COOPERATIVE VARIANCE ESTIMATION AND BAYESIAN NEURAL NETWORKS DISENTANGLE ALEATORIC AND EPISTEMIC UNCERTAINTIES

Anonymous authors

Paper under double-blind review

ABSTRACT

Real-world data contains aleatoric uncertainty – irreducible noise arising from imperfect measurements or from incomplete knowledge about the data generation process. Mean variance estimation (MVE) networks can learn this type of uncertainty but require ad-hoc regularization strategies to avoid overfitting and are unable to predict epistemic uncertainty (model uncertainty). Conversely, Bayesian neural networks predict epistemic uncertainty but are notoriously difficult to train due to the approximate nature of Bayesian inference. We propose to cooperatively train a variance network with a Bayesian neural network and empirically demonstrate that the resulting model disentangles aleatoric and epistemic uncertainties while improving the mean estimation. We demonstrate the effectiveness and scalability of this method across a diverse range of datasets, including a time-dependent heteroscedastic regression dataset we created where the aleatoric uncertainty is known. The proposed method is straightforward to implement, robust, and adaptable to various model architectures.

1 Introduction

Non-probabilistic neural networks that only estimate the mean (expected value) tend to be over-confident and vulnerable to adversarial attacks (Guo et al., 2017; 2019). Quantifying aleatoric (or data) uncertainty alleviates these issues by characterizing noise (Skafte et al., 2019; Seitzer et al., 2022). Estimating epistemic (or model) uncertainty enables active learning and risk-sensitive decision-making (Kendall and Gal, 2017; Depeweg et al., 2018). Consequently, except in cases with negligible or homoscedastic data noise, simultaneously predicting aleatoric and epistemic uncertainties is essential for a wide range of safety-critical applications (Hüllermeier and Waegeman, 2021). In such cases, an outcome with good mean performance but large aleatoric uncertainty may be unacceptable. Therefore, the principle of reducing epistemic uncertainty behind active learning or decision-making needs to be balanced by the respective prediction of aleatoric uncertainty. This is particularly important when the aleatoric uncertainty is heteroscedastic (input-dependent) due to imperfect measurements, environmental variability, and other factors (Smith et al., 2024).

Summary of contributions. We propose a cooperative learning strategy for uncertainty disentanglement based on sequential training of (1) a mean network, (2) a variance network, and (3) a Bayesian neural network. Figure 1 illustrates the method for one-dimensional heteroscedastic regression, briefly describing it at the figure caption.

2 RELATED WORK

2.1 ALEATORIC UNCERTAINTY

Aleatoric uncertainty can be estimated by parametric or nonparametric models. The latter do not explicitly define the likelihood function and instead focus on learning to sample from the data distribution (Mohamed and Lakshminarayanan, 2016). Their ability to estimate nontrivial aleatoric uncertainty comes at the cost of training difficulties and sampling inefficiencies (Sensoy et al., 2020;

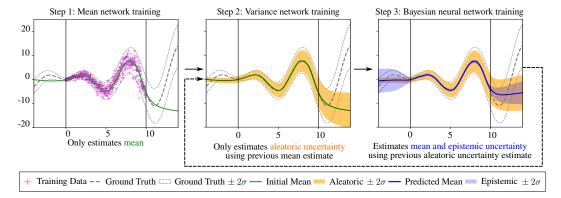


Figure 1: Illustration of the proposed cooperative training strategy leading to Variance estimation Bayesian Neural Networks (VeBNNs). The left figure shows the unseen ground truth mean as well as the respective training data. The method starts by training the mean network to only estimate the mean (green solid line in Step 1). Then, without updating the mean estimate, a variance network is trained to only predict aleatoric uncertainty (orange credible interval in Step 2). Subsequently, considering this aleatoric uncertainty estimation, a Bayesian neural network is trained to obtain an updated mean and corresponding epistemic uncertainty (solid blue line for the new mean, and shaded blue credible interval for the epistemic uncertainty in Step 3). If needed, the method can iterate between steps 3 and 2 to improve the disentanglement of uncertainties. Note the disentanglement of uncertainties together with the improvement of the mean estimation away from the data support (x < 0 and x > 10 in Step 3).

Harakeh et al., 2023). Therefore, parametric models are more common. They assume a parameterized observation distribution, usually a Gaussian as in Mean Variance Estimation (MVE) networks (Nix and Weigend, 1994), and learn the corresponding parameters by minimizing the Negative Log-Likelihood (NLL) loss with associated regularization. Similar parametric models have been developed, replacing the Gaussian distribution with other distributions (Meyer and Thakurdesai, 2020). However, training these models can be challenging. MVE networks have been observed to lead to good mean but overconfident variance estimations in regions within the data support, and exhibit generalization issues outside these regions (Skafte et al., 2019). As analyzed by different authors (Skafte et al., 2019; Seitzer et al., 2022; Immer et al., 2023; Sluijterman et al., 2024), these issues result from the loss gradients having very different magnitudes when calculated with respect to the mean or the aleatoric variance, as seen in Equation (3), which creates imbalances when minimizing the NLL.

Different solutions have been proposed to improve the parametric estimation of aleatoric uncertainty, including a Bayesian treatment of variance (Stirn and Knowles, 2020), calibration by distribution matching via maximum mean discrepancy loss (Cui et al., 2020), or modifying the loss function to include an additional balance hyperparameter (Seitzer et al., 2022). However, Sluijterman et al. (2024) demonstrated that training an MVE network that simultaneously predicts mean and variance leads to significantly worse predictions for both estimations compared to separately training, even when considering the above-mentioned modified losses (we independently reached the same conclusion while conducting our work; see Appendix A.1). Nevertheless, we note that parametric deterministic models are insufficient if they only estimate mean and aleatoric uncertainty, without accounting for epistemic uncertainty. This is also visible in Figure 1 (Step 2), as the model has an overconfident mean outside the data support (see x > 10). Furthermore, the inability to estimate epistemic uncertainty is problematic in itself, as discussed in Section 1.

2.2 Epistemic uncertainty

Epistemic uncertainty is usually estimated by probabilistic models that impose a prior distribution on the model parameters (Abdar et al., 2021). Instead of finding point estimates as in deterministic models, they predict a posterior predictive distribution (PPD). Unfortunately, accurate and computationally tractable determination of the PPD is challenging in most cases, except for a few models like Gaussian processes, where inference can be done exactly under strict assumptions and with limited data scalability (Rasmussen and Williams, 2005). Bayesian neural networks (BNNs) are more scalable, but the accuracy and scalability are strongly dependent on the type of Bayesian inference (Wilson and Izmailov, 2020; Wenzel et al., 2020; Izmailov et al., 2021).

Bayesian inference is commonly done by Markov Chain Monte Carlo (MCMC) sampling methods (Neal, 1995; Welling and Teh, 2011; Li et al., 2016) or Variational Inference (VI) methods that are faster to train but less accurate (Graves, 2011; Blundell et al., 2015). Avoiding formal Bayesian inference is also possible by adopting ensemble methods such as Monte Carlo (MC) Dropout (Gal and Ghahramani, 2016) and deep ensembles (Lakshminarayanan et al., 2017; Fort et al., 2020). Although not strictly Bayesian, MC-Dropout can be interpreted as a Variational Bayesian approximation that has additional scalability but even lower accuracy (no free lunch).

The practical difficulties of training BNNs have limited their widespread use. Therefore, they are often trained by disregarding aleatoric uncertainty or assuming homoscedastic noise, which can be set as a hyperparameter but is impractical for complex problems (Abdar et al., 2021). Instead, the end-to-end training (also referred to as second order uncertainty quantification) that trains an MVE network by MC-Dropout (Kendall and Gal, 2017) or deep ensembling (Wang et al., 2025) is reported to approximately disentangle aleatoric and epistemic uncertainties for heteroscedastic regression and classification. However, the essential challenges faced when training MVE networks are not solved by ensembling them, and the lack of accuracy associated with MC-Dropout raises questions about their ability to make high-quality predictions and truly disentangle epistemic and aleatoric uncertainties (Valdenegro-Toro and Mori, 2022; Mucsányi et al., 2024). This has motivated other authors to propose different solutions. A non-Bayesian solution to separate uncertainties was proposed by introducing a high-order evidential distribution, i.e., considering priors over the likelihood function instead of over network weights (Amini et al., 2020). Another proposal has been to use a natural reparameterization combined with an approximate Laplace expansion to estimate epistemic uncertainty (Immer et al., 2023). Still, a recent investigation (Mucsányi et al., 2024) has shown that no current method achieves reliable uncertainty disentanglement.

3 METHODOLOGY

3.1 Preliminaries

Consider a dataset $\mathcal{D} = \{\mathbf{x}_n, y_n\}_{n=1}^N$ with *i.i.d.* data points, where $\mathbf{x}_n \in \mathbb{R}^d$ represents the input features and $y_n \in \mathbb{R}$ the corresponding outputs¹. The heteroscedastic regression problem can be formulated as:

$$y = f(\mathbf{x}) + s(\mathbf{x}) \tag{1}$$

where $f(\mathbf{x})$ denotes the underlying noiseless ground truth function (expected mean), and $s(\mathbf{x})$ is the corresponding heteroscedastic noise (aleatoric uncertainty). If the noise is Gaussian, then $s(\mathbf{x}) \sim \mathcal{N}(0, \varepsilon^2(\mathbf{x}))$ where $\varepsilon^2(\mathbf{x})$ represents its ground truth input-dependent variance.

3.2 CHALLENGES OF ENSEMBLING TRAINING OF MEAN VARIANCE NETWORKS

The NLL loss used for training MVE networks resulting from a Gaussian observation distribution with heteroscedastic aleatoric uncertainty is:

$$\mathcal{L}_1(\boldsymbol{\theta}, \boldsymbol{\phi}) = \sum_{n=1}^{N} \left[\frac{(y_n - \mu(\mathbf{x}_n; \boldsymbol{\theta}))^2}{2\sigma_a^2(\mathbf{x}_n; \boldsymbol{\phi})} + \frac{1}{2} \log \left(\sigma_a^2(\mathbf{x}_n; \boldsymbol{\phi}) \right) \right]$$
(2)

where $\sigma_a^2(\mathbf{x}; \boldsymbol{\phi}) > 0$ is the parameterized aleatoric variance, and $\mu(\mathbf{x}; \boldsymbol{\theta})$ the mean. The derivatives with respect to $\mu(\mathbf{x}; \boldsymbol{\theta})$ and $\sigma_a^2(\mathbf{x}; \boldsymbol{\phi})$ become:

$$\nabla_{\mu} \mathcal{L}_{1} = \sum_{n=1}^{N} \left(\frac{\mu(\mathbf{x}_{n}; \boldsymbol{\theta}) - y_{n}}{\sigma_{a}^{2}(\mathbf{x}_{n}; \boldsymbol{\phi})} \right), \ \nabla_{\sigma_{a}^{2}} \mathcal{L}_{1} = \sum_{n=1}^{N} \left(\frac{\sigma_{a}^{2}(\mathbf{x}_{n}; \boldsymbol{\phi}) - (y_{n} - \mu(\mathbf{x}_{n}; \boldsymbol{\theta}))^{2}}{\left(\sigma_{a}^{2}(\mathbf{x}_{n}; \boldsymbol{\phi})\right)^{2}} \right). \tag{3}$$

According to Equation (3), the gradient with respect to $\sigma_a^2(\mathbf{x}; \phi)$ has a high-order term in the denominator that makes the optimization more challenging and causes an imbalance when minimizing

¹The equations become simpler when writing for one-dimensional outputs, but this article includes examples with multi-dimensional outputs $\mathbf{y}_n \in \mathbb{R}^m$, as shown later.

the loss. As a result, the end-to-end training with of an MVE network would include the following unexpected cases: (1) if the MVE network overfits the response, i.e., $(y_n - \mu(\mathbf{x}_n; \boldsymbol{\theta}))^2 \to 0$, $\sigma_a^2(\mathbf{x}_n; \boldsymbol{\phi}) \to 0$, the loss function \mathcal{L}_1 may become undefined; (2) if the MVE network leads to large variance estimates when quantifying aleatoric uncertainty, then $\sigma_a^2(\mathbf{x}_n; \boldsymbol{\phi}) \to \infty$ and the gradients with respect to $\sigma_a^2(\mathbf{x}_n; \boldsymbol{\phi})$ vanish. The loss function may also contain multiple saddle points that translate into low quality point estimates for the parameters and slow convergence, as training may become trapped in degenerate regions. While balancing the loss is possible, as proposed in Seitzer et al. (2022) by the β -NLL loss, this introduces additional hyperparameters and still does not address the difficulty in conducting Bayesian inference when trying to simultaneously predict aleatoric and epistemic uncertainties (Mucsányi et al., 2024).

