



HAL
open science

Encoding the Latent Posterior of Bayesian Neural Networks for Uncertainty Quantification

Gianni Franchi, Andrei Bursuc, Emanuel Aldea, Séverine Dubuisson, Isabelle Bloch

► **To cite this version:**

Gianni Franchi, Andrei Bursuc, Emanuel Aldea, Séverine Dubuisson, Isabelle Bloch. Encoding the Latent Posterior of Bayesian Neural Networks for Uncertainty Quantification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024, 46 (4), pp.2027-2040. 10.1109/TPAMI.2023.3328829 . hal-04320979

HAL Id: hal-04320979

<https://hal.science/hal-04320979v1>

Submitted on 3 Apr 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Encoding the latent posterior of Bayesian Neural Networks for uncertainty quantification

Gianni Franchi, Andrei Bursuc, Emanuel Aldea, *Member, IEEE*, Séverine Dubuisson and Isabelle Bloch

Abstract—Bayesian Neural Networks (BNNs) have long been considered an ideal, yet unscalable solution for improving the robustness and the predictive uncertainty of deep neural networks. While they could capture more accurately the posterior distribution of the network parameters, most BNN approaches are either limited to small networks or rely on constraining assumptions, e.g., parameter independence. These drawbacks have enabled prominence of simple, but computationally heavy approaches such as Deep Ensembles, whose training and testing costs increase linearly with the number of networks. In this work we aim for efficient deep BNNs amenable to complex computer vision architectures, e.g., ResNet-50 DeepLabv3+, and tasks, e.g., semantic segmentation and image classification, with fewer assumptions on the parameters. We achieve this by leveraging variational autoencoders (VAEs) to learn the interaction and the latent distribution of the parameters at each network layer. Our approach, called Latent-Posterior BNN (LP-BNN), is compatible with the recent BatchEnsemble method, leading to highly efficient (in terms of computation and memory during both training and testing) ensembles. LP-BNNs attain competitive results across multiple metrics in several challenging benchmarks for image classification, semantic segmentation, and out-of-distribution detection.

Index Terms—Uncertainty estimation, Deep Neural Network ensembles, Bayesian Neural Network.

1 INTRODUCTION

MOST top-performing approaches for predictive uncertainty estimation with Deep Neural Networks (DNNs) [1], [2], [3], [4] are essentially based on ensembles, in particular Deep Ensembles (DE) [1], which have been shown to display many strengths: stability, mode diversity, good calibration, *etc.* [5]. In addition, through the Bayesian lens, ensembles enable a more straightforward separation and quantification of the sources and forms of uncertainty [1], [6], [7], which in turn allows for a better communication of the decisions to humans [8], [9] or to connected modules in an autonomous system [10]. This is crucial for real-world decision making systems. Although originally introduced as simple and scalable alternative to Bayesian Neural Networks (BNNs) [11], [12], DE still have notable drawbacks in terms of computational cost for both training and testing, that often make them prohibitive in practical applications.

In this work we address uncertainty estimation with BNNs, the departure point of DE. BNNs propose an intuitive and elegant formalism suited for this task by estimating the posterior distribution over the parameters of a network conditioned on training data. Performing exact inference BNNs is intractable and most approaches require approximations. The most common one is the mean-field assumption [13], i.e., the weights are assumed to be independent of each other and factorized by their own distribution, usually Gaussian [14], [15], [16], [17], [18], [19]. However, this approximation can be damaging [11], [20] as a more complex organization can emerge within network layers, and that higher level correlations contribute to better performance and generalization [21], [22], [23]. Yet, even under such settings, BNNs are challenging to train at scale on modern DNN architectures [24], [25]. In response, researchers have looked into structured-covariance approximations [19], [26], [27], [28],

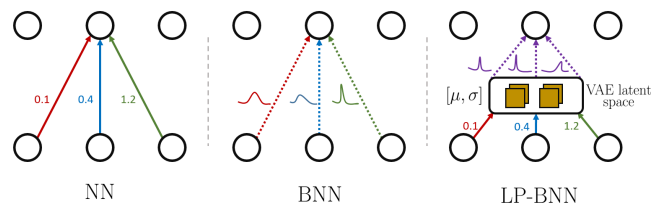


Figure 1: In a standard NN each weight has a fixed value. In most BNNs all weights follow Gaussian distributions and are assumed to be mutually independent: each weight is factorized by a Gaussian distribution. For LP-BNN in each layer, weights follow a multivariate Gaussian distribution with a latent weight space composed of independent Gaussian distributions. This enables computing expressive weight distributions in a lower dimensional space.

however they further increase memory and time complexity over the original mean-field approximation.

Here, we revisit BNNs in a pragmatic manner. We propose an approach to estimate the posterior of a BNN with layer-level inter-weight correlations, in a stable and computationally efficient manner, compatible with modern DNNs and complex computer vision tasks, e.g., semantic segmentation. We advance a novel deep BNN model, dubbed *Latent Posterior BNN* (LP-BNN), where the posterior distribution of the weights at each layer is encoded with a variational autoencoder (VAE) [29] into a lower-dimensional latent space that follows a Gaussian distribution (see Figure 1). We switch from the inference of the posterior in the high dimensional space of the network weights to a lower dimensional space which is easier to learn and already encapsulates weight interaction information. LP-BNN is naturally compatible with the recent BatchEnsemble (BE) approach [30] that enables learning a more diverse posterior from the weights of the BE sub-networks. Their combination outperforms most of related approaches across a breadth of benchmarks and metrics. In particular, LP-BNN is competitive with DE and has

Gianni Franchi is affiliated with Institut Polytechnique de Paris, Andrei Bursuc with valeo.ai, Emanuel Aldea with Paris Saclay University, Séverine Dubuisson with Aix Marseille University, and Isabelle Bloch with Sorbonne Université, CNRS, LIP6, Paris.

significantly lower costs for training and prediction.

Contributions. To summarize, the contributions of our work are: **(1)** We introduce a scalable approach for BNNs to implicitly capture *layer-level weight correlations* enabling more expressive posterior approximations, by foregoing the limiting mean-field assumption. LP-BNN scales to high capacity DNNs (e.g., 50+ layers and 30M parameters for DeepLabv3+), while still training on a single V100 GPU. **(2)** We propose to leverage VAEs for computing the posterior distribution of the weights by projecting them in the latent space. This improves significantly training stability while ensuring diversity of the sampled weights. **(3)** We extensively evaluate our method on a range of computer vision tasks and settings: image classification for *in-domain uncertainty*, *out-of-distribution (OOD) detection*, *robustness to distribution shift*, and semantic segmentation (high-resolution images, strong class imbalance) for *OOD detection*. We demonstrate that LP-BNN achieves similar performances with high-performing Deep Ensembles, while being substantially more efficient computationally.

In the rest of the paper, we will first introduce the background in Section 2, then we present our approach in Section 3, we describe all the notations used in the paper in Section 4, then we present the related works in Section 5. Finally we present our experiments and results in Section 6.

2 BACKGROUND

In this section, we present the chosen formalism for this work and a short background on BNNs. All the notations used in the paper are referenced in Table 1.

2.1 Preliminaries

We consider a training dataset $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ with n samples and labels, corresponding to two random variables $X \sim \mathcal{P}_X$ and $Y \sim \mathcal{P}_Y$. Without loss of generality we represent $\mathbf{x}_i \in \mathbb{R}^d$ as a vector, and y_i as a scalar label. We process the input data \mathbf{x}_i with a neural network $f_{\Theta}(\cdot)$ with parameters Θ , that outputs a classification or regression prediction. We view the neural network as a probabilistic model with $f_{\Theta}(\mathbf{x}_i) = P(Y = y_i | X = \mathbf{x}_i, \Theta)$. In the following, when there are no ambiguities, we discard the random variable from notations. For classification, $P(y_i | \mathbf{x}_i, \Theta)$ is a categorical distribution over the set of classes over the domain of Y , typically corresponding to the cross-entropy loss function, while for regression $P(y_i | \mathbf{x}_i, \Theta)$ is a Gaussian distribution of real values over the domain of Y when using the squared loss function. For simplicity we unroll our reasoning for the classification task.

In supervised learning, we leverage gradient descent for learning Θ that minimizes the cross-entropy loss, which is equivalent to finding the parameters that maximize the likelihood estimation (MLE) $P(\mathcal{D} | \Theta)$ over the training set $\Theta_{\text{MLE}} = \arg \max_{\Theta} \sum_{(\mathbf{x}_i, y_i) \in \mathcal{D}} \log P(y_i | \mathbf{x}_i, \Theta)$, or equivalently minimize the following loss function:

$$\mathcal{L}_{\text{MLE}}(\Theta) = - \sum_{(\mathbf{x}_i, y_i) \in \mathcal{D}} \log P(y_i | \mathbf{x}_i, \Theta). \quad (1)$$

The Bayesian approach enables adding prior information on the parameters Θ , by placing a prior distribution $\mathcal{P}(\Theta)$ upon them. This prior represents some expert knowledge about the dataset and the model. Instead of maximizing the likelihood, we can now find the maximum a posteriori (MAP) weights for $\mathcal{P}(\Theta | \mathcal{D}) \propto \mathcal{P}(\mathcal{D} |$

$\Theta)\mathcal{P}(\Theta)$ to compute $\Theta_{\text{MAP}} = \arg \max_{\Theta} \sum_{(\mathbf{x}_i, y_i) \in \mathcal{D}} \log P(y_i | \mathbf{x}_i, \Theta) + \log \mathcal{P}(\Theta)$, i.e. to minimize the following loss function:

$$\mathcal{L}_{\text{MAP}}(\Theta) = - \sum_{(\mathbf{x}_i, y_i) \in \mathcal{D}} \log P(y_i | \mathbf{x}_i, \Theta) - \log \mathcal{P}(\Theta), \quad (2)$$

inducing a specific distribution over the functions computed by the network and a regularization of the weights. For a Gaussian prior, Eq. (2) reads as L_2 regularization (weight decay).

2.2 Bayesian Neural Networks

In most neural networks only the Θ_{MAP} weights computed during training are kept for predictions. Conversely, in BNNs we aim to find the posterior distribution $P(\Theta | \mathcal{D})$ of the parameters given the training dataset, not only the values corresponding to the MAP. Here we can make a prediction y on a new sample \mathbf{x} by computing the expectation of the predictions from an infinite ensemble corresponding to different configurations of the weights sampled from the posterior distribution:

$$P(y | \mathbf{x}, \mathcal{D}) = \int P(y | \mathbf{x}, \Theta) P(\Theta | \mathcal{D}) d\Theta, \quad (3)$$

which is also known as Bayes ensemble. The integral in Eq. (3), which is calculated over the domain of Θ , is intractable, and in practice it is approximated by averaging predictions from a limited set $\{\Theta_1, \dots, \Theta_J\}$ of J weight configurations sampled from the posterior distribution:

$$P(y | \mathbf{x}, \mathcal{D}) \approx \frac{1}{J} \sum_{j=1}^J P(y | \mathbf{x}, \Theta_j). \quad (4)$$

Although BNNs are elegant and easy to formulate, their inference is non-trivial and has been subject to extensive research across the years [11], [12], [14]. Early approaches relied on Markov chain Monte Carlo variants for inference, while progress in variational inference (VI) [13] has enabled a recent revival of BNNs [15], [16], [17]. VI turns posterior inference into an optimization problem. In detail, VI finds the parameters ν of a distribution $Q_{\nu}(\Theta)$ on the weights that approximates the true Bayesian posterior distribution of the weights $P(\Theta | \mathcal{D})$ through KL-divergence minimization. This is equivalent to minimizing the following loss function, also known as expected lower bound (ELBO) loss [16], [29]:

$$\mathcal{L}_{\text{BNN}}(\Theta, \nu) = - \sum_{(\mathbf{x}_i, y_i) \in \mathcal{D}} \mathbb{E}_{\Theta \sim Q_{\nu}(\Theta)} \log (P(y_i | \mathbf{x}_i, \Theta)) + D_{\text{KL}}(Q_{\nu}(\Theta) || P(\Theta)). \quad (5)$$

The loss function $\mathcal{L}_{\text{BNN}}(\Theta, \nu)$ is composed of two terms: the KL term depends on the weights and the prior $P(\Theta)$, while the likelihood term is data dependent. This function strives to simultaneously capture faithfully the complexity and diversity of the information from data \mathcal{D} , while preserving the simplicity of the prior $P(\Theta)$. To optimize this loss function, Blundell *et al.* [16] proposed leveraging the *re-parameterization trick* [29], [31], foregoing the expensive MC estimates.

