91.15 / 87.24</b>	— / —	— / —	79.04 / 75.28	— / —	— / —