3.3 VEBNN: COOPERATIVE VARIANCE ESTIMATION BAYESIAN NEURAL NETWORKS

We propose a sequential training strategy (see Algorithm 1) that starts with training a mean network, then a variance estimation network that predicts aleatoric uncertainty, followed by a BNN that updates both mean and epistemic uncertainty for the previously determined aleatoric uncertainty. By separating the roles of each network and ensuring their cooperative training, we empirically demonstrate that the resulting BNN can learn and improve all three estimates. Importantly, training each network separately is easier, and we show that inference of the BNN is also facilitated due to the presence of a good estimate of aleatoric uncertainty (that is not being learned at that stage).

Algorithm 1: Cooperative VeBNN training

Data: mean network $\mu(\mathbf{x}; \boldsymbol{\theta})$, variance network $\sigma_a^2(\mathbf{x}; \boldsymbol{\phi})$, dataset $\mathcal{D} = \{\mathbf{X}, \mathbf{y}\}$, number of iterations K **Step 1: Mean network training:** Minimize Equation (2) to find point estimate of $\hat{\mu}(\mathbf{x}; \boldsymbol{\theta})$ **for** i = 1 **to** K **do**

Step 2: variance network training (Aleatoric uncertainty): Minimize Equation (6) for fixed mean from Step 1 or Step 3.

Step 3: Bayesian network training (Epistemic uncertainty): Sample posterior from Equation (9) to determine mean and epistemic variance estimates for fixed aleatoric variance from Equation (8); and compute the log marginal likelihood: $\mathtt{LMglk}[i] = \log \mathbb{E}_{p(\boldsymbol{\theta}|\mathcal{D})} \left[p(\boldsymbol{y} \mid \boldsymbol{\theta}^{(i)}, \boldsymbol{\phi}^{(i)}) \right]$

end

Identify the optimal parameters for θ^* and ϕ^* by $i^* = \arg \max_i \text{LMglk}[i]$

The proposed Algorithm 1 starts by not considering aleatoric uncertainty, i.e. $\sigma_a^2(\mathbf{x}; \phi) = \text{constant}$, and conventionally training the mean network by finding the maximum a posteriori (MAP) estimate of the parameters $\boldsymbol{\theta}$ of Equation (2), hence determining only the mean.

3.3.1 Variance estimation (VE) network training

Once the mean is obtained, we then train the variance estimation (Ve) network for this fixed mean. There is, however, an important detail that facilitates training of this network (Step 2 in Algorithm 1). The variance estimation network does not directly output $\sigma_a^2(\mathbf{x}; \phi)$, and so its parameters are not directly determined by minimizing Equation (2) and keeping the mean fixed. Instead, the variance network outputs the residual $r = (\mu(\mathbf{x}; \theta) - y)^2$, which follows a Gamma distribution because y is Gaussian.

Assumption 3.1. $(\mu(\mathbf{x}; \boldsymbol{\theta}) - f(\mathbf{x}))^2$ is finite and tends to 0 when $N \to \infty$. This follows from assuming unbiased or consistent estimations for the ground truth $f(\mathbf{x})$, regardless of noise.

Proof. We aim to demonstrate that from Assumption 3.1 and for y following a Gaussian distribution, the squared residual $r = (\mu(\mathbf{x}; \boldsymbol{\theta}) - y)^2$ follows a Gamma distribution.

Since $y \mid \mathbf{x} \sim \mathcal{N}(f(\mathbf{x}), \varepsilon^2(\mathbf{x}))$, the residual $\mu(\mathbf{x}; \boldsymbol{\theta}) - y \sim \mathcal{N}(0, \varepsilon^2(\mathbf{x}))$. By standardizing, we obtain:

$$Z = \frac{\mu(\mathbf{x}; \boldsymbol{\theta}) - y}{\varepsilon(\mathbf{x})} \sim \mathcal{N}(0, 1). \tag{4}$$

Thus, the squared residual can be obtained: $r = (\mu(\mathbf{x}; \boldsymbol{\theta}) - y)^2 = \varepsilon^2(\mathbf{x})Z^2$. Since $Z \sim \mathcal{N}(0, 1)$, we have $Z^2 \sim \chi^2(1)$, a specific case of Gamma distribution $Z^2 \sim \text{Gamma}\left(\frac{1}{2}, \frac{1}{2}\right)$. Since a Gamma

random variable $W \sim \text{Gamma}(\alpha, \lambda)$ scaled by a constant c > 0 gives $cW \sim \text{Gamma}(\alpha, \frac{\lambda}{c})$, then:

$$r \sim \text{Gamma}\left(\frac{1}{2}, \frac{1}{2\varepsilon^2(\mathbf{x})}\right)$$
 (5)

Showing that the mean of the Gamma distribution becomes $\frac{\alpha}{\lambda} = \varepsilon^2(\mathbf{x})$, i.e., the aleatoric variance.

We propose the Gamma likelihood to model the squared residual r, leading to the corresponding NLL loss to train the variance network:

$$\mathcal{L}_{2}(\phi) = \sum_{n=1}^{N} \left[\alpha(\mathbf{x}_{n}; \phi) \log \frac{\lambda(\mathbf{x}_{n}; \phi)}{\Gamma(\alpha(\mathbf{x}_{n}; \phi))} - (\alpha(\mathbf{x}_{n}; \phi) - 1) \log r_{n} + \frac{\lambda(\mathbf{x}_{n}; \phi)}{r_{n}} \right]$$
(6)

where r are the above-mentioned residuals, and with the shape and rate parameters of the Gamma distribution, $\alpha(\mathbf{x}; \phi) > 0$ and $\lambda(\mathbf{x}; \phi) > 0$, being the outputs of the network. These parameters are estimated via MAP, and their gradients with respect to α and λ contain no high-order terms, as shown below:

$$\nabla_{\alpha(\mathbf{x})} \mathcal{L}_2 = \sum_{n=1}^{N} \left[\log \frac{\lambda(\mathbf{x})}{\Gamma(\alpha(\mathbf{x}))} - \alpha(\mathbf{x}) \psi(\alpha(\mathbf{x})) - \log r_n \right], \nabla_{\lambda(\mathbf{x})} \mathcal{L}_2 = \sum_{n=1}^{N} \left[\frac{\alpha(\mathbf{x})}{\lambda(\mathbf{x})} + \frac{1}{r_n} \right]. \tag{7}$$

The expected value of the Gamma distribution becomes the desired heteroscedastic variance:

$$\sigma_a^2(\mathbf{x}; \boldsymbol{\phi}) = \frac{\alpha(\mathbf{x}; \boldsymbol{\phi})}{\lambda(\mathbf{x}; \boldsymbol{\phi})}$$
(8)

3.3.2 BAYESIAN NEURAL NETWORK TRAINING

Having determined the mean and aleatoric variance estimates, we then train a BNN with a warm-start for the mean, and assuming fixed aleatoric variance that is obtained from Equation (8). In principle, the same network architecture can be used for the BNN and the mean network². The posterior of the BNN is determined by Bayes' rule, and it is proportional to the product of likelihood and prior:

$$p(\boldsymbol{\theta} \mid \mathcal{D}) \propto p(\mathcal{D} \mid \boldsymbol{\theta}) p(\boldsymbol{\theta}) \tag{9}$$

where $p(\theta)$ is the prior over the neural network parameters, and $p(\mathcal{D} \mid \theta)$ is the likelihood for the observations. In this work, we consider a Gaussian prior with zero mean and unit variance, and the likelihood is given by Equation (2), i.e., it arises from a Gaussian observation distribution with heteroscedastic noise. The logarithm of the posterior becomes:

$$\log p\left(\boldsymbol{\theta} \mid \mathcal{D}\right) = \sum_{n=1}^{N} \left[\log \frac{1}{\sqrt{2\pi\sigma_a^2(\mathbf{x}_n; \boldsymbol{\phi})}} - \frac{\left(\mathbf{y}_n - \mu(\mathbf{x}_n; \boldsymbol{\theta})\right)^2}{2\sigma_a^2(\mathbf{x}_n; \boldsymbol{\phi})} \right] + m \log \frac{1}{\sqrt{2\pi/\kappa}} - \frac{\kappa}{2} \|\boldsymbol{\theta}\|^2$$
 (10)

where κ is the precision of the prior distribution and m is the length of θ . Note that Equation (10) is the same as Equation (2), but now we explicitly include the prior terms. Appendix B summarizes common Bayesian inference strategies, including the MCMC-based methods that sample directly from Equation (10), and VI methods that minimize the evidence lower bound (ELBO) (Appendices B.1 and B.2, respectively). From experience, preconditioned Stochastic Gradient Langevin Dynamics (pSGLD) (Li et al., 2016) is expected to be a good choice for this BNN because the aleatoric uncertainty is fixed, and the likelihood and prior both Gaussian, leading to a Gaussian posterior. Importantly, since the aleatoric variance $\sigma_a^2(\mathbf{x}_n; \phi)$ is fixed from step 2 and is regarded as constant, it converges to the correct mode $\mu(\mathbf{x}_n; \theta)$, and has better and faster convergence (no gradient issue).

From the PPD of the BNN, we then estimate the mean, aleatoric, and epistemic variances (disentangled) for any unseen point \mathbf{x}' based on Bayesian model averaging:

$$\mathcal{N}\left(\mathbf{y}' \middle| \underbrace{\mathbb{E}_{p(\boldsymbol{\theta}|\mathcal{D})}\left[\mu(\mathbf{x}';\boldsymbol{\theta})\right]}_{\text{Predictive Mean}}, \underbrace{\sigma_a^2(\mathbf{x}';\boldsymbol{\phi})}_{\text{Aleatoric}} + \underbrace{\mathbb{V}_{p(\boldsymbol{\theta}|\mathcal{D})}\left[\mu(\mathbf{x}';\boldsymbol{\theta})\right]}_{\text{Epistemic}}\right)$$
(11)

²Estimating epistemic uncertainty with BNNs can require a network with wider hidden layers when compared to a deterministic network that only estimates the mean. So, choosing a smaller mean network in Step 1 is also possible. However, we like the simplicity of not introducing a new network.

4 EXPERIMENTS

Datasets We consider four distinct sets of datasets (a total of 18 datasets): (1) the previously discussed one-dimensional illustrative example (Skafte et al., 2019), with results presented in Appendix D.1; (2) UCI regression datasets (Hernandez-Lobato and Adams, 2015); (3) large-scale image regression datasets (Gustafsson et al., 2023); and (4) our own dataset obtained from computer simulations of materials undergoing history-dependent deformations (material plasticity law discovery dataset) with known ground-truth of aleatoric uncertainty. Readers interested in more details about the dataset are referred to Appendix F.

Baselines We compare against representative methods: (1) **ME** (**MSE**): standard Mean Estimation network; (2) **MVE** (β -**NLL**): MVE using the β -NLL loss (Seitzer et al., 2022), which is capable of simultaneous prediction of mean and aleatoric uncertainty. Concerning methods that aim at disentangling uncertainties, we considered (3) **Evidential**: Deep Evidential regression (Amini et al., 2020); (4) **MVE** (**Ensembles**): ensembling (Lakshminarayanan et al., 2017) MVE networks; (5) **MVE** (**MC-Dropout**): MC-Dropout (Kendall and Gal, 2017) training for MVE networks. (6) **BNN-Endto-End**: end-to-end training of BNN with different inference methods, i.e., **pSGLD** (Li et al., 2016) and Bayes By Backpropagation (**BBB**) (Blundell et al., 2015), indicated in a parenthesis for every Figure and Table. (7) **BNN-Homo**: BNN training with homoscedastic noise assumption. Finally, our proposed **VeBNN** method is also trained with all the above-mentioned inference methods.

Details about accuracy metrics, additional results, and training are presented in Appendix C, Appendix D, and Appendix E, respectively.

4.1 REAL WORLD DATASETS WITH UNKNOWN ALEATORIC UNCERTAINTY

An important challenge in assessing the performance of methods that disentangle uncertainties is that existing datasets do not provide the ground-truth aleatoric uncertainty. In this section, we consider 16 different datasets that have been used in different papers on the topic, despite having **unknown aleatoric uncertainty**. Therefore, we caution that **the datasets of this section can only be used to assess the quality of the mean and total uncertainty estimations**.