2.3 Discussion about BNNs

BNNs are particularly appealing for uncertainty quantification thanks to the ensemble of predictions from multiple weight configurations sampled from the posterior distribution. However this brings an increased computational and memory cost. For instance,

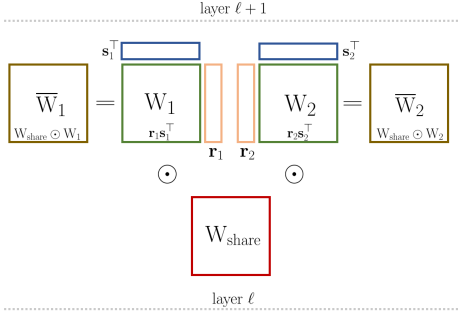


Fig. 2: **Diagram of a BatchEnsemble layer** that generates for an ensemble of size $J=2$, the ensemble weights $\overline{\mathbf{W}}_j$ from shared weights $\mathbf{W}_{\text{share}}$ and fast weights $\mathbf{W}_j=\mathbf{r}_j\mathbf{s}_j^T$, with $j \in [1, J]$.

the simplest variant of BNNs with fully factorized Gaussian approximation distributions [15], [16], *i.e.* each weight consists of a Gaussian mean and variance, carries a double amount of parameters. In addition, recent works [24], [25] point out that BNNs often underfit, and need multiple tunings to stabilize training dynamics involved by the loss function and the variance from weight samplings at each forward pass. Due to computational limitations, most BNN approaches assume that parameters are not correlated. This hinders their effectiveness [20], as empirical evidence has shown that encouraging weight collaboration improves training stability and generalization [22], [23], [32].

In order to calculate a tractable weight correlation aware posterior distribution, we propose to calculate the covariance matrix implicitly. This implicit calculation is performed by considering that each weight of the DNN has a latent variable following a distribution that we want to estimate. Hence we propose a new layer where a VAE encodes each layer to compute compressed latent distributions from which we can sample new weight configurations. We rely on the recent BatchEnsemble (BE) method [30] to further improve the parameter-efficiency of BNNs. We now proceed to describe BE and then derive our approach.

2.4 BatchEnsemble

Deep Ensembles (DEs) [1] are a popular and pragmatic alternative to BNNs. While DEs boast outstanding accuracy and predictive uncertainty, their training and testing cost increases linearly with the number of networks. This drawback has motivated the emergence of a recent stream of works proposing efficient ensemble methods [2], [3], [4], [30], [33]. One of the most promising ones is BatchEnsemble [30], which mimics in a parameter-efficient manner one of the main strengths of DE, *i.e.* diverse predictions [5].

In a nutshell, BE builds up an ensemble from a single base network (shared among ensemble members) and a set of layer-wise weight matrices specific to each member. At each layer, the weight of each ensemble member is generated from the Hadamard product between a weight shared among all ensemble members, called “*slow weight*”, and a Rank-1 matrix that varies among all members, called “*fast weight*”. Formally, let $\mathbf{W}_{\text{share}} \in \mathbb{R}^{m \times p}$ be the slow weights in a neural network layer with input dimension m and with p outputs. Each member j from an ensemble of size J owns a fast weight matrix $\mathbf{W}_j \in \mathbb{R}^{m \times p}$. \mathbf{W}_j is a Rank-1 matrix computed from a tuple of trainable vectors $\mathbf{r}_j \in \mathbb{R}^m$ and $\mathbf{s}_j \in \mathbb{R}^p$, with $\mathbf{W}_j = \mathbf{r}_j\mathbf{s}_j^T$. BE generates from them a family of ensemble weights as follows: $\overline{\mathbf{W}}_j = \mathbf{W}_{\text{share}} \odot \mathbf{W}_j$, where \odot is the Hadamard product. Each $\overline{\mathbf{W}}_j$ member of the ensemble is

essentially a Rank-1 perturbation of the shared weights $\mathbf{W}_{\text{share}}$ (see Figure 2). The sequence of operations during the forward pass reads:

$$h = a \left((\mathbf{W}_{\text{share}}^T (\mathbf{x} \odot \mathbf{s}_j)) \odot \mathbf{r}_j \right), \quad (6)$$

where a is an activation function and h the output activations.

The operations in BE can be efficiently vectorized, enabling each member to process in parallel the corresponding subset of samples from the mini-batch. $\mathbf{W}_{\text{share}}$ is trained in a standard manner over all samples in the mini-batch. A BE network $f_{\Theta^{\text{BE}}}$ is parameterized by an extended set of parameters $\Theta^{\text{BE}} = \{\theta^{\text{slow}} : \{\mathbf{W}_{\text{share}}\}, \theta^{\text{fast}} : \{\mathbf{r}_j, \mathbf{s}_j\}_{j=1}^J\}$.

With its multiple sub-networks parameterized by a reduced set of weights, BE is a practical method that can potentially improve the scalability of BNNs. We take advantage of the small size of the fast weights to capture efficiently the interactions between units and to compute a latent distribution of the weights. We detail our approach below.

3 EFFICIENT BAYESIAN NEURAL NETWORKS

3.1 Encoding the posterior weight distribution

Most BNN variants assume full independence between weights, both inter- and intra-layer. Modeling precisely weight correlations in modern high capacity DNNs with thousands to millions of parameters per layer [34] is however a daunting endeavor due to computational intractability. Yet, multiple strategies aiming to boost weight collaboration in one way or another, *e.g.* Dropout [23], WeightNorm [22], Weight Standardization [32], have proven to improve training speed, stability and generalization. Ignoring weight correlations might partially explain the shortcomings of BNNs in terms of underfitting [24], [25]. This motivates us to find a scalable way to compute the posterior distribution of the weights without discarding their correlations.

Li *et al.* [35] have recently found that the *intrinsic* dimension of DNNs can be in the order of hundreds to a few thousands. The good performances of BE, that builds on weights from a low-rank subspace, further confirm this finding. For efficiency, we leverage the Rank-1 subspace decomposition in BE and estimate here the distribution of the weights, leading to a novel form of BNNs. Formally, instead of computing the posterior distribution $P(\Theta | \mathcal{D})$, we aim now for $P(\theta^{\text{fast}} | \mathcal{D})$.

A first approach would be to compute Rank-1 weight distributions by using \mathbf{r}_j and \mathbf{s}_j as variational layers, place priors on them and compute their posterior distributions in a similar manner to [16]. Dusenberry *et al.* [25] show that these Rank-1 BNNs stabilize training by reducing the variance of the sampled weights, due to sampling only from Rank-1 variational distributions instead of full weight matrices. However this raises the memory cost significantly, as training is performed simultaneously over all J sub-networks: on CIFAR-10 for ResNet-50 with $J=4$, the authors use 8 TPUv2 cores with mini-batches of size 64 per core.

We argue that a more efficient way of computing the posterior distribution of the fast weights would be to learn instead the posterior distribution of the lower dimensional latent variables of $\{\mathbf{r}, \mathbf{s}\} \in \theta^{\text{fast}}$.¹ This can be efficiently done with a VAE [29] that can find a variational approximation $Q_\phi(\mathbf{z} | \mathbf{r})$ to the intractable posterior distribution $P_\psi(\mathbf{z} | \mathbf{r})$, with \mathbf{z} an unobserved latent random variable. We assume that $Q_\phi(\mathbf{z} | \mathbf{r})$ comes from

1. To facilitate readability we discard the j indices.

parametric families of distributions differentiable with respect to the parameter ϕ . VAEs can be seen as a generative model that can deal with complicated dependencies between input dimensions via a probabilistic encoder that projects the input into a latent space following a specific prior distribution. For simplicity and clarity, from here onward we derive our formalism only for \mathbf{r} at a single layer and consider weights \mathbf{s} to be deterministic. Here the input to the VAE are the weights \mathbf{r} and we rely on it to learn the dependencies between weights and encode them into the latent representation.

In detail, for each layer of the network $f_{\Theta}(\cdot)$ we introduce a VAE composed of a one layer encoder $g_{\phi}^{\text{enc}}(\cdot)$ with variational parameters ϕ and a one layer decoder $g_{\psi}^{\text{dec}}(\cdot)$ with parameters ψ . Let the prior over the latent variables be a centered isotropic Gaussian distribution $P_{\psi}(\mathbf{z}) = \mathcal{N}(\mathbf{z}; 0, \mathbf{I})$. Like common practice, we let the variational approximate posterior distribution $Q_{\phi}(\mathbf{z} | \mathbf{r})$ be a multivariate Gaussian with diagonal covariance. The encoder takes as input a mini-batch of size J (the size of the ensemble) composed of all the \mathbf{r}_j weights of this layer and outputs as activations $(\boldsymbol{\mu}_j, \boldsymbol{\sigma}_j^2)$. We sample a (realization) latent variable \mathbf{z}_j of a random variable distributed from $\mathcal{N}(\boldsymbol{\mu}_j, \boldsymbol{\sigma}_j^2 \mathbf{I})$ and feed it to the decoder, which in turn outputs the reconstructed weights $\hat{\mathbf{r}}_j = g_{\psi}^{\text{dec}}(\mathbf{z}_j)$. In other words, at each forward pass, we sample new fast weights $\hat{\mathbf{r}}_j$ from the latent posterior distribution to be further used for generating the ensemble. The weights of each member of the ensemble $\overline{\mathbf{W}}_j = \mathbf{W}_{\text{share}} \odot (\hat{\mathbf{r}}_j \mathbf{s}_j^{\top})$ are now random variables depending on $\mathbf{W}_{\text{share}}$, \mathbf{s}_j and \mathbf{z}_j . Note that while in practice we sample J weight configurations, this approach allows us to generate larger ensembles by sampling multiple times from the same latent distribution. We illustrate an overview of an LP-BNN layer in Figure 3.

The VAE modules are trained in the standard manner with the ELBO loss function [29] jointly with the rest of the network. The final loss function is:

$$\begin{aligned} \mathcal{L}_{\text{LP-BNN}}(\Theta^{\text{LP-BNN}}) = & - \sum_{(\mathbf{x}_i, y_i) \in \mathcal{D}} \mathbb{E}_{\mathbf{z} \sim Q_{\phi}(\mathbf{z} | \mathbf{r})} \log (P(y_i | \mathbf{x}_i, \Theta^{\text{LP-BNN}}, \mathbf{z})) \\ & + 1/L (D_{\text{KL}}(Q_{\phi}(\mathbf{z} | \mathbf{r}) || P_{\psi}(\mathbf{z})) + \|\mathbf{r} - \hat{\mathbf{r}}\|^2), \quad (7) \end{aligned}$$

where $\Theta^{\text{LP-BNN}} = \{\theta^{\text{slow}}, \theta^{\text{fast}}: \{\mathbf{r}_j, \mathbf{s}_j\}_{j=1}^J, \theta^{\text{variational}}: \{\phi, \psi\}\}$ and L the number of layers. The loss function is applied to all J members of the ensemble.