UCI regression datasets The results for the UCI regression datasets are summarized in Table 1 (and also in Table 2 of Appendix D.2) according to the TLL and RMSE metrics, respectively. The mean estimates (RMSE) are broadly similar, as expected, although they improve with our proposed strategy for every inference method used when compared to training a single network with both BNN-End-to-End and BNN-Homo. We also note that the mean estimates are better for VeBNN (pSGLD) than the ME (MSE). The improvements in terms of TLL are more noteworthy, where the best performance is also for VeBNN (pSGLD).

Table 1: TLL (\uparrow) results of UCI regression datasets. The best performance for each dataset is underlined and bold, and the second best is in bold.

Methods	Carbon (10721, 7,1)	Concrete (1030, 8,1)	Energy (768,8,2)	Boston (506,13,1)	Power (9568,4,1)	Superconduct (21263,81,1)	Wine-Red (1599, 11,1)	Yacht (308,6,1)
$\begin{aligned} & \text{MVE} \ (\beta_{\text{NLL}} = 1.0) \\ & \text{MVE} \ (\text{Ensemble}) \\ & \text{Evidential} \end{aligned}$	3.79 (0.25)	-3.58 (1.74)	-1.14 (0.26)	-2.93 (0.66)	-2.82 (0.12)	-3.81 (0.22)	-0.97 (0.12)	-1.98 (1.27)
	-10.41 (49.67)	-3.45 (0.55)	- 0.93 (0.34)	-4.39 (1.44)	-2.77 (0.08)	-3.43 (0.04)	-1.72 (0.31)	-1.04 (1.01)
	2.05 (0.31)	-3.60 (0.48)	-2.39 (0.43)	-3.57 (0.60)	-3.60 (0.65)	-4.02 (0.39)	-1.46 (0.54)	-2.93 (0.54)
MVE (MC-Dropout) Ours: VeBNN (MC-Dropout)	2.06 (0.12)	-2.99 (0.12)	-1.63 (1.63)	-2.52 (0.22)	-2.82 (0.03)	-3.51 (0.44)	-1.00 (0.24)	-0.61 (0.23)
	2.49 (0.03)	- 2.95 (0.15)	-1.43 (0.06)	-2.63 (0.41)	-2.80 (0.03)	-3.39 (0.12)	-1.01 (0.12)	-1.24 (0.08)
BNN-Homo (BBB)	1.15 (0.03)	-3.17 (0.13)	-2.04 (0.15)	-2.72 (0.21)	-2.89 (0.02)	-4.04 (0.04)	-0.98(0.09)	-2.75(0.05)
BNN-End-to-End (BBB)	-1.99 (0.00)	-6.36 (0.35)	-6.03 (0.42)	-7.99 (0.28)	-5.72 (1.00)	-7.48 (1.89)	-3.12 (0.49)	-7.57 (0.31)
Ours: VeBNN (BBB)	3.90 (0.33)	-3.10 (0.25)	-1.06 (0.11)	-3.17 (1.01)	-2.79 (0.04)	-3.57 (0.10)	-0.93 (0.08)	-1.45 (0.18)
BNN-Homo (pSGLD)	1.46(0.00)	-3.03(0.09)	-2.15(0.01)	-2.62(0.18)	-2.86(0.02)	-3.83(0.03)	-0.95(0.08)	-2.64(0.03)
BNN-End-to-End (pSGLD)	0.20 (5.96)	-3.03 (0.35)	-1.16 (0.51)	-2.57 (0.26)	-2.75 (0.06)	-3.48 (0.54)	-0.94 (0.13)	-0.64 (0.28)
Ours: VeBNN (pSGLD)	3.95 (0.42)	-2.91 (0.25)	- 0.83 (0.08)	-2.92 (0.59)	-2.74 (0.04)	-3.40 (0.09)	-0.93 (0.09)	-1.29 (0.09)

Image regression datasets We considered 8 large-scale image regression datasets under real-world distribution shifts (Gustafsson et al., 2023) that evaluate the mean and total uncertainty estimations. These problems were solved using a ResNet 34 network architecture (He et al., 2016). Due to space limitations, all results are presented in Appendix D.3. The conclusions are similar to

those obtained from the UCI regression datasets. Our cooperative training strategy improves performance when compared to training a single network, whether considering an MVE network with MC-Dropout or a BNN trained end-to-end with pSGLD – see Table 3. As before, the best results are obtained from the proposed VeBNN (pSGLD) method. Compared with the strongest baseline, MVE (Ensembles), we observe similar performance in estimating the mean and total uncertainty when considering all metrics in Table 3. This is very encouraging because training is faster for VeBNN (pSGLD) when compared to MVE (Ensembles), as the former collects samples in the same training procedure, while ensembling requires collecting one sample per training of a single MVE (i.e., training restarts many times).

4.2 NEW DATASET WITH KNOWN ALEATORIC UNCERTAINTY: MATERIAL PLASTICITY LAW DISCOVERY

Compare to baselines The results for this dataset are shown in Figure 2. Note that results using BBB inference are not available due to a lack of convergence. The reader is also referred to Figure 21 in Appendix F that shows a typical ground truth response with one input sequence (material deformation path) and corresponding output sequence (stochastic material response along that path) where each path results from 100 points obtained from a Physics simulator.

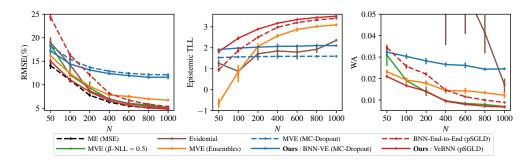


Figure 2: Accuracy metrics (RMSE \downarrow , Epistemic TLL \uparrow , and WA \downarrow) obtained for the plasticity law discovery dataset considering a training set with different number N of training sequences. The Wasserstein distance (WA) represents the closeness of the estimated aleatoric uncertainty distribution to the ground truth distribution. All metrics result from repeating the training of each method 5 times by randomly resampling points from the training datasets.

Figure 2 shows the accuracy metrics obtained with different training sequences taken from the training dataset where N increases from 50 to 1000. As expected, more training data improves accuracy for all methods, but it is clear that the proposed VeBNN (pSGLD) method achieves better predictions for all metrics. Figure 2 also demonstrates the performance improvements obtained from the proposed cooperative training strategy (solid lines) compared to others (dashed lines) for both inference types: MC-Dropout and pSGLD. Note that the proposed VeBNN (pSGLD) has comparable performance with the ME (MSE) when estimating the mean (RMSE metric), and similar Wasserstein distance to the MVE (β -NLL = 0.5) network, the best method among MVEs, when estimating the aleatoric uncertainty.

Figure 3 shows the predictions of the proposed VeBNN (pSGLD) and its correct identification of the disentangled data uncertainties, which improves as the training data increases (upper column to bottom column). As expected, the proposed VeBNN (pSGLD) outperforms BNN-End-to-End (pSGLD) for the same number of training sequences, and the estimated epistemic uncertainty decreases with increasing training data. The MVE (Ensembles) performs well when considering the larger training dataset (800 training sequences), but its prediction is significantly worse than VeBNN (pSGLD) for 50 training sequences (this is also seen in Figure 2). On the contrary, Evidential has significant difficulties with this problem. Predictions for other methods are given in Figure 16, in which the proposed cooperative training strategy also considerably improves predictions when using MC-Dropout as inference.

We also observe in Figure 2 that the aleatoric uncertainty estimated by the VeBNN (pSGLD) converges for N > 400, since WA remains stable. Furthermore, the VeBNN also correctly disentangles

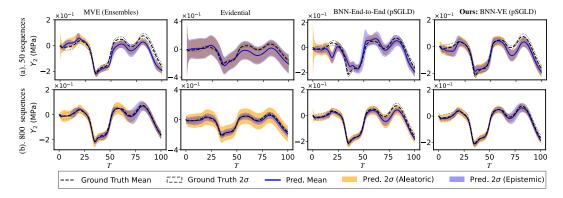


Figure 3: **Predictions of different methods on plasticity law discovery dataset.** We randomly pick one test point from the dataset and show the entire third component y_2 of 100-time steps under the 50 and 800 training sequences, respectively.

the uncertainties, as shown in Figure 3. Unfortunately, we have not found other datasets reporting the ground-truth aleatoric uncertainty like the dataset we created herein. Nevertheless, we find these results very encouraging, and we hope our dataset will motivate future developments in the field.

Ablation study for Gamma Loss and number of iterations K An ablation study is conducted for this dataset (plasticity law discovery) because the ground-truth aleatoric uncertainty is known. The results are presented in Figure 4. The proposed training strategy significantly improves upon the state-of-the-art even when training the networks only once (i.e., K=1). Interestingly, we also noticed that training converged for every dataset with only one additional iteration (K=2). To further explore this, we show the trajectory of all essential components of VeBNN with respect to the iteration K from 1 to 5 for 800 in Figure 18. These results confirm that K is merely the iteration count until convergence rather than a tunable hyperparameter.

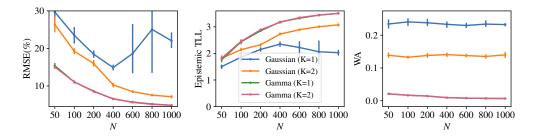


Figure 4: Ablation study of the Gamma loss and iteration parameter K using VeBNN (pSGLD). Curves labeled *Gamma* and *Gaussian* correspond to Step 2 training with the proposed Gamma loss (Equation (6)) and the original Gaussian NLL loss (Equation (2)), respectively, while fixing $\mu(\mathbf{x}; \boldsymbol{\theta})$ from Step 1.

Finding $\sigma_a^2(\mathbf{x}; \boldsymbol{\phi})$ with a Gaussian likelihood loss in Step 2 leads to worst performance due to the high-order term in the gradient calculation (Equation (3)) that can lead to $\sigma_a^2(\mathbf{x}) \to \infty$ and a degenerate local optimum for which the gradient tends to zero (Equation (3)). In contrast, the proposed Gamma NLL loss addresses this and yields a stable optimization process (Figure 4).

Ablation study for the regularization of the Gaussian likelihood of Step 3 Figure 5 compares VeBNN (pSGLD) with BNNs trained under different likelihood regularization schemes. While β -NLL achieves good mean and aleatoric performance (Figure 2), it fails to disentangle uncertainties in a Bayesian setting (Figure 5) and is sensitive to the choice of β . This is consistent with prior work on posterior tempering Wenzel et al. (2020), which shows that likelihood scaling does not necessarily improve inference. In contrast, MSE and Homo-NLL yield poor uncertainty estimates, as they rely on unrealistic unit-variance or constant-variance assumptions in heteroscedastic regression. Finally,

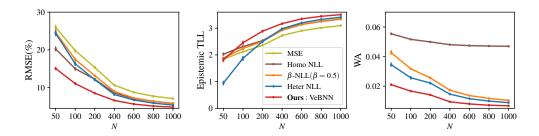


Figure 5: Comparison with BNN end-to-end (pSGLD) training under different likelihood regularization schemes. MSE: unit variance; Homo-NLL: constant variance via grid search; β -NLL: best β via grid search; Heter-NLL: original likelihood.

VeBNN (pSGLD) performs Bayesian inference using the variance learned in Step 2, achieving better mean and epistemic uncertainty—even under the original NLL—due to the correct identification of aleatoric uncertainty.

5 DISCUSSIONS

Size of variance estimation network The variance estimation network used for estimating aleatoric uncertainty is trained separately from the BNN. However, we believe that this is an advantage, as the architecture required to learn aleatoric uncertainty separately is simpler than when training for everything at once. We considered 8 different variance network configurations and observed robust estimations of aleatoric uncertainty, as shown in Figure 19 (see Appendix D.6).

Computational cost As shown in Experiments, our cooperative training converges quickly: the mean network is trained once, while the variance network and Bayesian inference are performed twice. The mean network quickly finds a posterior mode, facilitating MCMC sampling without a long burn-in phase. For the best inference method (pSGLD), the number of training epochs is comparable to ME and MVE training, as shown in Appendix E. In contrast, considering homoscedastic noise and treating it as a hyperparameter typically requires many full training runs and performs worse. Compared to Deep Ensembles, which yield only one posterior sample per run, our method is more efficient. End-to-end training has similar efficiency, but our method consistently achieves better accuracy. Since deterministic ME and VE training is inexpensive relative to BNN training, the additional cost of our strategy is negligible and benefits from the warm start of the mean.

5.1 LIMITATIONS

Although the proposed method facilitates Bayesian inference for BNNs, training BNNs remains more challenging than training deterministic networks. When the posterior is multimodal, non-smooth, or only partially known, inference of epistemic uncertainty under a fixed Gaussian prior–likelihood pair can become unreliable. Recent advances in Imprecise Probabilistic Machine Learning (IPML), such as Credal Bayesian Deep Learning (CBDL) (Caprio et al., 2023; Wang et al., 2024) provide broader epistemic coverage but at significantly higher computational cost. Our framework may serve as a lightweight foundation that can be extended with IPML techniques in future work. Importantly, our results show that Bayesian inference is easier when aleatoric uncertainty is determined separately, as proposed herein.

6 CONCLUSION

We propose a novel Bayesian heteroscedastic regression strategy based on cooperative training of variance estimation and Bayesian neural networks that disentangles epistemic and aleatoric uncertainties. The proposed method is efficient, robust, and straightforward to implement because it is simpler to train each network in isolation, while ensuring complementarity in their training. As a Bayesian method, it does not require validation data. We believe the method is scalable and applicable to real-world problems involving active learning and Bayesian optimization, considering data-scarce and data-rich scenarios, unlike Gaussian process regression.

DECLARATION OF GENERATIVE AI AND AI-ASSISTED TECHNOLOGIES IN THE WRITING PROCESS

The authors used CHATGPT to IMPROVE LANGUAGE AND READABILITY. After using these tools/services, the authors reviewed and edited the content as needed and assumed full responsibility for the content of the publication.

ETHICS STATEMENT

This research does not involve human participants, personally identifiable data, or sensitive social impacts. The authors are not aware of any ethical concerns associated with the methods or experiments presented.

REPRODUCIBILITY STATEMENT

All experimental settings, hyperparameters, and architectures are described in detail in the paper and appendix. We also provide core code to replicate the illustrative example in the supplementary material to facilitate full reproducibility.

REFERENCES

- Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural networks. In *International conference on machine learning*, pages 1321–1330. PMLR, 2017.
- Chuan Guo, Jacob Gardner, Yurong You, Andrew Gordon Wilson, and Kilian Weinberger. Simple black-box adversarial attacks. In *International conference on machine learning*, pages 2484–2493. PMLR, 2019.
- Nicki Skafte, Martin Jø rgensen, and Sø ren Hauberg. Reliable training and estimation of variance networks. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- Maximilian Seitzer, Arash Tavakoli, Dimitrije Antic, and Georg Martius. On the pitfalls of heteroscedastic uncertainty estimation with probabilistic neural networks, 2022.
- Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for computer vision? In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- 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*, pages 1184–1193. PMLR, 2018.
- Eyke Hüllermeier and Willem Waegeman. Aleatoric and epistemic uncertainty in machine learning: An introduction to concepts and methods. *Machine learning*, 110(3):457–506, 2021.
- Freddie Bickford Smith, Jannik Kossen, Eleanor Trollope, Mark van der Wilk, Adam Foster, and Tom Rainforth. Rethinking aleatoric and epistemic uncertainty, 2024. URL https://arxiv.org/abs/2412.20892.
- Shakir Mohamed and Balaji Lakshminarayanan. Learning in implicit generative models. *arXiv preprint arXiv:1610.03483*, 2016.
- Murat Sensoy, Lance Kaplan, Federico Cerutti, and Maryam Saleki. Uncertainty-aware deep classifiers using generative models. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 5620–5627, 2020.
- Ali Harakeh, Jordan Sir Kwang Hu, Naiqing Guan, Steven Waslander, and Liam Paull. Estimating regression predictive distributions with sample networks. *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(6):7830–7838, Jun. 2023. doi: 10.1609/aaai.v37i6.25948.
- David A Nix and Andreas S Weigend. Estimating the mean and variance of the target probability distribution. In *Proceedings of 1994 ieee international conference on neural networks (ICNN'94)*, volume 1, pages 55–60. IEEE, 1994.

- Gregory P Meyer and Niranjan Thakurdesai. Learning an uncertainty-aware object detector for autonomous driving. In 2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pages 10521–10527. IEEE, 2020.
 - Alexander Immer, Emanuele Palumbo, Alexander Marx, and Julia E Vogt. Effective bayesian heteroscedastic regression with deep neural networks. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
 - Laurens Sluijterman, Eric Cator, and Tom Heskes. Optimal training of mean variance estimation neural networks. *Neurocomputing*, page 127929, 2024.
 - Andrew Stirn and David A. Knowles. Variational variance: Simple, reliable, calibrated heteroscedastic noise variance parameterization, 2020.
 - Peng Cui, Wenbo Hu, and Jun Zhu. Calibrated reliable regression using maximum mean discrepancy. *Advances in Neural Information Processing Systems*, 33:17164–17175, 2020.
 - Moloud Abdar, Farhad Pourpanah, Sadiq Hussain, Dana Rezazadegan, Li Liu, Mohammad Ghavamzadeh, Paul Fieguth, Xiaochun Cao, Abbas Khosravi, U. Rajendra Acharya, Vladimir Makarenkov, and Saeid Nahavandi. A review of uncertainty quantification in deep learning: Techniques, applications and challenges. *Information Fusion*, 76:243–297, 2021. ISSN 1566-2535. doi: https://doi.org/10.1016/j.inffus.2021.05.008.
 - Carl Edward Rasmussen and Christopher K. I. Williams. *Gaussian Processes for Machine Learning*. The MIT Press, 11 2005. ISBN 9780262256834. doi: 10.7551/mitpress/3206.001.0001.
 - Andrew G Wilson and Pavel Izmailov. Bayesian deep learning and a probabilistic perspective of generalization. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, Advances in Neural Information Processing Systems, volume 33, pages 4697–4708. Curran Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper/2020/file/322f62469c5e3c7dc3e58f5a4dlea399-Paper.pdf.
 - Florian Wenzel, Kevin Roth, Bastiaan S. Veeling, Jakub Świątkowski, Linh Tran, Stephan Mandt, Jasper Snoek, Tim Salimans, Rodolphe Jenatton, and Sebastian Nowozin. How good is the bayes posterior in deep neural networks really?, 2020. URL https://arxiv.org/abs/2002.02405.
 - Pavel Izmailov, Sharad Vikram, Matthew D. Hoffman, and Andrew Gordon Wilson. What are bayesian neural network posteriors really like?, 2021.
 - Radford M Neal. Bayesian learning for neural networks, volume 118. Springer Science & Business Media, 1995.
 - Max Welling and Yee W Teh. Bayesian learning via stochastic gradient langevin dynamics. In *Proceedings of the 28th international conference on machine learning (ICML-11)*, pages 681–688, 2011.
 - Chunyuan Li, Changyou Chen, David Carlson, and Lawrence Carin. Preconditioned stochastic gradient langevin dynamics for deep neural networks. In *Proceedings of the AAAI conference on artificial intelligence*, volume 30, 2016.
 - Alex Graves. Practical variational inference for neural networks. In J. Shawe-Taylor, R. Zemel, P. Bartlett, F. Pereira, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 24. Curran Associates, Inc., 2011. URL https://proceedings.neurips.cc/paper/2011/file/7eb3c8be3d411e8ebfab08eba5f49632-Paper.pdf.
 - Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in neural networks, 2015. URL https://arxiv.org/abs/1505.05424.
 - Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In Maria Florina Balcan and Kilian Q. Weinberger, editors, *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pages 1050–1059, New York, New York, USA, 20–22 Jun 2016. PMLR.
 - Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
 - Stanislav Fort, Huiyi Hu, and Balaji Lakshminarayanan. Deep ensembles: A loss landscape perspective, 2020. URL https://arxiv.org/abs/1912.02757.

- Kaizheng Wang, Fabio Cuzzolin, Keivan Shariatmadar, David Moens, and Hans Hallez. Credal wrapper of model averaging for uncertainty estimation in classification. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=cv2iMNWCsh.
 - Matias Valdenegro-Toro and Daniel Saromo Mori. A deeper look into aleatoric and epistemic uncertainty disentanglement. In 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW), pages 1508–1516, 2022. doi: 10.1109/CVPRW56347.2022.00157.
 - Bálint Mucsányi, Michael Kirchhof, and Seong Joon Oh. Benchmarking uncertainty disentanglement: Specialized uncertainties for specialized tasks. In *ICML 2024 Workshop on Structured Probabilistic Inference & Generative Modeling*, 2024. URL https://openreview.net/forum?id=58Lh94RuFQ.
 - Alexander Amini, Wilko Schwarting, Ava Soleimany, and Daniela Rus. Deep evidential regression. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 14927–14937. Curran Associates, Inc., 2020.
 - Jose Miguel Hernandez-Lobato and Ryan Adams. Probabilistic backpropagation for scalable learning of bayesian neural networks. In Francis Bach and David Blei, editors, *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 1861–1869, Lille, France, 07–09 Jul 2015. PMLR.
 - Fredrik K. Gustafsson, Martin Danelljan, and Thomas B. Schön. How reliable is your regression model's uncertainty under real-world distribution shifts?, 2023. URL https://arxiv.org/abs/2302.03679.
 - Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
 - Michele Caprio, Souradeep Dutta, Kuk Jin Jang, Vivian Lin, Radoslav Ivanov, Oleg Sokolsky, and Insup Lee. Credal bayesian deep learning. *arXiv preprint arXiv:2302.09656*, 2023.
 - Kaizheng Wang, Fabio Cuzzolin, Keivan Shariatmadar, David Moens, Hans Hallez, et al. Credal deep ensembles for uncertainty quantification. Advances in Neural Information Processing Systems, 37:79540–79572, 2024.
 - Kevin P Murphy. *Machine learning: a probabilistic perspective*. MIT press, 2021.
 - Sebastian Ruder. An overview of gradient descent optimization algorithms, 2017. URL https://arxiv.org/abs/1609.04747.
 - Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua Bengio and Yann LeCun, editors, 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings, 2015. URL http://arxiv.org/abs/1412.6980.
 - L. V. Kantorovich. Mathematical methods of organizing and planning production. *Management Science*, 6(4): 366–422, 1960. doi: 10.1287/mnsc.6.4.366. URL https://doi.org/10.1287/mnsc.6.4.366.
 - Anonymous. Anonymous submission, 2025.
 - A.R. Melro, P.P. Camanho, and S.T. Pinho. Generation of random distribution of fibres in long-fibre reinforced composites. *Composites Science and Technology*, 68(9):2092–2102, 2008. ISSN 0266-3538. doi: https://doi.org/10.1016/j.compscitech.2008.03.013.
 - Dassault Systèmes. Abaqus 2024, 2024. URL https://www.3ds.com/products/simulia/abaqus. Computer software.
 - Kyunghyun Cho. Learning phrase representations using rnn encoder-decoder for statistical machine translation. arXiv preprint arXiv:1406.1078, 2014.
 - Zhe Gan, Chunyuan Li, Changyou Chen, Yunchen Pu, Qinliang Su, and Lawrence Carin. Scalable bayesian learning of recurrent neural networks for language modeling, 2017.

A MEAN VARIANCE ESTIMATION NETWORKS

A.1 COMPARISON BETWEEN TRAINING A SINGLE NETWORK AND SEPARATE TRAINING OF TWO NETWORKS

As discussed in Section 2.1, we can train a mean variance estimation (MVE) network that predicts both mean and variance together, or we can train two separate networks: a mean estimation network and a variance estimation network. Figure 6 shows the difference when they are trained for the one-dimensional dataset.

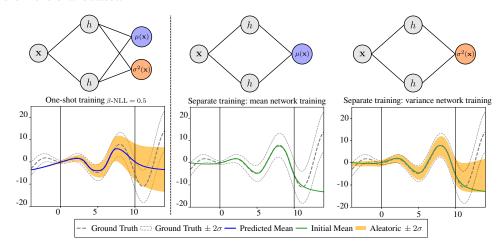


Figure 6: Comparison between MVE network training (left figure) and separate training of ME network and variance estimation network (right figure). The MVE network outputs $\mu(\mathbf{x})$ and $\sigma_a^2(\mathbf{x})$ and is trained by minimizing Equation (2). The ME network is trained first using Equation (2) assuming constant aleatoric uncertainty, and then the VE network is trained using Equation (6) and assuming a fixed mean obtained from the ME network.