At a first glance, the loss function $\mathcal{L}_{\text{LP-BNN}}$ bears some similarities with \mathcal{L}_{BNN} (Eq. 5). Both functions include likelihood and KL terms. The likelihood in \mathcal{L}_{BNN} , *i.e.* the cross-entropy loss, depends on input data \mathbf{x}_i and on the parameters Θ sampled from $Q_{\nu}(\Theta)$, while $\mathcal{L}_{\text{LP-BNN}}$ depends on \mathbf{z}_j , a latent variable of a random variable distributed from $Q_{\phi}(\mathbf{z}_j | \mathbf{r}_j)$ that lead to the fast weights $\hat{\mathbf{r}}_j$. It guides the weights towards useful values for the main task. The KL term in \mathcal{L}_{BNN} enforces the per-weight prior, while in $\mathcal{L}_{\text{LP-BNN}}$ it preserves the consistency and simplicity of the common latent distribution of the weights \mathbf{r}_j . In addition, $\mathcal{L}_{\text{LP-BNN}}$ has an input weight reconstruction loss (last term in Eq. 7) ensuring that the generated weights $\hat{\mathbf{r}}_j$ are still compatible with the rest of the parameters of the network and do not cause high variance and instabilities during training, as typically occurs in standard BNNs [25].

At test time, we generate the LP-BNN ensemble on the fly by sampling the weights $\hat{\mathbf{r}}_j$ from the encodings of \mathbf{r}_j to compute

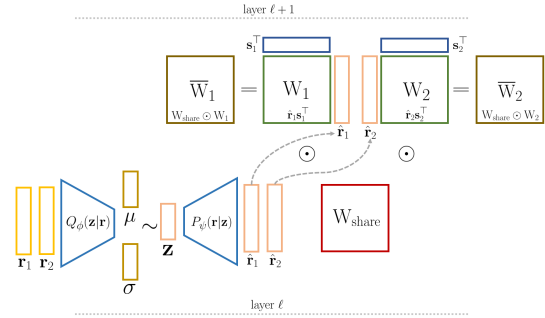


Fig. 3: **Diagram of a LP-BNN layer** that generates for an ensemble of size $J=2$, ensemble weights $\overline{\mathbf{W}}_j$ from shared weights $\mathbf{W}_{\text{share}}$ and fast weights \mathbf{s}_j and $\hat{\mathbf{r}}_j$, the latter sampled and decoded from the corresponding latent projection \mathbf{z}_j of \mathbf{r}_j , with $j \in [1, J]$.

$\overline{\mathbf{W}}_j$. For the final prediction we compute the empirical mean of the likelihoods of the ensemble:

$$P(y_i | \mathbf{x}_i) = \frac{1}{J} \sum_{j=1}^J P(y_i | \mathbf{x}_i, \theta^{\text{slow}}, \mathbf{s}_j, \hat{\mathbf{r}}_j) \quad (8)$$

3.2 Discussion on LP-BNN

We discuss here the quality of the uncertainty from LP-BNN. The predictive uncertainty of a DNN stems from two main types of uncertainty [36]: *aleatoric uncertainty* and *epistemic uncertainty*. The former is related to randomness, typically due to the noise in the data. The latter concerns finite size training datasets. The epistemic uncertainty captures the uncertainty in the DNN parameters and their lack of knowledge on the model that has generated the training data.

In BNN approaches, through likelihood marginalization over weights, the prediction is computed by integrating the outputs from different DNNs weighted by the posterior distribution (Eq. 3), allowing us to conveniently capture both types of uncertainties [7]. The quality of the uncertainty estimates depends on the diversity of predictions and views provided by the BNN. DE [1] achieve excellent diversity [5] by mimicking BNN ensembles through training of multiple individual models. Recently, Wilson and Izmailov [37] proposed to combine DE and BNNs towards improving diversity further. However, as DE are already computationally demanding, we argue that BE is a more pragmatic choice for increasing the diversity of our BNN, leading to better uncertainty quantification.

Figure 4 shows a qualitative comparison of the prediction diversity from different methods. We compare LP-BNN, BE, and DE based on WRN-28-10 [38] trained on CIFAR-10 [39] and analyze predictions on CIFAR-10, CIFAR10-C [40], and SVHN [41] test images. SVHN contains digits which have a different distribution from the training data, *i.e.*, predominant epistemic uncertainty, while CIFAR10-C displays a distribution shift via noise corruption, *i.e.*, more aleatoric uncertainty. The expected behavior is that individual DNNs in an ensemble would predict different classes for OOD images and have higher entropy on the corrupted ones, reducing the confidence score of the ensemble. We can see that the diversity of BE is lower for CIFAR10-C and SVHN, leading to poorer results in Table 8.

In the next section we discuss about our posterior covariance matrix (§ 3.3), the link between the size of the ensemble and the covariance matrix approximation (§ 3.4), the computational

complexity (§ 3.5), the training stability (§ 3.6) and stability (§ 3.7) of LP-BNN, and diversity of LP-BNN.

3.3 Covariance Priors of Bayesian Neural Network

We consider a data sample (\mathbf{x}, y) , with $\mathbf{x} \in \mathbb{R}^d$ and $y \in \mathbb{R}$. We process the input data \mathbf{x} with a MLP network $f_{\Theta}(\cdot)$ with parameters Θ composed of one hidden layer of h neurons. We detail the composing operations of the function $f_{\Theta}(\cdot)$ associated to this network: $f_{\Theta}(x) = W_2^T \sigma(W_1^T \mathbf{x})$, where $\sigma(\cdot)$ is an element-wise activation function. For simplicity, we ignore the biases in this example. $W_1 \in \mathbb{R}^{d \times h}$ is the weight matrix associated with the first fully connected layer and $W_2 \in \mathbb{R}^{h \times 1}$ the weights of the second layer. In BNNs, W_1 and W_2 represent random variables, while for classic DNNs, they are simply singular realizations of the distribution sought by BNN inference. Most works exploiting in some way the statistics of the network parameters assume, for tractability reasons, that all weights of W_1 and W_2 follow independent Gaussian distributions. Hence this leads to the following covariance matrix for W_1 :

$$\text{diag}([\text{var}(W_i[1, 1]), \dots, \text{var}(W_i[d, h])]) \quad (9)$$

where $\text{var}(W_1[i, j])$ is the variance of the coefficient $[i, j]$ of matrix W_1 , *diag* is the diagonal operator which takes a vector as input and outputs a diagonal matrix whose diagonal values are equal to the vector coefficients. Similarly, for W_2 we will have a diagonal matrix.

Now, let us assume that W_1 and W_2 have a latent representation $Z_1 = [Z_1[1], Z_1[2], Z_1[3]] \in \mathbb{R}^3$ and $Z_2 = [Z_2[1], Z_2[2], Z_2[3]] \in \mathbb{R}^3$, respectively, such that for every coefficient $[i, j]$ of W_1 and W_2 there exist real weights $\{\alpha_1^{[i,j]}[k]\}_{k=1}^3$ and $\{\alpha_2^{[i,j]}[k]\}_{k=1}^3$ such that: $W_1[i, j] = \sum_{k=1}^3 \alpha_1^{[i,j]}[k] Z_1[k]$ and $W_2 = \sum_{k=1}^3 \alpha_2^{[i,j]}[k] Z_2[k]$, respectively. In the case of LP-BNN, we consider that each coefficient of Z_1 and Z_2 represents an independent random variable. Thus, in contrast to approaches based on the mean-field approximation directly on the weights of the DNN, we can have for each layer a non-diagonal covariance matrix with the following variance and covariance terms for W_1 :

$$\text{var}(W_1[i, j]) = \sum_{k=1}^3 (\alpha_1^{[i,j]}[k])^2 \text{var}(Z_1[k]) \quad (10)$$

$$\text{cov}(W_1[i, j], W_1[i', j']) = \sum_{k=1}^3 \alpha_1^{[i,j]}[k] \alpha_1^{[i',j']}[k] \text{var}(Z_1[k]) \quad (11)$$

This allows us to leverage the lower-dimensional parameters of the distributions of Z_1 and Z_2 for estimating the higher-dimensional distributions of W_1 and W_2 . In this manner, in LP-BNN we model an **implicit covariance** of weights at each layer.

We note that several approaches for modeling correlation between weights have been proposed under certain settings and assumptions. For instance, Karaletesos and Bui [42] model correlations between weights within a layer and across layers thanks to a Gaussian process-based approach working in the function space via hierarchical priors instead of directly on the weights. Albeit elegant, this approach is still limited to relatively shallow MLPs (e.g., one hidden layer with 100 units [42]) and cannot scale up yet to deep architectures considered in this work (e.g., ResNet-50). Other approaches [26], [27] model layer-level weight correlations through

Matrix Variate Gaussian (MVG) prior distributions, increasing the expressiveness of the inferred posterior distribution at the cost of further increasing the computational complexity w.r.t. mean-field approximated BNNs [15], [16]. By contrast, LP-BNN does not **explicitly model the covariance** thanks to LP-BNN’s fully connected layers, which allow it to project the weights into a low-dimensional latent space in which it infers the posterior distribution there. This strategy leads to a lighter BNN that is competitive in terms of computation and performance for complex computer vision tasks.

3.4 The utility of Rank-1 perturbations

One could ask why using the Rank-1 perturbation formalism from BE [30], instead of simply feeding the weights of a layer to the VAE to infer the latent distribution. Rank-1 perturbations significantly reduce the number of weights upon which we train the VAE, due to the decomposition of the fast weights into \mathbf{r} and \mathbf{s} . This further allows us to consider multiple such weights at each forward pass enabling faster training of the VAE as its training samples are more numerous and more diverse.

Next, we establish connections between the cardinality J of the ensemble and the posterior covariance matrix. The mechanism of placing a prior distribution over the latent space enables an implicit modeling of correlations between weights in their original space. This is a desirable property due to its superior expressiveness [20], [42] but which can be otherwise computationally intractable or difficult to approximate. The covariance matrix of our prior in the original weight space is a Rank-1 matrix. Thanks to the Eckart-Young theorem (Theorem 5.1 in [43]), we can quantify the error of approximating the covariance by a Rank-1 matrix, based on the second up to the last singular values.

Let us denote by $\Theta_1, \dots, \Theta_J$ the J weights trained by our algorithm, $\Theta_{\text{avg}} = \frac{1}{J} \sum_{j=1}^J \Theta_j$ and $\Delta_j = \Theta_j - \Theta_{\text{avg}}$. The differences and the sum in the previous equations are calculated element-wise on all the weights of the DNNs. Then, for each new data sample \mathbf{x} , the prediction of the DNN $f_{\Theta_{\text{avg}}}(\cdot)$ is equivalent to the average of the DNNs $f_{\Theta_j}(\cdot)$ applied on \mathbf{x} :

$$\frac{1}{J} \sum_{j=1}^J f_{\Theta_j}(\mathbf{x}) = f_{\Theta_{\text{avg}}}(\mathbf{x}) + \mathcal{O}(\|\Delta\|^2) \quad (12)$$

with $\|\Delta\| = \max_j \|\Delta_j\|$. The L_2 norm is computed over all weights. We refer the reader to the proof in §3.5 of [44]. One can see that we do not learn a Rank-1 matrix, but an up to Rank- J covariance matrix, if all the $\mathbf{s}_j, \mathbf{r}_j$ are independent. Hence the choice of J acts as an approximation factor of the covariance matrix. Wen *et al.* [30] tested different values of J and found that $J = 4$ was the best compromise, which we also use here.