MVE achieves worse predictions (left in Figure 6) in the region (x > 7) with higher aleatoric variance prediction. This occurs due to the high-order variance term that appears in the denominator of the gradient of Equation (2), as shown in Equation (3). Conversely, this should not be surprising in light of the vast empirical evidence throughout the literature on the successes of training deterministic models that only predict the mean shown in the middle figure. Interestingly, the variance prediction obtained from separate training aligns well with the ground truth.

A.2 GAUSSIAN LIKELIHOOD LOSS FUNCTION WITH GRADIENT STOPPING OPERATION

 β -NLL loss (Seitzer et al., 2022) was proposed by introducing an additional variance-weighting term, which is given as follows:

$$\mathcal{L}_{\beta\text{-NLL}} = \frac{1}{N} \lfloor \sigma_a^{2\beta}(\mathbf{x}_n) \rfloor \sum_{i=n}^{N} \left[\frac{(y_n - \mu(\mathbf{x}_n))^2}{2\sigma_a^2(\mathbf{x}_n)} + \frac{\log(\sigma_a^2(\mathbf{x}_n))}{2} \right]$$
(12)

where $\lfloor \cdot \rfloor$ represents the stop gradient operation. Under this setup, the gradient of the β -NLL loss becomes:

$$\nabla_{\mu} \mathcal{L}_{\beta\text{-NLL}} = \frac{1}{N} \sum_{n=1}^{N} \left(\frac{\mu(\mathbf{x}_n) - y_n}{\sigma_a^{2-2\beta}(\mathbf{x}_n)} \right)$$
 (13)

$$\nabla_{\sigma_a^2} \mathcal{L}_{\beta\text{-NLL}} = \frac{1}{2N} \sum_{n=1}^N \left(\frac{\sigma_a^2(\mathbf{x}_n) - (y_n - \mu(\mathbf{x}_n))^2}{\sigma_a^{4-2\beta}(\mathbf{x}_n)} \right)$$
(14)

According to Equation (12), the variance-weighting term β can interpolate between the original NLL loss and equivalent MSE, recovering the original NLL loss when $\beta=0$. The gradient of Equation (13) is equivalent to MSE for $\beta=1$.

B BAYESIAN NEURAL NETWORKS

In the Bayesian formalism, the posterior parameter distribution is defined as:

$$p(\boldsymbol{\theta} \mid \mathcal{D}) = \frac{p(\mathcal{D} \mid \boldsymbol{\theta}) p(\boldsymbol{\theta})}{p(\mathcal{D})}, p(\mathcal{D}) = \int p(\mathcal{D} \mid \boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta}$$
(15)

where $p(\theta)$ is the prior, typically a Gaussian distribution, $p(\mathcal{D} \mid \theta)$ is the likelihood, $p(\mathcal{D})$ is the marginal likelihood (or evidence), which integrates the likelihood over θ , and $p(\theta \mid \mathcal{D})$ is the posterior distribution of the parameters θ .

The BNN posterior predictive distribution (PPD) for a new data point is computed as follows:

$$p(y' \mid \mathbf{x}', \mathcal{D}) = \int p(y' \mid \mathbf{x}', \boldsymbol{\theta}) p(\boldsymbol{\theta} \mid \mathcal{D}) d\boldsymbol{\theta}$$
 (16)

where x' denotes the features of the unknown point, and y' the corresponding target.

B.1 MCMC SAMPLING FOR BNNS

Markov Chain Monte Carlo (MCMC) methods remain the gold standard in Bayesian inference. Usually, the prior is enforced on the weights and biases of the neural network. The most common choice is a multivariate Gaussian prior:

$$p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta} \mid \mathbf{0}, \lambda^{-1}\mathbf{I}) \tag{17}$$

where λ is the precision (inverse variance) of the Gaussian prior, which in deterministic networks is referred to as regularization, and \mathbf{I} is the identity matrix.

The likelihood function $p(\mathcal{D} \mid \boldsymbol{\theta})$ quantifies how well the neural network output $\mu(\mathbf{x}; \boldsymbol{\theta})$ explains the observations y. Commonly, the observation distribution is assumed to be Gaussian:

$$p(\mathcal{D} \mid \boldsymbol{\theta}) = \prod_{n=1}^{N} \mathcal{N}(\mathbf{y}_n \mid \mu(\mathbf{x}_n; \boldsymbol{\theta}), \sigma_a^2(\mathbf{x}_n; \boldsymbol{\phi}))$$
(18)

where $\mu(\mathbf{x}; \boldsymbol{\theta})$ is the predictive mean of the neural network and $\sigma_a^2(\mathbf{x}; \boldsymbol{\phi})$ is the data noise variance (assumed constant when the noise is homoscedastic, or learned by another neural network parameterized by $\boldsymbol{\phi}$ as described in Section 3.3.1).

As shown in Equations (9) and (15), it is not possible to get the analytical solution because it requires integrating the posterior and marginal likelihood. However, we can obtain the following relation by omitting the denominator:

$$p(\boldsymbol{\theta} \mid \mathcal{D}) \propto p(\mathbf{y} \mid \boldsymbol{\theta}, \mathbf{x}) p(\boldsymbol{\theta})$$
 (19)

Substituting the prior (Equation (17)) and likelihood (Equation (18)), we obtain:

$$p(\boldsymbol{\theta} \mid \mathcal{D}) \propto \prod_{n=1}^{N} \mathcal{N}\left(\mathbf{y}_{n} \mid \mu(\mathbf{x}_{n}; \boldsymbol{\theta}), \sigma_{a}^{2}(\mathbf{x}_{n}; \boldsymbol{\phi})\right) \cdot \mathcal{N}(\boldsymbol{\theta} \mid \mathbf{0}, \lambda^{-1}\mathbf{I})$$
(20)

$$= \prod_{n=1}^{N} \frac{1}{\sqrt{2\pi\sigma_a^2(\mathbf{x}_n; \boldsymbol{\phi})}} \exp\left(-\frac{(\mathbf{y}_n - \mu(\mathbf{x}_n; \boldsymbol{\theta}))^2}{2\sigma_a^2(\mathbf{x}_n; \boldsymbol{\phi})}\right) \cdot \frac{\lambda^{m/2}}{(2\pi)^{m/2}} \exp\left(-\frac{\lambda}{2} \|\boldsymbol{\theta}\|^2\right)$$
(21)

where m is the number of parameters within θ . By taking the logarithmic form, we obtain

$$\log p\left(\boldsymbol{\theta} \mid \mathcal{D}\right) = \sum_{n=1}^{N} \left[\log \frac{1}{\sqrt{2\pi\sigma_a^2(\mathbf{x}_n; \boldsymbol{\phi})}} - \frac{1}{2} \frac{(\mathbf{y}_n - \mu(\mathbf{x}_n; \boldsymbol{\theta}))^2}{\sigma_a^2(\mathbf{x}_n; \boldsymbol{\phi})} \right] + m \log \frac{1}{\sqrt{2\pi/\lambda}} - \frac{\lambda}{2} \|\boldsymbol{\theta}\|^2$$
(22)

where the first two terms relate to the likelihood and the final two terms only relate to the neural network parameters. In the deterministic setting, the last two often act as regularization or weight decay.

MCMC is a class of algorithms that can be used to sample from a probability distribution (Murphy, 2021). The most straightforward approach to get the posterior distribution is the random walk Metropolis-Hasting algorithm (Neal, 1995), although it suffers from high rejection rates for high-dimensional problems. To accelerate the mixing rate of the MCMC methods, the Hamiltonian Monte Carlo (HMC) was developed by Neal et al. (Neal, 1995), which leverages the concept of Hamiltonian mechanics, where the *gradient* of the target probability density function is properly utilized. Therefore, the mixing rate can be improved tremendously compared with the random walk Metropolis-Hasting algorithm. To address the scalability bottleneck of HMC, Welling et al. (Welling and Teh, 2011) developed a novel Bayesian inference approach called Stochastic Gradient Langevin Dynamics (SGLD) leveraging Stochastic Gradient Decent (SGD) (Ruder, 2017) and Langevin Monte Carlo (LMC) (Murphy, 2021).

In essence, SGLD starts with a random guess for the unknown neural network parameters $\theta^{(0)}$ in Equation (22) and then updates the position according to the following rule (Welling and Teh, 2011):

$$\Delta \boldsymbol{\theta}_{t} = \frac{\eta_{t}}{2} \left(\nabla \log p\left(\boldsymbol{\theta}_{t}\right) + \nabla \frac{N}{M} \sum_{n=1}^{M} \log p(\mathcal{D}_{n} | \boldsymbol{\theta}_{t}) \right) + \boldsymbol{\zeta}_{t}$$
(23)

where $\zeta_t \sim \mathcal{N}(\mathbf{0}, \eta_t)$, η is the step size (learning rate), N is the number of the total data points, M is the batch size.

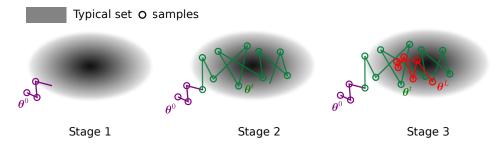


Figure 7: **Schematic of MCMC-based Bayesian inference.** The shaded area is the typical set, which is the region that covers the high probabilities of the posterior distribution. The samplers start with a random initialization and try to integrate the typical set. It has three stages: the first stage is to search for the typical set; the second stage is the most effective phase, exploring the typical set rapidly; and in the third stage, the sampler might converge to a single mode.

However, SGLD assumes that all parameters θ have the same step size, which can lead to slow convergence or even divergence in cases where the components of θ have different curvatures. A refined version of SGLD called preconditioned SGLD (pSGLD) (Li et al., 2016), was proposed to address this issue. In pSGLD (Li et al., 2016), the update rule incorporates a user-defined preconditioning matrix $G(\theta_t)$, which adjusts the gradient updates and the noise term adaptively:

$$\Delta \theta_t = \frac{\eta_t}{2} \left[G(\boldsymbol{\theta}_t) \left(\nabla \log p(\boldsymbol{\theta}_t) + \nabla \frac{N}{M} \sum_{n=1}^M p(\mathcal{D}_n | \boldsymbol{\theta}_t) \right) + \chi(\boldsymbol{\theta}_t) \right] + \zeta_t G(\boldsymbol{\theta}_t)$$
 (24)

where $\chi_n = \sum_j \frac{\partial G_{n,j}(\boldsymbol{\theta})}{\partial \theta_j}$; $j = 0, \dots, d$ describes how the preconditioner changes with respect to $\boldsymbol{\theta}$.

B.2 VARIATIONAL INFERENCE FOR BNNs

VI is another type of Bayesian inference method that approximates the posterior distribution by a proposed distribution from a parametric family (Blundell et al., 2015; Murphy, 2021). The goal of VI is to minimize the Kullback-Leibler divergence between the true posterior distribution $p(\theta|\mathcal{D})$ and the proposed distribution $q(\theta)$ as illustrated in Figure 8.

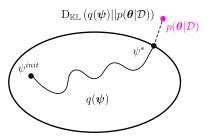


Figure 8: Schematic of Variational Inference (Murphy, 2021). The ellipse represents the search space of the proposed distribution, $p(\theta|\mathcal{D})$ is the true posterior distribution, the KL divergence is the distance between the two distributions defined as $D_{\mathbb{KL}}\left(q(\psi)||p(\theta|\mathcal{D})\right)$, and the goal is to minimize it

Practically, ψ is regarded as the variational parameter from a parametric family Ω ; therefore, the optimal ψ^* can be obtained by minimizing the KL divergence as follows (Murphy, 2021):

$$\psi^* = \arg\min_{\boldsymbol{\psi}} \mathbb{D}_{\mathbb{KL}} \left(q(\boldsymbol{\psi}) || p(\boldsymbol{\theta} | \mathcal{D}) \right)$$

$$= \arg\min_{\boldsymbol{\psi}} \mathbb{E}_{q(\boldsymbol{\psi})} \left[\log q(\boldsymbol{\psi}) - \log p(\boldsymbol{\theta} | \mathcal{D}) \right]$$

$$= \arg\min_{\boldsymbol{\psi}} \mathbb{E}_{q(\boldsymbol{\psi})} \left[\log q(\boldsymbol{\psi}) - \log \left(\frac{p(\mathcal{D} | \boldsymbol{\theta}) p(\boldsymbol{\theta})}{p(\mathcal{D})} \right) \right]$$

$$= \arg\min_{\boldsymbol{\psi}} \mathbb{E}_{q(\boldsymbol{\psi})} \left[\log q(\boldsymbol{\psi}) - \log p(\mathcal{D} | \boldsymbol{\theta}) - \log p(\boldsymbol{\theta}) + \log p(\mathcal{D}) \right]$$

$$= \arg\min_{\boldsymbol{\psi}} \mathbb{E}_{q(\boldsymbol{\psi})} \left[\log q(\boldsymbol{\psi}) - \log p(\mathcal{D} | \boldsymbol{\theta}) - \log p(\boldsymbol{\theta}) \right] + \log p(\mathcal{D})$$

$$(25)$$

According to Equation (25), the optimization problem can be decomposed into two terms: the first term is the negative evidence lower bound (ELBO), and the second term is the log evidence. The log evidence is a constant term that does not depend on ψ ; therefore, we only need to optimize the first term:

$$\psi^* = \arg\min_{\psi} \mathbb{E}_{q(\psi)} \left[\log q(\psi) - \log p(\mathcal{D}|\boldsymbol{\theta}) - \log p(\boldsymbol{\theta}) \right]$$

$$= \arg\min_{\psi} \mathbb{E}_{q(\psi)} \left[-\log p(\mathcal{D}|\boldsymbol{\theta}) \right] + \mathcal{D}_{\mathbb{KL}}(q(\psi)||p(\boldsymbol{\theta}))$$
(26)

Similarly as Equation (22), the first term is the negative log-likelihood of the given dataset, and the second term is the KL divergence between the proposed distribution and the prior distribution, which again can be regarded as a regularizer. The optimization problem can be solved by any optimization algorithm, such as SGD (Ruder, 2017) or Adam (Kingma and Ba, 2015). After obtaining the best ψ , the variational distribution $q(\psi)$ can be used as the posterior distribution.

C PERFORMANCE METRICS

We use the Root Mean Square Error (RMSE) and Test Log-Likelihood (TLL) for synthetic, UCI regression, and plasticity law datasets. In addition, as we know the ground truth aleatoric uncertainty for the last dataset (material plasticity), we also use the Wasserstein distance (Kantorovich, 1960) to evaluate the correctness of the learned aleatoric uncertainty for that problem. Details of those metrics are listed:

• Root Mean Square Error (RMSE)

- For MLP:

$$RMSE = \sqrt{\frac{1}{N_{test}} \sum_{i=1}^{N_{test}} (\bar{\boldsymbol{\mu}}_i - \mathbf{y}_i)^T (\bar{\boldsymbol{\mu}}_i - \mathbf{y}_i)}$$
 (27)

where N_{test} is the number of test data points, $\bar{\mu}$ is the predictive mean, and y is the observation values of the test points.

- For RNN:

$$RMSE = \frac{1}{N_{test}T} \sum_{i=1}^{N_{test}} \sum_{j=1}^{T} \left(\frac{\|\bar{\boldsymbol{\mu}}_{i,j} - \bar{\mathbf{y}}_{i,j}\|_{F}}{\|\bar{\mathbf{y}}_{i,j}\|_{F}} \right) \cdot 100\%, \tag{28}$$

where N_{test} is the number of testing points, $\|\cdot\|_F$ is a Frobenius norm, $\bar{\mathbf{y}}_{i,j}$ and $\bar{\mu}_{i,j}$ are ground truth and the predictive mean of i^{th} test sample and j^{th} time step, correspondingly.

• Test Log-Likelihood (TLL)

$$TLL = \frac{1}{N_{test}T} \sum_{i=1}^{N_{test}} \sum_{j=1}^{T} \log p(\bar{\mathbf{y}}_{i,j}|\boldsymbol{\theta})$$
 (29)

We also clarify that we use the ground truth $\bar{\mathbf{y}}$ for the plasticity law dataset. Observation values \mathbf{y} are used for synthetic and UCI regression datasets and T=1.

• Wasserstein distance (WA)

$$WA = \frac{1}{N_{\text{test}}T} \sum_{i=1}^{N_{\text{test}}} \sum_{j=1}^{T} W_2 \left(p\left(\bar{\boldsymbol{\mu}}_{i,j}, \boldsymbol{\sigma}_{a_{i,j}}^2 \mathbf{I} \mid \boldsymbol{\theta}, \boldsymbol{\phi} \right), q\left(\bar{\mathbf{y}}_{i,j}\right) \right)$$
(30)

where $p\left(\bar{\boldsymbol{\mu}}_{i,j}, \boldsymbol{\sigma}_{a_{i,j}}^2 \mathbf{I} \mid \boldsymbol{\theta}, \boldsymbol{\phi}\right)$ and $q\left(\bar{\mathbf{y}}_{i,j}\right)$ are the predictive and ground truth distributions, respectively. Since Gaussian assumption is applied, we can rewrite Equation (30) into:

$$WA = \frac{1}{N_{\text{test}}T} \sum_{i=1}^{N_{\text{test}}} \sum_{j=1}^{T} \sqrt{(\bar{\mathbf{y}}_{i,j} - \bar{\boldsymbol{\mu}}_{i,j})^{T} (\bar{\mathbf{y}}_{i,j} - \bar{\boldsymbol{\mu}}_{i,j}) + (\boldsymbol{\sigma}_{i,j}^{2} - \boldsymbol{\sigma}_{ai,j}^{2})^{T} (\boldsymbol{\sigma}_{i,j}^{2} - \boldsymbol{\sigma}_{ai,j}^{2})}$$
(31)

Regarding the Image regression datasets, as introduced by respective authors, the aim is to evaluate the reliability of regression uncertainty estimation under distribution shift. In this case, the prediction is first calibrated with a validation dataset (in distribution). Then the test coverage for the test dataset, which has a distribution shift, is evaluated (Gustafsson et al., 2023).

• Mean absolute error (MAE)

$$MAE = \frac{1}{N} \sum_{i=1}^{N} |\bar{\boldsymbol{\mu}}_i - \mathbf{y}_i|$$
 (32)

where N is the number of points for validation or testing.

· val Interval length

$$\hat{C}_{\alpha}(\mathbf{x}, \mathcal{D}_{val}) = [L_{\alpha}(\mathbf{x}, \mathcal{D}_{val}) - Q_{1-\alpha}(E, \mathcal{D}_{val}), \ U_{\alpha}(\mathbf{x}, \mathcal{D}_{val}) + Q_{1-\alpha}(E, \mathcal{D}_{val})], \ \ (33)$$
 where L_{α} and U_{α} are the lower and upper bounds of confidence interval under $\alpha = 0.1$. $E = \{\max(L_{\alpha}(x_i) - y_i, \ y_i - U_{\alpha}(x_i)) : i \in \mathcal{D}_{val}\}$ are conformity scores of the validation dataset, and $Q_{1-\alpha}(E, \mathcal{D}_{val})$ is the $(1-\alpha)$ -th quantile.

Test coverage

Following the same procedure, we can get the predictive confidence interval for the test dataset calibrated by the validation dataset.

$$\hat{C}_{\alpha}(\mathbf{x}, \mathcal{D}_{test}) = [L_{\alpha}(\mathbf{x}, \mathcal{D}_{test}) - Q_{1-\alpha}(E, \mathcal{D}_{val}), \ U_{\alpha}(\mathbf{x}, \mathcal{D}_{test}) + Q_{1-\alpha}(E, \mathcal{D}_{val})], \ (34)$$

With the predictive confidence interval, we can obtain the test coverage with the following formula:

Test coverage =
$$\frac{1}{N_{test}} \sum_{i=1}^{N_{test}} \mathbb{I}\{y_i' \in \hat{C}_{\alpha}(\mathbf{x}, \mathcal{D}_{test})\}. \tag{35}$$

D ADDITIONAL EXPERIMENTS RESULTS

D.1 ILLUSTRATIVE EXAMPLE: ONE-DIMENSIONAL DATASET

Comparison of different inference method The predictions for the one-dimensional example by our method with pSGLD inference – VeBNN (pSGLD) – are shown in Figure 1. A direct comparison with MVE (β -NLL) assuming the best value for the β hyperparameter that we found, $\beta = 0.5$, is shown in the Appendix in Figure 6. It is clear that our strategy of separately training the mean and aleatoric variance leads to better results, and avoids the need for an extra hyperparameter (β).

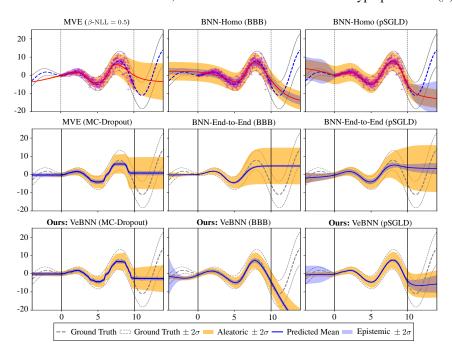


Figure 9: Heteroscedastic regression by our method (bottom) compared to existing BNN training with homoscedastic assumption (top) and End-to-End training methods (middle) for each inference type.

Figure 9 also shows the same example when compared to other Bayesian methods with homoscedastic and heteroscedastic assumptions that are capable of predicting both uncertainties. We see a consistent improvement in the predictions with the cooperative training strategy we proposed, independently of the inference method that is chosen. BNN training with homoscedastic assumption leads to inferior results on the aleatoric uncertainty estimation, since a constant variance σ_a^2 is not adequate. Furthermore, the End-to-End training strategy starts to overestimate aleatoric uncertainty approximately after x>8, regardless of the Bayesian inference method. This effect becomes more pronounced in the extrapolation region (x>10), which reflects the difficulties that arise from the joint training process. There is a clear tendency to overestimate the aleatoric uncertainty by the End-to-End strategy that also affects the mean estimation. In contrast, the proposed strategy improves predictions according to all metrics. As expected, Bayesian inference with pSGLD is the best.

Heteroscedastic noise We show the regression results for increasing the number of training data points from 5 to 500. The accuracy metrics obtained by running each case 5 times for randomly sampled training sets are shown in Figure 10, and we also visualize the fitting performance of selected methods under an arbitrary run in Figure 11.

The best value for the hyperparameter β was 0.5. The experimental results reveal that MVE with $\beta=0.5$ converges to the same mean but slower than the proposed VeBNN (pSGLD). In contrast, BNN-End-to-End (pSGLD) has difficulty in converging to either mean or aleatoric uncertainty. The proposed strategy successfully combines the strengths of MVE and BNN to effectively model aleatoric uncertainty, mean predictions, and epistemic uncertainty. The results highlight the effectively

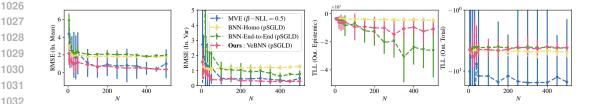


Figure 10: RMSE and TLL convergence curves under 5 different seeds for heteroscedastic data generation. The first figure shows the RMSE between the data values and predictive mean within the interpolation region; the second figure shows the RMSE between the noise standard deviation and the predictive one in the interpolation region since the ground truth of aleatoric uncertainty is known; and the third and fourth figures show the Epistemic TLL and Total TLL values for extrapolation test points.

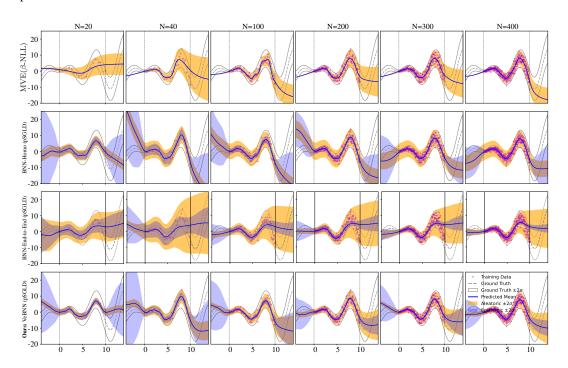


Figure 11: Predictions of MVE (β -NLL) with $\beta = 0.5$, BNN-Homo (pSGLD), BNN-End-to-End (pSGLD), and Ours: VeBNN (pSGLD) with different data points under heteroscedastic data **noise.** In this plot, the first to fourth rows represent the prediction of MVE (β -NLL = 0.5), BNN-Homo (pSGLD), BNN-End-to-End (pSGLD), and **Ours:** VeBNN (pSGLD) with different training data points under the same seed for data generation.

tiveness of cooperative training instead of End-to-End training, especially in the case of a data-scarce scenario. Notably, the proposed method maintains a stable TLL value in the extrapolation region, outperforming alternative approaches in this critical area.