3.5 Computational complexity

Recent works [5], [37] studied the weight modes computed by Deep Ensembles under a BNN lens, yet these approaches are computationally prohibitive at the scale required for practical computer vision tasks. Recently, Dusenberry *et al.* [25] proposed a more scalable approach for BNNs, which can still be subject to high instabilities as the ELBO loss is applied over a high-dimensional parameter space, all BatchEnsemble parameters. Increased stability can be achieved by leveraging large mini-batches that bring more robust feature and gradient statistics, at significantly higher computational cost (large virtual mini-batches are obtained through

distributed training over multiple TPUs). In comparison, our approach has a smaller memory overhead since we encode \mathbf{r}_j in a lower dimensional space (we found empirically that a latent space of size only 32 provides an appealing compromise between accuracy and compactness). The ELBO loss here is applied over this lower-dimensional space which is easier to optimize. The only additional cost in terms of parameters and memory used w.r.t. BE is related to the compact VAEs associated with each layer.

In addition to the lower number of parameters, LP-BNN training is more stable than Rank-1 BNN [25] due to the reconstruction term $\|\mathbf{r}_j - \hat{\mathbf{r}}_j\|_2^2$ which regularizes the $\mathcal{L}_{\text{LP-BNN}}$ loss in Eq. (7) by controlling the variances of the sampled weights. In practice, to train BNNs successfully, a range of carefully crafted heuristics are necessary, e.g., clipping, initialization from truncated Normal distributions, extra weight regularization to stabilize training [25]. For LP-BNN, training is overall straightforward even on complex and deep models, e.g., DeepLabV3+, thanks to the VAE module that is stable and trains faster.

3.6 Stability of Bayesian Neural Networks

In this section, we experiment on CIFAR-10 to evaluate the stability of LP-BNN versus a classic BNN. For this experiment, we use the LeNet-5 architecture and choose a weight decay of 10^{-4} along with a mini-batch size of 128. Our goal is to see whether both techniques are stable when we vary the learning rate. Both DNNs were trained under the exact same conditions for 80 epochs. In Table 2, we present two metrics for both DNNs. The first metric is the accuracy. The second metric is the epoch during which the training loss of the DNN explodes, i.e., is equal to infinity. This phenomenon may occur if the DNN is highly unstable to train.

We argue that LP-BNN is visibly more stable than standard BNNs during training. We can see from Table 2 that LP-BNN is more stable than the standard BNN as it does not diverge for a wide range of learning rates. Moreover, its accuracy is higher than that of a standard BNN implemented on the same architecture, a property that we attribute to the VAE regularization.

3.7 LP-BNN diversity

At test-time, BNNs and DE aggregate the different predictions. For BNNs, these predictions come from the different realizations of the posterior distribution, while, for the DE, these predictions are provided by several DNNs trained in parallel. As proved in [5], [45], the diversity among these different predictions is key to quantify the uncertainty of a DNN. Indeed, we want the different DNN predictions to have a high variance when the model is not accurate. Figure 4 highlights this desirable property of high variance on out-of-distribution samples exhibited by LP-BNN. Also, as in [45], we evaluate the ratio-error introduced in [46]. The ratio-error between two classifiers is the number of data samples on which only one classifier is wrong divided by the number of samples on which they are both wrong. A higher value means that the two classifiers are less likely to make the same errors. We also evaluated the Q-statistics [46], which measures the diversity between two classifiers. The value of the Q-statistics is between -1 and 1 and is defined as:

$$Q = \frac{N_{11}N_{00} - N_{10}N_{01}}{N_{11}N_{00} + N_{10}N_{01}} \quad (13)$$

where N_{11} and N_{00} are the numbers of data on which both classifiers are correct and incorrect, respectively. N_{10} and N_{01} are the number of data where just one of the two classifiers is

wrong. If the two classifiers are always wrong or right for all data, then $N_{10}=N_{01}=0$ and $Q=1$, while if both classifiers always make errors on different inputs, then $Q=-1$. The maximum diversity comes when Q is minimum.

Finally, we evaluated the correlation coefficient [46], which assesses the correlation between the error vectors of the two classifiers. Tables 3 and 4 illustrate that, for the normal case (CIFAR-10), LP-BNN displays similar diversity with DE, while in the corrupted case (CIFAR-10-C) LP-BNN achieves better diversity scores. We conclude that in terms of diversity metrics, LP-BNN has indeed the behavior that one would expect for uncertainty quantification purposes.

4 NOTATIONS

In Table 1, we summarize the main notations used in the paper. Table 1 should facilitate the understanding of Section 2 (the preliminaries) and Section 3 (the presentation of our approach) of the paper.

5 RELATED WORK

In this section, we discuss some of the recent works on uncertainty quantification and Deep Learning.

Bayesian Deep Learning. Bayesian approaches and neural networks have a long joint history [11], [12], [47]. Early approaches relied on Markov chain Monte Carlo variants for inference on BNNs, which were later replaced by variational inference (VI) [13] in the context of deeper networks. Hernández-Lobato *et al.* [17] propose a new backpropagation technique, called probabilistic backpropagation, which allows them to estimate the posterior of the BNN without variational inference. Gal *et al.* [18] use the dropout during inference time to approximate variational BNN. Most of the modern approaches make use of VI with the mean-field approximation [14], [15], [16], [19] which conveniently makes posterior inference tractable. However this limits the expressivity of the posterior [11], [20]. This drawback became subject of multiple attempts for structured-covariance approximations using matrix variate Gaussian priors [26], [27] or natural gradients [19], [28]. However they further increase memory and time complexity over the original mean-field approximation. Recent methods proposed more simplistic BNNs by performing inference with structured priors only over the first and last layer [48] or just the last layer [24], [49]. Full covariance can be computed for shallow networks thanks to a meta-prior in a low-dimensional space where the VI can be performed [42], [50]. Rossi *et al.* [51] use the Walsh-Hadamard-based factorization strategies to project the covariance on a space where it is a diagonal matrix, leading to a BNN that tracks the covariance in an easier way. Tran *et al.* [52] use a functional prior on a BNN to better estimate the posterior. Most variational inference based BNNs are still challenging to train, underfit and are difficult to scale to big DNNs [24], [25], while the issue of finding a proper prior is still open [53], [54]. Our approach builds upon the low dimensional fast weights from BE and on the stability of the VAEs, foregoing many of the shortcomings of BNN training.

Ensemble Approaches. Ensembles mimic and attain, to some extent, properties of BNNs [5]. Deep Ensembles [1] train multiple DNNs with different random initializations leading to excellent uncertainty quantification scores. The major inherent drawback in terms of computational and memory overhead has been subsequently addressed through multi-head networks [55],

| Notations | Meaning |
|---|---|
| $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ | the set of n data samples and the corresponding labels |
| Θ | the set of weights of a DNN |
| $P(\Theta)$ | the prior distribution over the weights of a DNN |
| $Q_\nu(\Theta)$ | the variational prior distribution over the weights of a DNN used in standard BNNs [16] |
| ν | the parameters of the variational prior distribution over the weights of a DNN used in standard BNNs [16] |
| $P(y_i \mathbf{x}_i, \Theta)$ | the likelihood that DNN outputs y_i following a prediction over input image \mathbf{x}_i |
| J | the number of ensembling DNNs |
| $\theta^{\text{slow}} = \{W_{\text{share}}\}$ | the shared “slow” weights of the network |
| $\theta^{\text{fast}} = \{W_j\}_{j=1}^J = \{(\mathbf{r}_j, \mathbf{s}_j)\}_{j=1}^J$ | the set of individual “fast” weights of BatchEnsemble for ensembling of J networks |
| $\hat{\theta}^{\text{fast}} = \{(\hat{\mathbf{r}}_j, \hat{\mathbf{s}}_j)\}_{j=1}^J$ | the set of fast weights of LP-BNN for ensembling of J networks. $\hat{\mathbf{r}}_j$ are sampled from the latent weight space of weights \mathbf{r}_j . |
| $\theta^{\text{variational}} = \{(\phi_j, \psi_j)\}_{j=1}^J$ | the parameters of the VAE for computing the low dimensional latent distribution of \mathbf{r}_j |
| $g_\phi^{\text{enc}}(\cdot)$ | the encoder of the VAE applied on \mathbf{r} |
| $g_\psi^{\text{dec}}(\cdot)$ | the decoder of the VAE for reconstructing $\hat{\mathbf{r}}$ from latent code of \mathbf{r} |
| $Q_\phi(\mathbf{z} \mathbf{r})$ | the variational distribution over the weights \mathbf{r} to approximate the intractable posterior $P_\psi(\mathbf{z} \mathbf{r})$ |
| $(\boldsymbol{\mu}_j, \boldsymbol{\sigma}_j) = g_\phi^{\text{enc}}(\mathbf{r}_j)$ | encoder output that parameterize a multivariate Gaussian with diagonal covariance |
| $\mathbf{z}_j \sim Q_\phi(\mathbf{z} \mathbf{r}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}_j, \boldsymbol{\sigma}_j^2 \mathbf{I})$ | sampling a latent code \mathbf{z} from the latent distribution |
| $P_\psi(\mathbf{z}_j) = \mathcal{N}(\mathbf{z}_j; 0, \mathbf{I})$ with $j \in [1, J]$ | the prior distribution on \mathbf{z}_j |
| $\hat{\mathbf{r}}_j = g_\psi^{\text{dec}}(\mathbf{z}_j)$ | the reconstruction of \mathbf{r}_j from its latent distribution, <i>i.e.</i> the variational fast weights |
| $\bar{W}_j = W_{\text{share}} \odot (\mathbf{r}_j \mathbf{s}_j^\top)$ | the weight of a BatchEnsemble network j computed from slow and fast weights where \odot is the Hadamard product and $(\mathbf{r}_j \mathbf{s}_j^\top)$ the inner product between these two vectors. |
| $\bar{W}_j = W_{\text{share}} \odot (\hat{\mathbf{r}}_j \mathbf{s}_j^\top)$ | the weight of LP-BNN network network j computed from slow and variational fast weights |

TABLE 1: Summary of the main notations of the paper.

| Learning Rate | 0.2 | 0.1 | 0.05 | 0.01 | 0.005 |
|------------------|-------|-------|-------|-------|-------|
| BNN accuracy | 22.48 | 44.60 | 49.83 | 48.70 | 56.69 |
| BNN epoch div | 3 | 25 | 65 | None | None |
| LP-BNN accuracy | 20.02 | 55.04 | 59.68 | 63.73 | 64.41 |
| LP-BNN epoch div | 3 | None | None | None | None |

TABLE 2: Stability analysis of BNNs. Stability experiment with LeNet 5 architecture and 80 epochs on CIFAR-10. On the epoch divergence row, *None* means that the DNN does not diverge.

| | ratio errors \uparrow | Q-statistic \downarrow | correlation coefficient \downarrow |
|--------|-------------------------|--------------------------|--------------------------------------|
| DE | 0.9825 | 0.9877 | 0.6583 |
| BE | 0.5915 | 0.9946 | 0.7634 |
| LP-BNN | 0.8390 | 0.9842 | 0.6601 |

TABLE 3: Comparative results of diversity scores for image classification on the CIFAR-10 dataset.

snapshot-ensembles from intermediate training checkpoints [56], [57], efficient computation of the posterior distribution from weight trajectories during training [3], [4], use of multiple Dropout masks at test time [18], multiple random perturbations to the weights of a pre-trained network [4], [33], [58], multiple perturbation of the input image [2], multiple low-rank weights tied to a backbone network [30], simultaneous processing of multiple images by the same DNN [59]. Most approaches still have a significant computational overhead for training or for prediction, while struggling with diversity [5].

Dirichlet Networks (DNs). DNs [7], [60], [61], [62], [63] bring a promising line of approaches that estimate uncertainty from a single network by parameterizing a Dirichlet distribution over its predictions. However, most of these methods [7], [60] use OOD samples during training, which may be unrealistic in many

| | ratio errors \uparrow | Q-statistic \downarrow | correlation coefficient \downarrow |
|--------|-------------------------|--------------------------|--------------------------------------|
| DE | 0.4193 | 0.9690 | 0.7568 |
| BE | 0.2722 | 0.9874 | 0.8352 |
| LP-BNN | 0.4476 | 0.9595 | 0.7332 |

TABLE 4: Comparative results of diversity scores for image classification on the CIFAR-10-C dataset.

applications [62], or do not scale to bigger DNNs [64]. DNs have been developed only for classification tasks and extending them to regression requires further adjustments [65], unlike LP-BNN that can be equally used for classification and regression.

Stochastic MCMC BNN. Among the different solutions for sampling according to the posterior of the DNN we also consider the methods based on Markov chain Monte Carlo (MCMC), which lead to quantification of uncertainty. For example, one could use traditional sampling strategies [66], however, these techniques suffer from the curse of dimensionality. One technique that is often used is based on stochastic Langevin gradient dynamics (SGLD) [67] which relies on adding noise to the stochastic gradient descent and an adapted Metropolis hasting algorithm. This method does not take into account the moment of the Langevin gradient which is a crucial term of the Hamiltonian Monte Carlo (HMC). This term was added in stochastic gradient Hamiltonian Monte Carlo (SGHMC) [68], [69]. Unfortunately, as pointed out by Dauphin et al. (2014), DNNs often exhibit pathological curvature and their posterior has numerous modes. Thus the authors in [70] proposed a strategy to estimate the multi modes of the posterior based on the flat histogram algorithm [71]. Kim et al. [72] proposed to use an adaptive drift in the gradient to exit the saddle points of the loss function. Marceau-Caron et al. [73] used the natural Langevin dynamics combined with the natural gradient descent to have a better estimate of the posterior samples.

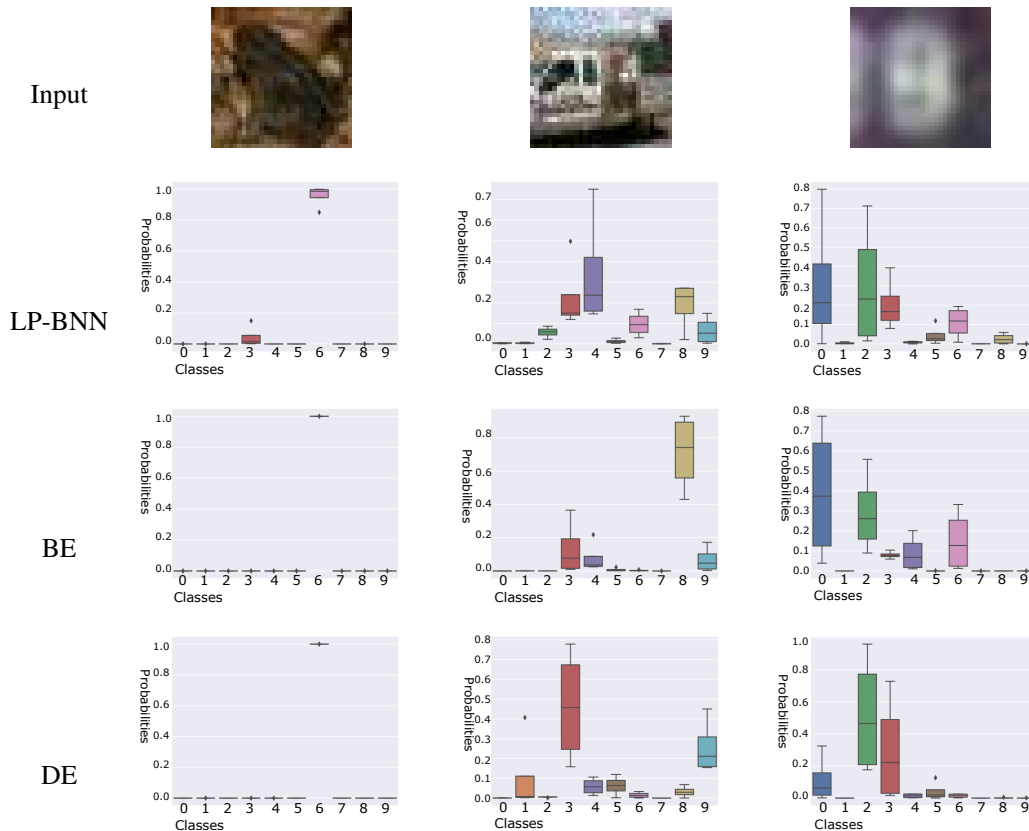


Fig. 4: **Diversity of predictions of different ensemble methods.** The first row contains in order two images from the test set of CIFAR-10, of CIFAR-10-C and of SVHN, respectively. The next three rows represent the corresponding outputs of the different sub models for the three ensembling algorithms being considered: LP-BNN, BatchEnsemble and Deep Ensembles.

6 EXPERIMENTS AND RESULTS

6.1 Implementation details

We evaluate the performance of LP-BNN in assessing the uncertainty of its predictions. For our benchmark, we evaluate LP-BNN on different scenarios against several strong baselines with different advantages in terms of performance, training or runtime: BE [30], DE [1], Maximum Class Probability (MCP) [74], MC Dropout [18], TRADI [4], EDL [61], DUQ [75], MIMO [59], Subspace Inference (SI) [76], Stochastic gradient Hamiltonian Monte Carlo (SGHMC) [68], BBB [16], FBNN [77], and GPI-G prior [52].

First, we evaluate the predictive performance in terms of accuracy for image classification, MSE for regression, and mIoU [78] for semantic segmentation, respectively. Secondly, we evaluate the quality of the confidence scores provided by the DNNs by means of Expected Calibration Error (ECE) [79]. For ECE we use M -bin histograms of confidence scores and accuracy, and compute the average of M bin-to-bin differences between the two histograms. Similarly to [79] we set $M = 15$. To evaluate the robustness to dataset shift via corrupted images, we first train the DNNs on CIFAR-10 [39] or CIFAR-100 [39] and then test on the corrupted versions of these datasets [40]. The corruptions include different types of noise, blurring, and some other transformations that alter the quality of the images. The following corruptions are applied with different levels of severity: Gaussian Noise, Shot Noise, Impulse Noise, Defocus Blur, Frosted Glass Blur, Motion Blur, Zoom Blur, Snow, Frost, Fog, Brightness, Contrast, Elastic, Pixelate, JPEG. More details about these corruptions can be found

in [40].

For this scenario, similarly to [80], we use as evaluation measures the Corrupted Accuracy (cA) and Corrupted Expected Calibration Error (cE), that offer a better understanding of the behavior of our DNN when facing shift of data distribution and aleatoric uncertainty.

In order to evaluate the epistemic uncertainty, we propose to assess the OOD detection performance. This scenario typically consists in training a DNN over a dataset following a given distribution, and testing it on data coming from this distribution and data from another distribution, not seen during training. We quantify the confidence of the DNN predictions in this setting through their prediction scores, i.e., output softmax values. We use the same indicators of the accuracy of detecting OOD data as in [74]: AUC, AUPR, and the FPR-95%-TPR. These indicators measure whether the DNN model lacks knowledge regarding some specific data and how reliable are its predictions. For the regression, we use the Negative Log Likelihood (NLL) as a performance metric for the uncertainty.

For our implementation, we use PyTorch [81]. Our code is available at https://github.com/giannifranchi/LP_BNN. In the following we share the hyper-parameters for our experiments on image classification and semantic segmentation. Please be aware that the values in bold within the tables are linked to the top result associated with the given metrics.

6.2 Regression

To verify whether our work also generalizes to trivial tasks, we propose to study regression tasks on the UCI dataset [17]. We

| Method | CIFAR-10 | | | | CIFAR-100 | | | | | | |
|--------------------------------|--------------|---------------|---------------|--------------|---------------|--------------|---------------|--------------|---------------|--------------|---------------|
| | Acc ↑ | AUC ↑ | AUPR ↑ | FPR-95-TPR ↓ | ECE ↓ | cA ↑ | cE ↓ | Acc ↑ | ECE ↓ | cA ↑ | cE ↓ |
| MCP [74] | 96.33 | 0.9600 | 0.9767 | 0.115 | 0.0207 | 32.98 | 0.6167 | 80.19 | 0.1228 | 19.33 | 0.7844 |
| MC dropout [18] | 96.50 | 0.9273 | 0.9603 | 0.242 | 0.0117 | 32.35 | 0.6403 | 77.92 | 0.0672 | 27.66 | 0.5909 |
| SI [76] | 96.31 | 0.9681 | 0.9842 | 0.101 | 0.0105 | 73.30 | 0.1079 | 81.89 | 0.0736 | 51.78 | 0.2283 |
| DUQ [75] [†] | 92.9 | 0.9338 | 0.9600 | 0.150 | 0.1572 | 69.10 | 0.2712 | - | - | - | - |
| DUQ Resnet18 [75] [‡] | 92.8 | 0.9138 | 0.9421 | 0.20 | 0.064 | 67.21 | 0.5912 | - | - | - | - |
| EDL [61] [†] | 85.73 | 0.9002 | 0.9198 | 0.247 | 0.0904 | 59.54 | 0.3412 | - | - | - | - |
| SGHMC [68] | 76.59 | 0.6023 | 0.7820 | 0.8367 | 0.1439 | 56.40 | 0.1640 | - | - | - | - |
| BBB [16] | 75.11 | 0.6820 | 0.8079 | 0.7075 | 0.071 | 49.12 | 0.3647 | - | - | - | - |
| MIMO [59] | 95.73 | 0.7740 | 0.8990 | 0.175 | 0.0261 | 67.95 | 0.2530 | 79.12 | 0.0718 | 47.12 | 0.2901 |
| Deep Ensembles [1] | 96.74 | 0.9803 | 0.9896 | 0.071 | 0.0093 | 68.75 | 0.1414 | 83.01 | 0.0673 | 47.35 | 0.2023 |
| BatchEnsembles [30] | 96.48 | 0.9540 | 0.9731 | 0.132 | 0.0167 | 71.67 | 0.1928 | 81.27 | 0.0912 | 47.44 | 0.2909 |
| LP-BNN (ours) | 95.02 | 0.9691 | 0.9836 | 0.103 | 0.0094 | 69.51 | 0.1197 | 79.3 | 0.0702 | 48.40 | 0.2224 |

TABLE 5: **Comparative results for image classification tasks.** We evaluate on CIFAR-10 and CIFAR-100 for the tasks: in-domain classification, out-of-distribution detection with SVHN (CIFAR-10 only), robustness to distribution shift (CIFAR-10-C, CIFAR-100-C). We run all methods ourselves in similar settings using publicly available code for related methods. Results are averaged over three seeds. [†]: We did not manage to scale these methods to WRN-28-10 on CIFAR-100. A similar finding for EDL was reported in [64]. [‡] DUQ does not scale on CIFAR-100 and it does not perfectly scale to WRN-28-10 on CIFAR-10 so we train it with ResNet-18 [34] architecture like in the original paper.