Homoscedastic noise We execute a similar experiment to test homoscedastic noise cases where the data size changes from 5 to 200. Figure 12 and Figure 13 highlight similar conclusions as Appendix D.1. Again, the MVE training experiences great difficulty when the training data is scarce. The proposed cooperative training strategy still has the best performance even in the case where 5 points are used for training. In contrast, BNN-End-to-End (pSGLD) gives really large uncertainties in the cases of N < 30. It should also be noticed that the extrapolation behaviors for both VeBNN (pSGLD) and BNN-End-to-End (pSGLD) are similar, as shown in Figure 13, while BNN-End-to-End gets slightly larger values for TLLs in Figure 12.

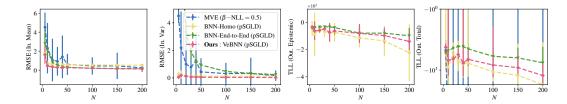


Figure 12: RMSE and TLL convergence curves under 5 different seeds for homoscedastic data generation.

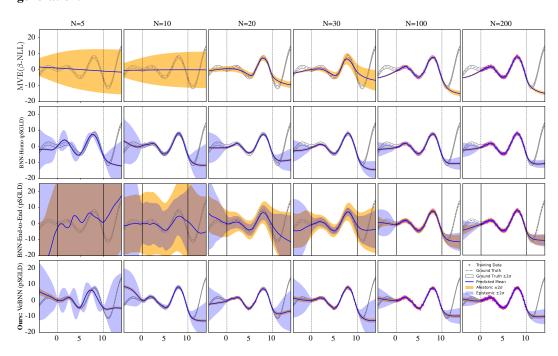


Figure 13: Predictions of MVE (β -NLL = 0.5), BNN-Homo (pSGLD) and Ours: VeBNN (pS-GLD) with different data points under homoscedastic data noise.

D.2 RMSE RESULTS FOR UCI REGRESSION DATASETS

In Section 4.1, we only presented the TLL performance where **Ours:** VeBNN (pSGLD) has 5 datasets with the best and 1 dataset with the second best among all 8 datasets. Moreover, we can see more clearly that the proposed training strategy has improved within the same Bayesian inference method. In this section, we summarize the results concerning RMSE in Table 2.

Overall, VeBNN (pSGLD) performs well for all datasets, having the best performance in 4 datasets and being the second best for 2 datasets. It should be mentioned that the deterministic network, ME (MSE), is competitive because Table 2 only evaluates the error in estimating the mean. We raise awareness that the proposed cooperative training strategy consistently outperforms end-to-end training for all Bayesian inference methods that we consider (MC-Dropout, BBB and pSGLD). This reinforces that the proposed strategy is clearly beneficial, independent of the Bayesian inference method of choice. Nevertheless, as we discuss in the main text, we recommend the use of pSGLD for the types of distributions considered herein.

D.3 COMPLETE RESULTS FOR IMAGE REGRESSION DATASETS

We report results only for the problems involving distribution shift in Figure 14, while the complete results across all problems are summarized in Table 3. As stated in (Gustafsson et al., 2023), the purpose of this dataset is to evaluate the uncertainty quantification capabilities of deep learning

Table 2: RMSE (↓) results on UCI Regression Datasets. The best performance for each dataset is underlined and bold, and the second best is in bold.

Methods	Carbon (10721,7,1)	Concrete (1030,8,1)	Energy (768,8,2)	Boston (506,13,1)
ME (MSE)	0.0069 (0.0028)	4.59 (0.86)	0.74 (0.07)	3.56 (0.81)
MVE (β_{NLL} =1.0)	0.0070 (0.0029)	5.38 (0.84)	0.78 (0.08)	3.60 (0.91)
MVE (Ensemble)	0.0066 (0.0029)	5.11 (0.67)	0.79(0.12)	4.79 (0.94)
Evidential	0.0068 (0.0029)	5.96 (0.61)	2.27 (0.46)	4.01 (1.02)
MVE (MC-Dropout)	0.0159 (0.0050)	5.43 (0.70)	1.66 (0.41)	3.85 (1.01)
Ours: VeBNN (MC-Dropout)	0.0105 (0.0022)	5.05 (0.81)	1.29 (0.12)	3.08 (0.73)
BNN-Homo (BBB)	0.0070 (0.0035)	5.69 (0.74)	1.35 (0.13)	3.81 (0.95)
BNN-End-to-End (BBB)	0.1307 (0.0384)	58.13 (47.72)	54.53 (40.59)	573.39 (238.16)
Ours: VeBNN (BBB)	0.0069 (0.0027)	5.43 (0.69)	1.20 (0.15)	3.56 (0.83)
BNN-Homo (pSGLD)	0.0068 (0.0028)	4.91 (0.68)	1.04 (0.10)	3.56 (0.87)
BNN-End-to-End (pSGLD)	0.0066 (0.0029)	5.58 (0.62)	2.08 (0.37)	3.81 (0.93)
Ours: VeBNN (pSGLD)	0.0065 (0.0030)	4.65 (0.97)	0.70 (0.08)	3.57 (0.95)
	Power	Superconduct	Wine-Red	Yacht
Methods	(9568,4,1)	(21263,81,1)	(1599,11,1)	(308,6,1)
ME (MSE)	3.86 (0.15)	11.70 (0.55)	0.65 (0.05)	0.67 (0.24)
MVE (β_{NLL} =1.0)	3.90 (0.13)	12.52 (0.77)	0.63 (0.05)	$\overline{0.83 (0.38)}$
MVE (Ensemble)	3.86 (0.14)	11.82 (0.34)	0.79(0.09)	0.96 (0.35)
Evidential	3.92 (0.14)	13.93 (0.44)	0.64 (0.06)	3.26 (2.21)
MVE (MC-Dropout)	4.06 (0.12)	13.08 (0.92)	0.64 (0.06)	0.94 (0.31)
Ours: VeBNN (MC-Dropout)	4.00 (0.13)	12.28 (0.50)	0.63 (0.05)	0.78 (0.28)
BNN-Homo (BBB)	4.08 (0.13)	13.27 (0.39)	0.64 (0.05)	2.20 (0.53)
BNN-End-to-End (BBB)	8.82 (3.27)	404.30 (414.21)	2.11 (1.56)	261.08 (148.66)
Ours: VeBNN (BBB)	3.99 (0.13)	13.65 (0.46)	0.63 (0.05)	1.09 (0.26)
BNN-Homo (pSGLD)	3.79 (0.14)	11.20 (0.39)	0.62 (0.05)	1.54 (0.46)
BNN-End-to-End (pSGLD)	3.83 (0.15)	13.65 (0.31)	$\overline{0.64\ (0.06)}$	1.56 (0.76)
Ours: VeBNN (pSGLD)	3.77 (0.14)	12.04 (0.43)	0.62 (0.05)	0.70 (0.29)

models under distribution shift. MVE (Ensembles) consistently serves as the strongest baseline across all problems, aligning with findings from (Gustafsson et al., 2023). As shown in Table 3, *Cells* and *ChairAngle* are the two baseline tasks without any distribution shift between training and test sets. In these cases, the test coverage reaches the target of 90%, as expected. However, our proposed VeBNN (pSGLD) performs better because it achieves smaller MAE and narrower interval lengths on the validation set.

For the shifted tasks *Cells-Gap* and *Cells-Tails*, VeBNN (pSGLD) attains test coverage closer to the 90% target compared to other methods, while also achieving lower validation MAE and shorter interval lengths. For the remaining tasks, VeBNN (pSGLD) performs comparably to MVE (Ensembles). Specifically, in *ChairAngle-Gap*, VeBNN (pSGLD) yields slightly higher test coverage but also slightly larger validation MAE and interval length. In *ChairAngle-Tail*, VeBNN (pSGLD) achieves significantly better test coverage with only a minor increase in validation MAE and interval length. For the last two problems, VeBNN (pSGLD) maintains comparable test coverage while producing tighter prediction intervals.

It is also worth noting that MC-Dropout fails to maintain adequate coverage across these image regression datasets. In contrast, VeBNN (MC-Dropout) generally outperforms MVE (MC-Dropout), with the exception of the *SkinLesion* task on the validation set.

Although these datasets cannot be used to evaluate the quality of aleatoric uncertainty estimation, we noticed that the proposed VeBNN (pSGLD) method tends to have higher epistemic uncertainty in extrapolation regions (datasets ending with "-TAIL"), while the predicted aleatoric uncertainty remains stable even in these regions. This provides some indication that aleatoric uncertainty might be disentangled more effectively than with existing methods – see Figure 15.

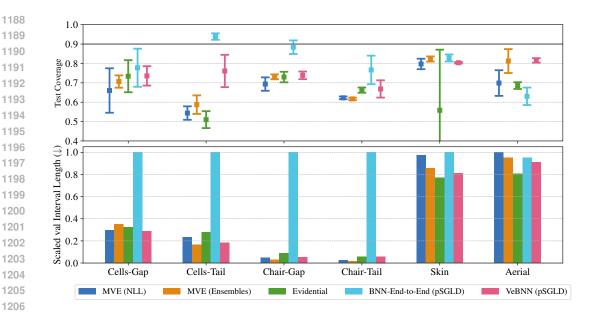


Figure 14: Accuracy metrics of select methods for six problems with the distribution shift test dataset. The upper figure illustrates the test coverage, where values approaching 0.9 indicate better calibration. The lower figure presents the scaled interval length, normalized by the maximum value across all methods; smaller values signify tighter and more precise prediction intervals.

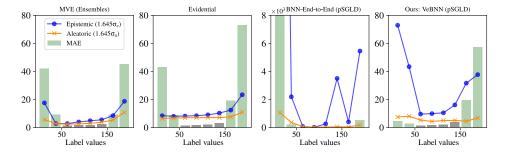


Figure 15: Uncertainty predictions versus Absolute Error for the *Cells-Tail* problem. The central four gray bars represent the Absolute Error of test points from the in-distribution data. The two bars on each side correspond to the Absolute Error of test points under distribution shift relative to the training dataset.

D.4 PLASTICITY LAWS DATASETS

Additional plots of predictions for correct identification Given the limited space in the main text, we present the predictions of MVE (β -NLL), MVE (MC-Dropout), and VeBNN (MC-Dropout) on the plasticity law discovery dataset in Figure 16. We first observe that MVE (β -NLL) provides reliable aleatoric uncertainty estimates when sufficient training data is available. However, its predictions—both in terms of the mean and aleatoric uncertainty—deteriorate when the training size is reduced to N=50. For MVE (MC-Dropout), the predictions remain suboptimal, even with N=800 training sequences. In contrast, the proposed VeBNN (MC-Dropout) improves prediction quality, particularly in estimating aleatoric uncertainty.

It is important to note that predictions for plasticity laws are expected to be smooth. The inherent bumpiness of MC-Dropout-based approaches poses challenges for deployment in Finite Element Analysis in real-world applications.

Table 3: Complete results for the image regression datasets. The first two columns are the MAE and interval length for the validation dataset, and the third and fourth columns are the MAE and test coverage results of the test dataset. For each problem, we highlight the best method in green but we also highlight in yellow methods with similar performance. Note that all metrics must be considered in the method evaluation: the objective is to have test coverage as close to 0.90 as possible (for a low validation interval length), and to have low MAE (for validation and testing sets). Test coverage is the most important metric, but only when MAE is low.