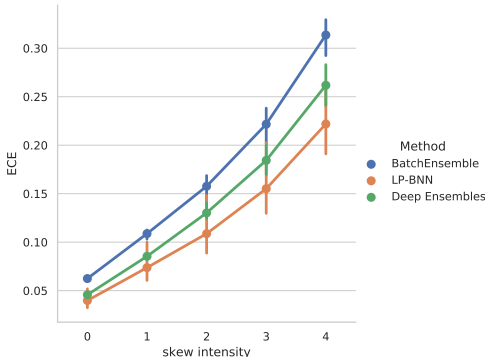


Fig. 5: **Calibration at different levels of corruption.** We report ECE scores for LP-BNN, BE [30], and DE [1] on CIFAR-10-C.

reproduce the framework developed in [17] and [1]. Similarly to [17], we consider a MLP neural network with one hidden layer, composed of 50 hidden units trained for 40 epochs with 20-fold cross-validation.

Based on [82], we use the Gaussian NLL in Equation (14) using networks with two output neurons which estimate the parameters of a heteroscedastic Gaussian distribution [82], [83]. Hence we consider that the outputs of the DNN are the parameters of a Gaussian distribution (mean and variance).

$$\mathcal{L}(\mu_{\theta_m}(\mathbf{x}_i), \sigma_{\theta_m}(\mathbf{x}_i)^2, y_i) = \frac{(y_i - \mu_{\theta_m}(\mathbf{x}_i))^2}{2\sigma_{\theta_m}(\mathbf{x}_i)^2} + \frac{1}{2} \log \sigma_{\theta_m}(\mathbf{x}_i)^2 + \frac{1}{2} \log 2\pi \quad (14)$$

We compare LP-BNN, BBB [16], FBNN [77], GPI-G prior [52], Deep Ensembles [1] on the UCI datasets in Table 6. The results show that LP-BNN provides equivalent results to state-of-the-art algorithms on most datasets.

6.3 Image classification with CIFAR-10/100 [39]

Protocol. Here we train on CIFAR-10 [39] composed of 10 classes. For CIFAR-10 we consider as OOD the SVHN dataset [41]. Since SVHN is a color image dataset of digits, it guarantees that the OOD data comes from a distribution different from those of CIFAR-10. We use WRN-28-10 [38] for all methods, a popular architecture for this dataset, and evaluate on CIFAR-10-C [40]. For CIFAR-100 [39] we use again WRN-28-10 and evaluate on the test sets of CIFAR-100 and CIFAR-100-C [40]. Note that for all DNNs, even

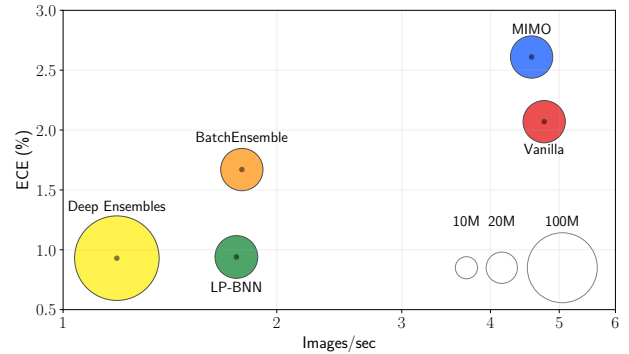


Fig. 6: **ECE vs. network size and inference time.** We report ECE scores for LP-BNN, BE [30], MIMO [59], and DE [1] on CIFAR-10. The size of each circle corresponds to the number of parameters of the model.

for DE, we average results over three random seeds for statistical relevance. We use cutout [84] as data augmentation, as commonly used for these datasets. Please find in the supplementary material the hyperparameters for this experiment.

Discussion. We illustrate results for this experiment in Table 8. We notice that DE with cutout outperforms other methods on most of the metrics except ECE, cA, and cE on CIFAR-10, and cA on CIFAR-100, where LP-BNN achieves state of the art results. This means that LP-BNN is competitive for aleatoric uncertainty estimation. In fact, ECE is calculated on the test set of CIFAR-10 and CIFAR-100, so it mostly measures the reliability of the confidence score in the training distribution. cA and cE are evaluated on corrupted versions of CIFAR-10 and CIFAR-100, which amounts to quantifying the aleatoric uncertainty. We can see that for this kind of uncertainty, LP-BNN achieves state of the art performance. On the other hand, for epistemic uncertainty, we can see that DE always attain best results. Overall, our LP-BNN is more computationally efficient while providing better results for the aleatoric uncertainty. Computation wise, DE takes 52 hours to train on CIFAR-10, while LP-BNN needs 2 times less, 26 hours and 30 minutes. In Figure 5 and Table 5, we observe that our method exhibits top ECE score on CIFAR-10-C, as well as for the stronger corruptions.

6.4 Semantic segmentation

Next, we evaluate semantic segmentation, a task of interest for autonomous driving, where high capacity DNNs are used for

| Datasets | RMSE | | | | | | NLL | | | | | |
|------------------------|----------------|---------------|---------------|----------------------|----------------|----------------------|----------------|-------------|--------------|----------------------|----------------|-----------------------|
| | LP-BNN | BBB | FBNN | GPI-G prior | SGHMC | Deep Ensembles | LP-BNN | BBB | FBNN | GPI-G prior | SGHMC | Deep Ensembles |
| Boston housing | 2.995 ± 0.545 | 3.171±0.149 | 2.378±0.104 | 2.850 ± 1.007 | 4.496 ± 2.272 | 2.219 ± 0.098 | 2.611 ± 0.175 | 2.602±0.031 | 2.301±0.038 | 2.469 ± 0.160 | 3.097 ± 0.465 | 2.047 ± 0.028 |
| Concrete | 5.900 ± 0.653 | 5.678±0.087 | 4.935±0.180 | 4.781 ± 0.443 | 7.995 ± 2.414 | 5.167 ± 0.234 | 3.250 ± 0.133 | 3.149±0.018 | 3.096±0.016 | 3.007 ± 0.057 | 3.461 ± 0.226 | 2.885 ± 0.032 |
| Energy | 2.537 ± 0.760 | 0.565±0.018 | 0.412±0.017 | 0.370 ± 0.076 | 8.222 ± 23.197 | 1.712 ± 0.067 | 2.553 ± 0.227 | 1.500±0.006 | 0.684±0.020 | 0.425 ± 0.210 | 2.738 ± 1.077 | 1.553 ± 0.060 |
| Kin8nm | 0.07 ± 0.004 | - | - | 0.065 ± 0.002 | 0.083 ± 0.019 | 0.058 ± 0.003 | -1.241 ± 0.010 | - | - | -1.241 ± 0.015 | -1.156 ± 0.149 | -1.452 ± 0.010 |
| Naval Propulsion Plant | 0.002 ± 0.000 | 0.002 ± 0.000 | 0.001 ± 0.000 | 0.009 ± 0.000 | 0.001 ± 0.001 | 0.002 ± 0.000 | -4.938 ± 0.041 | - | -6.950±0.052 | -7.130±0.024 | -6.923 ± 0.062 | -3.337 ± 0.185 |
| Power Plant | 3.303 ± 0.02 | - | - | 3.936 ± 0.170 | - | 3.097 ± 0.020 | 2.660 ± 0.007 | - | - | 2.790 ± 0.043 | - | 2.600 ± 0.007 |
| Protein | 3.624 ± 0.0718 | 4.331±0.033 | 4.326±0.019 | 3.926 ± 0.019 | 3.822 ± 0.069 | 3.412 ± 0.017 | 2.533 ± 0.043 | 2.892±0.007 | 2.892±0.004 | 2.799 ± 0.004 | 2.650 ± 0.049 | 2.442 ± 0.015 |

TABLE 6: Comparison between the results obtained with LP-BNN on regression tasks

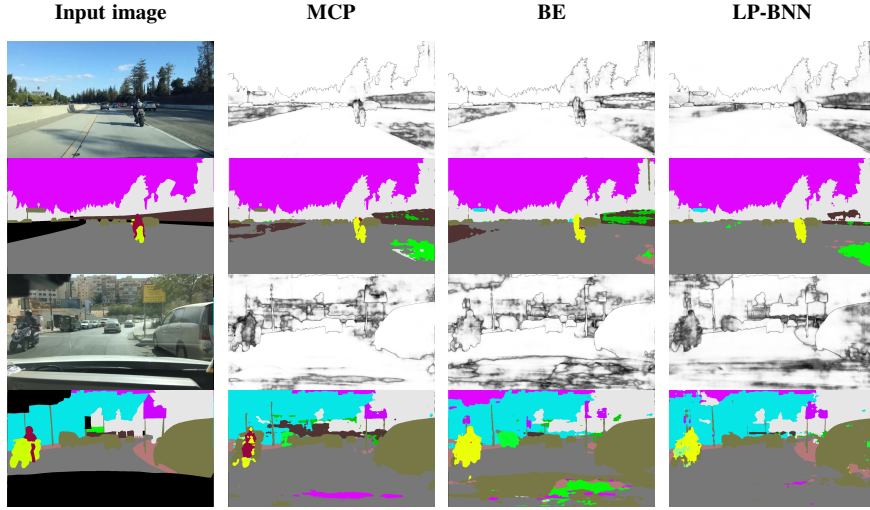


Figure 7: Visual assessment on two BDD-Anomaly test images containing a motorcycle (OOD class). For each image: on the first row, input image and confidence maps from MCP [74], BE [30], and LP-BNN; on the second row, ground-truth segmentation and segmentation maps from MCP, BE, and LP-BNN. LP-BNN is less confident on the OOD objects.

| | | Vanilla | BatchEnsemble | Deep Ensembles | TRADI | LP-BNN |
|-----------|--------------|---------|---------------|----------------|-------|--------|
| Training | Time (s) | 1,506 | 1,983 | 6,026 | 1,782 | 1,999 |
| | for 1 epochs | 1 | ×1.31 | ×4.0 | ×1.18 | ×1.33 |
| | Memory (MiB) | 8,848 | 9,884 | 35,392 | 9,040 | 9,888 |
| Testing | Time (s) | 0.21 | 0.56 | 0.84 | 0.84 | 0.57 |
| | on 1 image | 1 | ×2.67 | ×4.0 | ×4.0 | ×2.71 |
| | Memory (MiB) | 1,884 | 4,114 | 7,536 | 7,536 | 4,114 |
| Mult-Adds | | 5.95 | 23.82 | 23.82 | 21.24 | 23.81 |
| | | 1 | ×4.0 | ×4.0 | ×3.56 | ×4.0 |

TABLE 7: Runtime and memory analysis. Numbers correspond to StreetHazards images processed with DeepLabv3+ ResNet-50 with PyTorch on a PC: Intel Core i9-9820X and 1× GeForce RTX 2080Ti. Colored numbers are relative to vanilla approach. Mini-batch size for training is 4 and for testing 1. *Mult-Adds* corresponds to the inference cost, i.e., the number of giga multiply-add operations for a forward pass which is estimated with `torchinfo` (v1.7.1).

| Dataset | OOD method | mIoU ↑ | AUC ↑ | AUPR ↑ | FPR-95-TPR ↓ | ECE ↓ |
|--------------------|---------------------|---------------|---------------|---------------|---------------|---------------|
| StreetHazards | Baseline (MCP) [74] | 53.90 | 0.8660 | 0.0691 | 0.3574 | 0.0652 |
| | TRADI [4] | 52.46 | 0.8739 | 0.0693 | 0.3826 | 0.0633 |
| | Deep Ensembles [1] | 55.59 | 0.8794 | 0.0832 | 0.3029 | 0.0533 |
| | MIMO [59] | 55.44 | 0.8738 | 0.0690 | 0.3266 | 0.0557 |
| | BatchEnsemble [30] | 56.16 | 0.8817 | 0.0759 | 0.3285 | 0.0609 |
| | LP-BNN (ours) | 54.50 | 0.8833 | 0.0718 | 0.3261 | 0.0520 |
| LP-BNN + GN (ours) | 56.12 | 0.8908 | 0.0742 | 0.2999 | 0.0593 | |
| BDD-Anomaly | Baseline (MCP) [74] | 47.63 | 0.8515 | 0.0450 | 0.2878 | 0.1768 |
| | TRADI [4] | 44.26 | 0.8480 | 0.0454 | 0.3687 | 0.1661 |
| | Deep Ensembles [1] | 51.07 | 0.8480 | 0.0524 | 0.2855 | 0.1419 |
| | MIMO [59] | 47.20 | 0.8438 | 0.0432 | 0.3524 | 0.1633 |
| | BatchEnsemble [30] | 48.09 | 0.8427 | 0.0449 | 0.3017 | 0.1690 |
| | LP-BNN (ours) | 49.01 | 0.8532 | 0.0452 | 0.2947 | 0.1716 |
| LP-BNN + GN (ours) | 47.15 | 0.8553 | 0.0577 | 0.2866 | 0.1623 | |

TABLE 8: Comparative results on the OOD task for semantic segmentation. We run all methods ourselves in similar settings using publicly available code for related methods. Results are averaged over three seeds.

processing high resolution images with complex urban scenery with strong class imbalance.

StreetHazards [85]. StreetHazards is a large-scale dataset that consists of different sets of synthetic images of street scenes.

More precisely, this dataset is composed of 5,125 images for training and 1,500 test images. The training dataset contains pixel-wise annotations for 13 classes. The test dataset comprises 13 training classes and 250 OOD classes, unseen in the training set, making it possible to test the robustness of the algorithm when facing a diversity of possible scenarios. For this experiment, we used DeepLabv3+ [86] with a ResNet-50 encoder [34]. Following the implementation in [85], most papers use PSPNet [87] that aggregates predictions over multiple scales, an ensembling that can obfuscate in the evaluation the uncertainty contribution of a method. This can partially explain the excellent performance of MCP on the original settings [85]. We propose using DeepLabv3+ instead, as it enables a clearer evaluation of the predictive uncertainty. We propose two DeepLabv3+ variants as follows. DeepLabv3+ is composed of an encoder network and a decoder network; in the first version, we change the decoder by replacing all the convolutions with our new version of LP-BNN convolutions and leave the encoder unchanged. In the second variant we use weight standardization [88] on the convolutional layers of the decoder, replacing batch normalization [89] in the decoder with group normalization [90], to better balance mini-batch size and ensemble size. We denote the first version LP-BNN and the second one LP-BNN + GN.

BDD-Anomaly [85]. BDD-Anomaly is a subset of the BDD100K dataset [91], composed of 6,688 street scenes for training and 361 for the test set. The training set contains pixel-level annotations for 17 classes, and the test dataset is composed of the 17 training classes and 2 OOD classes: motor-cycle and train. For this experiment, we use DeepLabv3+ [86] with the experimental protocol from [85]. As previously we use a ResNet-50 encoder [34]. For this experiment, we use the LP-BNN and LP-BNN + GN variants.

Discussion. We emphasize that the semantic segmentation is more

challenging than the CIFAR classification since images are bigger and their content is more complex. The larger input size constrains to use smaller mini-batches. This is crucial since the fast weights of the ensemble layers are trained just on one mini-batch slice. In this experiment, we could use mini-batches of size 4 and train the fast weights on slices of size 1. Yet, despite these computational difficulties, with our technique, we achieve state-of-the-art results for most metrics. We can see in Table 8 that our strategies achieve state-of-the-art performance in detecting OOD data and are well calibrated. We can also see in Figure 7, where the OOD class is the motorcycle, that our DNN is less confident on this class. Hence LP-BNN allows us to have a more reliable DNN which is essential for real-world applications.

Table 7 and Figure 6 show the computational cost of LP-BNN and related methods. For training, LP-BNN takes only $\times 1.33$ more time than a vanilla approach, in contrast to DE that take much longer, while their performances are equivalent in most cases. Our technique allows for a lighter training than DE, which is interesting when using GPUs with low VRAM. At the same time, LP-BNN enables implicit modeling of weight correlations at every layer with limited overhead as it does not explicitly compute the covariances. To the best of our knowledge, LP-BNN is the first approach with the posterior distribution computed with variational inference successfully trained and applied for semantic segmentation.

7 CONCLUSION

We propose a new BNN framework able to quantify uncertainty in the context of deep learning. Owing to each layer of the network being tied to and regularized by a VAE, LP-BNNs are stable, efficient, and therefore easy to train compared to existing BNN models. The extensive empirical comparisons on multiple tasks show that LP-BNNs reach state-of-the-art levels with substantially lower computational cost. We hope that our work will open new research paths on effective training of BNNs. In the future we intend to explore new strategies for plugging more sophisticated VAEs in our models along with more in-depth theoretical studies.

ACKNOWLEDGMENTS

This work was performed using HPC resources from GENCI-IDRIS (Grant 2020-AD011011970) and (Grant 2021-AD011011970R1).

REFERENCES

- [1] B. Lakshminarayanan, A. Pritzel, and C. Blundell, "Simple and scalable predictive uncertainty estimation using deep ensembles," in *NeurIPS*, 2017. 1, 3, 4, 6, 8, 9, 10
- [2] A. Ashukha, A. Lyzhov, D. Molchanov, and D. Vetrov, "Pitfalls of in-domain uncertainty estimation and ensembling in deep learning," in *ICLR*, 2020. 1, 3, 7
- [3] W. J. Maddox, P. Izmailov, T. Garipov, D. P. Vetrov, and A. G. Wilson, "A simple baseline for bayesian uncertainty in deep learning," in *NeurIPS*, 2019. 1, 3, 7
- [4] G. Franchi, A. Bursuc, E. Aldea, S. Dubuisson, and I. Bloch, "Tradi: Tracking deep neural network weight distributions," in *ECCV*, 2020. 1, 3, 7, 8, 10
- [5] S. Fort, H. Hu, and B. Lakshminarayanan, "Deep ensembles: A loss landscape perspective," in *arXiv*, 2019. 1, 3, 4, 5, 6, 7
- [6] Y. Gal, "Uncertainty in deep learning." *PhD Thesis, University of Cambridge*, 2016. 1
- [7] A. Malinin and M. Gales, "Predictive uncertainty estimation via prior networks," in *NeurIPS*, 2018. 1, 4, 7
- [8] U. Bhatt, Y. Zhang, J. Antorán, Q. V. Liao, P. Sattigeri, R. Fogliato, G. G. Melançon, R. Krishnan, J. Stanley, O. Tickoo *et al.*, "Uncertainty as a form of transparency: Measuring, communicating, and using uncertainty," in *AIES*, 2021. 1
- [9] B. Kompa, J. Snoek, and A. L. Beam, "Second opinion needed: communicating uncertainty in medical machine learning," *NPJ Digital Medicine*, 2021. 1
- [10] R. McAllister, Y. Gal, A. Kendall, M. Van Der Wilk, A. Shah, R. Cipolla, and A. Weller, "Concrete problems for autonomous vehicle safety: Advantages of bayesian deep learning." *IJCAI*, 2017. 1
- [11] D. J. MacKay, "A practical bayesian framework for backpropagation networks," *Neural Computation*, 1992. 1, 2, 6
- [12] R. M. Neal, "Bayesian learning for neural networks," *PhD thesis, University of Toronto*, 1995. 1, 2, 6
- [13] M. I. Jordan, Z. Ghahramani, T. S. Jaakkola, and L. K. Saul, "An introduction to variational methods for graphical models," *ML*, 1999. 1, 2, 6
- [14] G. E. Hinton and D. Van Camp, "Keeping the neural networks simple by minimizing the description length of the weights," in *COLT*, 1993. 1, 2, 6
- [15] A. Graves, "Practical variational inference for neural networks," in *NeurIPS*, 2011. 1, 2, 3, 5, 6
- [16] C. Blundell, J. Cornebise, K. Kavukcuoglu, and D. Wierstra, "Weight uncertainty in neural networks," in *ICML*, 2015. 1, 2, 3, 5, 6, 7, 8, 9
- [17] J. M. Hernández-Lobato and R. Adams, "Probabilistic backpropagation for scalable learning of bayesian neural networks," in *ICML*, 2015. 1, 2, 6, 8, 9
- [18] Y. Gal and Z. Ghahramani, "Dropout as a bayesian approximation: Representing model uncertainty in deep learning," in *ICML*, 2016. 1, 6, 7, 8, 9
- [19] A. Mishkin, F. Kunstner, D. Nielsen, M. Schmidt, and M. E. Khan, "Slang: Fast structured covariance approximations for bayesian deep learning with natural gradient," in *NeurIPS*, 2018. 1, 6
- [20] A. Y. Foong, D. R. Burt, Y. Li, and R. E. Turner, "On the expressiveness of approximate inference in bayesian neural networks," in *NeurIPS*, 2020. 1, 3, 5, 6
- [21] Y. Bengio, I. Goodfellow, and A. Courville, *Deep learning*. MIT press, 2017. 1
- [22] T. Salimans and D. P. Kingma, "Weight normalization: A simple reparameterization to accelerate training of deep neural networks," in *NeurIPS*, 2016. 1, 3
- [23] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: A simple way to prevent neural networks from overfitting," *JMLR*, 2014. 1, 3
- [24] Y. Ovadia, E. Fertig, J. Ren, Z. Nado, D. Sculley, S. Nowozin, J. V. Dillon, B. Lakshminarayanan, and J. Snoek, "Can you trust your model's uncertainty? evaluating predictive uncertainty under dataset shift," in *NeurIPS*, 2019. 1, 3, 6
- [25] M. W. Dusenberry, G. Jerfel, Y. Wen, Y. Ma, J. Snoek, K. Heller, B. Lakshminarayanan, and D. Tran, "Efficient and scalable bayesian neural nets with rank-1 factors," in *ICML*, 2020. 1, 3, 4, 5, 6
- [26] C. Louizos and M. Welling, "Structured and efficient variational deep learning with matrix gaussian posteriors," in *ICML*, 2016. 1, 5, 6
- [27] S. Sun, C. Chen, and L. Carin, "Learning structured weight uncertainty in bayesian neural networks," in *AISTATS*, 2017. 1, 5, 6
- [28] G. Zhang, S. Sun, D. Duvenaud, and R. Grosse, "Noisy natural gradient as variational inference," in *ICML*, 2018. 1, 6
- [29] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," in *ICLR*, 2014. 1, 2, 3, 4
- [30] Y. Wen, D. Tran, and J. Ba, "Batchensemble: an alternative approach to efficient ensemble and lifelong learning," in *ICLR*, 2020. 1, 3, 5, 7, 8, 9, 10
- [31] D. J. Rezende, S. Mohamed, and D. Wierstra, "Stochastic backpropagation and approximate inference in deep generative models," in *ICML*, 2014. 2
- [32] S. Qiao, H. Wang, C. Liu, W. Shen, and A. Yuille, "Weight standardization," in *arXiv*, 2019. 3
- [33] A. Mehrtash, P. Abolmaesumi, P. Golland, T. Kapur, D. Wassermann, and W. Wells, "Pep: Parameter ensembling by perturbation," in *NeurIPS*, 2020. 3, 7
- [34] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016. 3, 9, 10
- [35] C. Li, H. Farkhoor, R. Liu, and J. Yosinski, "Measuring the intrinsic dimension of objective landscapes," in *ICLR*, 2018. 3
- [36] S. C. Hora, "Aleatory and epistemic uncertainty in probability elicitation with an example from hazardous waste management," *Reliability Engineering & System Safety*, 1996. 4
- [37] A. G. Wilson and P. Izmailov, "Bayesian deep learning and a probabilistic perspective of generalization," in *NeurIPS*, 2020. 4, 5
- [38] S. Zagoruyko and N. Komodakis, "Wide residual networks," in *BMVC*, 2016. 4, 9
- [39] A. Krizhevsky, "Learning multiple layers of features from tiny images," MIT, Tech. Rep., 2009. 4, 8, 9

- [40] D. Hendrycks and T. Dietterich, "Benchmarking neural network robustness to common corruptions and perturbations," in *ICLR*, 2018. 4, 8, 9
- [41] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng, "Reading digits in natural images with unsupervised feature learning," in *NeurIPS*, 2011. 4, 9
- [42] T. Karalestos and T. D. Bui, "Hierarchical gaussian process priors for bayesian neural network weights," in *NeurIPS*, 2020. 5, 6
- [43] J. Wang, *Geometric structure of high-dimensional data and dimensionality reduction*. Springer, 2012. 5
- [44] P. Izmailov, D. Podoprikin, T. Garipov, D. Vetrov, and A. G. Wilson, "Averaging weights leads to wider optima and better generalization," in *UAI*, 2018. 5
- [45] A. Rame and M. Cord, "Dice: Diversity in deep ensembles via conditional redundancy adversarial estimation," in *ICLR*, 2021. 6
- [46] M. Aksele, "Comparison of classifier selection methods for improving committee performance," in *IWMCS*, 2003. 6
- [47] D. J. MacKay, "Bayesian methods for adaptive models," Ph.D. dissertation, California Institute of Technology, 1992. 6
- [48] T. Pearce, A. Y. Foong, and A. Brintrup, "Structured weight priors for convolutional neural networks," in *ICMLW*, 2020. 6
- [49] C. Riquelme, G. Tucker, and J. Snoek, "Deep bayesian bandits showdown: An empirical comparison of bayesian deep networks for thompson sampling," in *ICLR*, 2018. 6
- [50] T. Karalestos, P. Dayan, and Z. Ghahramani, "Probabilistic meta-representations of neural networks," in *UAI*, 2018. 6
- [51] S. Rossi, S. Marmin, and M. Filippone, "Efficient approximate inference with Walsh-Hadamard variational inference," in *NeurIPS*, 2019. 6
- [52] B.-H. Tran, S. Rossi, D. Milios, and M. Filippone, "All you need is a good functional prior for bayesian deep learning," *JMLR*, 2022. 6, 8, 9
- [53] F. Wenzel, K. Roth, B. S. Veeling, J. Świątkowski, L. Tran, S. Mandt, J. Snoek, T. Salimans, R. Jenatton, and S. Nowozin, "How good is the bayes posterior in deep neural networks really?" in *ICML*, 2020. 6
- [54] V. Fortuin, A. Garriga-Alonso, F. Wenzel, G. Rätsch, R. Turner, M. van der Wilk, and L. Aitchison, "Bayesian neural network priors revisited," in *AABI*, 2021. 6
- [55] S. Lee, S. Purushwalkam, M. Cogswell, D. Crandall, and D. Batra, "Why M heads are better than one: Training a diverse ensemble of deep networks," *arXiv preprint arXiv:1511.06314*, 2015. 6
- [56] G. Huang, Y. Li, G. Pleiss, Z. Liu, J. E. Hopcroft, and K. Q. Weinberger, "Snapshot ensembles: Train 1, get M for free," in *ICLR*, 2017. 7
- [57] T. Garipov, P. Izmailov, D. Podoprikin, D. Vetrov, and A. G. Wilson, "Loss surfaces, mode connectivity, and fast ensembling of DNNs," in *NeurIPS*, 2018. 7
- [58] A. Atanov, A. Ashukha, D. Molchanov, K. Neklyudov, and D. Vetrov, "Uncertainty estimation via stochastic batch normalization," in *ICLRW*, 2018. 7
- [59] M. Havasi, R. Jenatton, S. Fort, J. Z. Liu, J. Snoek, B. Lakshminarayanan, A. M. Dai, and D. Tran, "Training independent subnetworks for robust prediction," in *ICLR*, 2021. 7, 8, 9, 10
- [60] A. Malinin and M. Gales, "Reverse kl-divergence training of prior networks: Improved uncertainty and adversarial robustness," in *NeurIPS*, 2019. 7
- [61] M. Sensoy, L. M. Kaplan, and M. Kandemir, "Evidential deep learning to quantify classification uncertainty," in *NeurIPS*, 2018. 7, 8, 9
- [62] B. Charpentier, D. Zügner, and S. Günnemann, "Posterior network: Uncertainty estimation without ood samples via density-based pseudo-counts," in *NeurIPS*, 2020. 7
- [63] T. Tsiligkaridis, "Information aware max-norm dirichlet networks for predictive uncertainty estimation," *Neural Networks*, 2021. 7
- [64] T. Joo, U. Chung, and M.-G. Seo, "Being bayesian about categorical probability," in *ICML*, 2020. 7, 9
- [65] A. Malinin, S. Chervontsev, I. Provilkov, and M. Gales, "Regression prior networks," in *ICLR*, 2020. 7
- [66] W. K. Hastings, "Monte carlo sampling methods using markov chains and their applications," 1970. 7
- [67] M. Welling and Y. W. Teh, "Bayesian learning via stochastic gradient langevin dynamics," in *ICML*, 2011. 7
- [68] T. Chen, E. Fox, and C. Guestrin, "Stochastic gradient hamiltonian monte carlo," in *ICML*, 2014. 7, 8, 9
- [69] R. Zhang, C. Li, J. Zhang, C. Chen, and A. G. Wilson, "Cyclical stochastic gradient MCMC for bayesian deep learning," in *ICLR*, 2020. 7
- [70] W. Deng, G. Lin, and F. Liang, "A contour stochastic gradient langevin dynamics algorithm for simulations of multi-modal distributions," in *NeurIPS*, 2020. 7
- [71] B. A. Berg and T. Neuhaus, "Multicanonical ensemble: A new approach to simulate first-order phase transitions," *Physical Review Letters*, 1992. 7
- [72] S. Kim, Q. Song, and F. Liang, "Stochastic gradient langevin dynamics algorithms with adaptive drifts," *arXiv preprint arXiv:2009.09535*, 2020. 7
- [73] G. Marceau-Caron and Y. Ollivier, "Natural langevin dynamics for neural networks," in *ICGSI*, 2017. 7
- [74] D. Hendrycks and K. Gimpel, "A baseline for detecting misclassified and out-of-distribution examples in neural networks," in *ICLR*, 2017. 8, 9, 10
- [75] J. Van Amersfoort, L. Smith, Y. W. Teh, and Y. Gal, "Uncertainty estimation using a single deep deterministic neural network," in *ICML*, 2020. 8, 9
- [76] P. Izmailov, W. J. Maddox, P. Kirichenko, T. Garipov, D. Vetrov, and A. G. Wilson, "Subspace inference for bayesian deep learning," in *UAI*, 2020. 8, 9
- [77] S. Sun, G. Zhang, J. Shi, and R. Grosse, "Functional variational bayesian neural networks," in *ICLR*, 2019. 8, 9
- [78] M. Everingham, S. A. Eslami, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, "The Pascal Visual Object Classes challenge: A retrospective," *IJCV*, 2015. 8
- [79] C. Guo, G. Pleiss, Y. Sun, and K. Q. Weinberger, "On calibration of modern neural networks," in *ICML*, 2017. 8
- [80] Y. Wen, G. Jerfel, R. Muller, M. W. Dusenberry, J. Snoek, B. Lakshminarayanan, and D. Tran, "Improving calibration of batchensemble with data augmentation," in *ICMLW*, 2020. 8
- [81] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga *et al.*, "Pytorch: An imperative style, high-performance deep learning library," in *NeurIPS*, 2019. 8
- [82] A. Kendall and Y. Gal, "What uncertainties do we need in bayesian deep learning for computer vision?" in *NeurIPS*, 2017. 9
- [83] D. Nix and A. Weigend, "Estimating the mean and variance of the target probability distribution," in *ICNN*, 1994. 9
- [84] T. DeVries and G. W. Taylor, "Improved regularization of convolutional neural networks with cutout," in *arXiv*, 2017. 9
- [85] D. Hendrycks, S. Basart, M. Mazeika, M. Mostajabi, J. Steinhardt, and D. Song, "A benchmark for anomaly segmentation," in *arXiv*, 2019. 10
- [86] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam, "Encoder-decoder with atrous separable convolution for semantic image segmentation," in *ECCV*, 2018. 10
- [87] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia, "Pyramid scene parsing network," in *CVPR*, 2017. 10
- [88] S. Qiao, H. Wang, C. Liu, W. Shen, and A. Yuille, "Rethinking normalization and elimination singularity in neural networks," in *arXiv*, 2019. 10
- [89] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *ICML*, 2015. 10
- [90] Y. Wu and K. He, "Group normalization," in *ECCV*, 2018. 10
- [91] F. Yu, H. Chen, X. Wang, W. Xian, Y. Chen, F. Liu, V. Madhavan, and T. Darrell, "BDD100K: A diverse driving dataset for heterogeneous multitask learning," in *CVPR*, 2020. 10



Gianni Franchi received a MSc degree in engineering science from Ecole Centrale Marseille in 2013, and a Ph.D. degree in applied mathematics (2016) from PSL University (Mines ParisTech). Then he did a Postdoc at Siegen University from October 2016 to December 2017. Then he did a Postdoc at Paris Saclay University. He is currently pursuing an assistant Professor at ENSTA Paris. His topics of interest are computer vision, machine learning, statistical learning, video understanding.



Andrei Bursuc is a research scientist at valeo.ai and Inria Astra-Vision team in Paris, France. He completed his PhD at Mines ParisTech in 2012. His current research interests concern computer vision and deep learning, in particular learning with limited supervision and predictive uncertainty quantification for autonomous driving.



Emanuel Aldea is associate professor at Paris-Saclay University. He received his BS and MS degrees in Computer Science from Ecole Polytechnique in 2005 and from Paris 6 University in 2006 respectively, and his PhD in image processing from Télécom ParisTech in 2009. His current research interests include image processing and computer vision for autonomous systems.



Séverine Dubuisson received a Master degree (1997) then Ph.D. degree (2000) in System Control from the University of Technology of Compiègne. She has been an associate professor in Sorbone University and is currently an associate professor in Aix-Marseille University, France. Her research interests include computer vision, visual tracking, probabilistic models for video sequence analysis and human interaction.



Isabelle Bloch graduated from the Ecole des Mines de Paris, Paris, France, in 1986, and received the master's degree from the University Paris 12, Paris, in 1987, the Ph.D. degree from the Ecole Nationale Supérieure des Télécommunications (Télécom Paris), Paris, in 1990, and the Habilitation degree from University Paris 5, Paris, in 1995. She has been a Professor at Télécom Paris until 2020 and is now a Professor at Sorbonne Université. Her current research interests include 3-D image understanding, computer vision, mathematical morphology, information fusion, fuzzy set theory, structural, graph-based, and knowledge-based object recognition, spatial reasoning, artificial intelligence, and medical imaging.