Problem	Method	val MAE (↓)	val Interval Length (↓)	test MAE (↓)	test Coverage (0.90)
	MVE (NLL)	3.6170 (1.1362)	14.5492 (4.4493)	3.3858 (1.1927)	0.9028 (0.0059)
	MVE (Ensembles)	2.7876 (1.1695)	17.1285 (4.3337)	2.8788 (0.4288)	0.9006 (0.0027)
	Evidential	2.1816 (0.5478)	12.6658 (2.1219)	12.6657 (0.5186)	0.8865 (0.0026)
Cells	MVE (MC-Dropout)	15.8976 (1.0837)	93.5508 (4.2047)	50.6258 (0.2451)	0.8979 (0.0165)
	Ours: VeBNN (MC-Dropout)	10.5583 (1.2076)	47.8594 (5.1651)	10.8850 (1.5156)	0.8948 (0.0148)
	BNN-End-to-End (pSGLD)	9.4833 (3.4854)	69.2175 (34.6855)	10.8901 (4.2301)	0.9019 (0.0035)
	Ours: VeBNN (pSGLD)	2.1714 (0.1681)	11.4478 (2.9070)	2.2307 (0.1593)	0.8959 (0.0056)
	MVE (NLL)	3,5309 (1,0619)	15.6396 (6.2345)	6.5164 (1.2826)	0.6603 (0.1147)
	` /	` ′	, ,	, ,	0.7066 (0.0318)
					0.7341 (0.0828)
Cells-Gap					0.7372 (0.0764)
•	Ours: VeBNN (MC-Dropout)	11.3806 (4.5054)	56.0703 (16.8834)	17.3838 (1.0463)	0.7181 (0.1721)
	BNN-End-to-End (pSGLD)	5.9391 (1.8787)	53.0860 (34.1252)	8.0690 (3.6231)	0.7779 (0.0984)
	Ours: VeBNN (pSGLD)	2.4714 (0.5216)	15.2334 (5.8290)	5.5104 (1.9681)	0.7358 (0.0501)
	MVE (NLL)	4.0545 (1.3315)	15.4321 (4.9880)	14.6865 (2.2122)	0.5440 (0.0348)
~					0.5874 (0.0479)
	Evidential	2.0857 (0.1780)	18.4345 (0.8247)	18.4345 (3.4760)	0.5100 (0.0436)
Cells-Tail			, ,	, ,	0.6942 (0.0217)
					0.5974 (0.0380)
Cells-Tail WVE (MC-l Ours: VeBN BNN-End-to Ours: VeBN WVE (NLL) WVE (Enser Evidential Chair Chair MVE (MC-l Ours: VeBN BNN-End-to Ours: VeBN WVE (MC-l Ours: VeBN BNN-End-to Ours: VeBN WVE (NLL) MVE (Enser		8.7096 (4.2490)	66.0935 (35.1367)	407.6956 (601.3301)	0.9387 (0.0166)
	Ours: VeBNN (pSGLD)	2.4510 (0.7056)	12.2763 (3.5848)	41.4620 (37.2828)	0.7611 (0.0834)
	MVE (NLL)	0.3767 (0.1719)	1.3776 (0.38225)	0.4805 (0.0274)	0.9016 (0.0031)
	MVE (Ensembles)	0.3618 (0.1653)	1.2044 (0.4424)	0.4051 (0.0799)	0.9104 (0.0032)
	Evidential	0.3413 (0.1762)	2.6932 (0.3732)	0.3350 (0.1756)	0.9066 (0.0032)
Chair	MVE (MC-Dropout)	9.7782 (0.6145)	54.9429 (4.3036)	12.5577 (0.2768)	0.6379 (0.0386)
	Ours: VeBNN (MC-Dropout)	10.1604 (1.1318)	44.9722 (4.0769)	21.5712 (0.7730)	0.5752 (4.5598)
	BNN-End-to-End (pSGLD)	11.0020 (4.2520)	84.5294 (29.1052)	12.0116 (4.9162)	0.9072 (0.0052)
	Ourse VoDNN (pCCLD)	0.2010 (0.0022)	1 1524 (0 2000)	0.0005 (0.0000)	0.0060 (0.0046)
	Ours: vebinin (psoled)	0.2918 (0.0823)	1.1524 (0.2098)	0.2895 (0.0838)	0.9060 (0.0046)
	MVE (NLL)	0.4545 (0.2802)	2.0921 (0.9338)	1.4182 (0.2737)	0.6936 (0.0346)
	MVE (NLL) MVE(Ensembles)	0.4545 (0.2802) 0.2264 (0.0677)	2.0921 (0.9338) 1.3059 (0.1362)	1.4182 (0.2737) 1.1909 (0.0607)	0.6936 (0.0346) 0.7310 (0.0126)
	MVE (NLL) MVE(Ensembles) Evidential	0.4545 (0.2802) 0.2264 (0.0677) 0.4683 (0.1803)	2.0921 (0.9338) 1.3059 (0.1362) 3.8559 (0.5422)	1.4182 (0.2737) 1.1909 (0.0607) 1.7800 (0.2434)	0.6936 (0.0346) 0.7310 (0.0126) 0.7299 (0.0272)
Chair-Gap	MVE (NLL) MVE(Ensembles) Evidential MVE (MC-Dropout)	0.4545 (0.2802) 0.2264 (0.0677) 0.4683 (0.1803) 15.5046 (2.3557)	2.0921 (0.9338) 1.3059 (0.1362) 3.8559 (0.5422) 99.6555 (14.3327)	1.4182 (0.2737) 1.1909 (0.0607) 1.7800 (0.2434) 15.4490 (1.7975)	0.6936 (0.0346) 0.7310 (0.0126) 0.7299 (0.0272) 0.7372 (0.0764)
Chair-Gap	MVE (NLL) MVE(Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout)	0.4545 (0.2802) 0.2264 (0.0677) 0.4683 (0.1803) 15.5046 (2.3557) 14.2239 (1.8959)	2.0921 (0.9338) 1.3059 (0.1362) 3.8559 (0.5422) 99.6555 (14.3327) 63.0351 (6.7923)	1.4182 (0.2737) 1.1909 (0.0607) 1.7800 (0.2434) 15.4490 (1.7975) 22.1338 (2.2755)	0.6936 (0.0346) 0.7310 (0.0126) 0.7299 (0.0272) 0.7372 (0.0764) 0.7331 (0.0463)
Chair-Gap	MVE (NLL) MVE(Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD)	0.4545 (0.2802) 0.2264 (0.0677) 0.4683 (0.1803) 15.5046 (2.3557) 14.2239 (1.8959) 7.0322 (0.9985)	2.0921 (0.9338) 1.3059 (0.1362) 3.8559 (0.5422) 99.6555 (14.3327) 63.0351 (6.7923) 44.2601 (1.2242)	1.4182 (0.2737) 1.1909 (0.0607) 1.7800 (0.2434) 15.4490 (1.7975) 22.1338 (2.2755) 7.4614 (1.5881)	0.6936 (0.0346) 0.7310 (0.0126) 0.7299 (0.0272) 0.7372 (0.0764) 0.7331 (0.0463) 0.8838 (0.0351)
Chair-Gap	MVE (NLL) MVE(Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout)	0.4545 (0.2802) 0.2264 (0.0677) 0.4683 (0.1803) 15.5046 (2.3557) 14.2239 (1.8959)	2.0921 (0.9338) 1.3059 (0.1362) 3.8559 (0.5422) 99.6555 (14.3327) 63.0351 (6.7923)	1.4182 (0.2737) 1.1909 (0.0607) 1.7800 (0.2434) 15.4490 (1.7975) 22.1338 (2.2755)	0.6936 (0.0346) 0.7310 (0.0126) 0.7299 (0.0272) 0.7372 (0.0764) 0.7331 (0.0463)
Chair-Gap	MVE (NLL) MVE(Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD)	0.4545 (0.2802) 0.2264 (0.0677) 0.4683 (0.1803) 15.5046 (2.3557) 14.2239 (1.8959) 7.0322 (0.9985)	2.0921 (0.9338) 1.3059 (0.1362) 3.8559 (0.5422) 99.6555 (14.3327) 63.0351 (6.7923) 44.2601 (1.2242)	1.4182 (0.2737) 1.1909 (0.0607) 1.7800 (0.2434) 15.4490 (1.7975) 22.1338 (2.2755) 7.4614 (1.5881)	0.6936 (0.0346) 0.7310 (0.0126) 0.7299 (0.0272) 0.7372 (0.0764) 0.7331 (0.0463) 0.8838 (0.0351)
Chair-Gap	MVE (NLL) MVE(Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD) Ours: VeBNN (pSGLD)	0.4545 (0.2802) 0.2264 (0.0677) 0.4683 (0.1803) 15.5046 (2.3557) 14.2239 (1.8959) 7.0322 (0.9985) 0.5123 (0.1273)	2.0921 (0.9338) 1.3059 (0.1362) 3.8559 (0.5422) 99.6555 (14.3327) 63.0351 (6.7923) 44.2601 (1.2242) 2.3708 (0.7181)	12.6657 (0.5186)	0.6936 (0.0346) 0.7310 (0.0126) 0.7299 (0.0272) 0.7372 (0.0764) 0.7331 (0.0463) 0.8838 (0.0351) 0.7387 (0.0198)
	Ours: VeBNN (MC-Dropout) 0.5583 (1.2076) 47.8594 (5.1651) 10.8850 (1.5156) 0.	0.6936 (0.0346) 0.7310 (0.0126) 0.7299 (0.0272) 0.7372 (0.0764) 0.7331 (0.0463) 0.8838 (0.0351) 0.7387 (0.0198) 0.6226 (0.0071) 0.6170 (0.0063) 0.6624 (0.0135)			
Chair-Gap Chair-Tail		0.6936 (0.0346) 0.7310 (0.0126) 0.7299 (0.0272) 0.7372 (0.0764) 0.7331 (0.0463) 0.8838 (0.0351) 0.7387 (0.0198) 0.6226 (0.0071) 0.6170 (0.0063) 0.6624 (0.0135)			
	MVE (NLL) MVE(Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD) Ours: VeBNN (pSGLD) MVE (NLL) MVE(Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout)	0.4545 (0.2802) 0.2264 (0.0677) 0.4683 (0.1803) 15.5046 (2.3557) 14.2239 (1.8959) 7.0322 (0.9985) 0.5123 (0.1273) 0.2412 (0.0917) 0.1337 (0.0190) 0.3828 (0.2221) 14.5111 (1.9255) 8.3647 (1.3729)	2.0921 (0.9338) 1.3059 (0.1362) 3.8559 (0.5422) 99.6555 (14.3327) 63.0351 (6.7923) 44.2601 (1.2242) 2.3708 (0.7181) 1.0941 (0.3829) 0.7691 (0.0814) 2.5303 (1.0144) 81.3095 (4.3530) 37.2592 (5.1038)	1.4182 (0.2737) 1.1909 (0.0607) 1.7800 (0.2434) 15.4490 (1.7975) 22.1338 (2.2755) 7.4614 (1.5881) 1.6514 (0.2616) 3.1580 (0.2162) 1.4892 (0.0454) 2.8544 (0.2497) 12.6523 (0.3463) 20.9121 (1.2496)	0.6936 (0.0346) 0.7310 (0.0126) 0.7299 (0.0272) 0.7372 (0.0764) 0.7331 (0.0463) 0.8838 (0.0351) 0.7387 (0.0198) 0.6226 (0.0071) 0.6170 (0.0063) 0.6624 (0.0135) 0.6942 (0.0217) 0.5090 (0.0686)
	MVE (NLL) MVE(Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD) Ours: VeBNN (pSGLD) MVE (NLL) MVE(Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD)	0.4545 (0.2802) 0.2264 (0.0677) 0.4683 (0.1803) 15.5046 (2.3557) 14.2239 (1.8959) 7.0322 (0.9985) 0.5123 (0.1273) 0.2412 (0.0917) 0.1337 (0.0190) 0.3828 (0.2221) 14.5111 (1.9255) 8.3647 (1.3729) 6.7768 (3.0065)	2.0921 (0.9338) 1.3059 (0.1362) 3.8559 (0.5422) 99.6555 (14.3327) 63.0351 (6.7923) 44.2601 (1.2242) 2.3708 (0.7181) 1.0941 (0.3829) 0.7691 (0.0814) 2.5303 (1.0144) 81.3095 (4.3530) 37.2592 (5.1038) 43.6364 (1.5572)	1.4182 (0.2737) 1.1909 (0.0607) 1.7800 (0.2434) 15.4490 (1.7975) 22.1338 (2.2755) 7.4614 (1.5881) 1.6514 (0.2616) 3.1580 (0.2162) 1.4892 (0.0454) 2.8544 (0.2497) 12.6523 (0.3463) 20.9121 (1.2496) 10.0926 (2.3311)	0.6936 (0.0346) 0.7310 (0.0126) 0.7299 (0.0272) 0.7372 (0.0764) 0.7331 (0.0463) 0.8838 (0.0351) 0.7387 (0.0198) 0.6226 (0.0071) 0.6170 (0.0063) 0.6624 (0.0135) 0.6942 (0.0217) 0.5090 (0.0686) 0.7665 (0.0737)
	MVE (NLL) MVE(Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD) Ours: VeBNN (pSGLD) MVE (NLL) MVE(Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD)	0.4545 (0.2802) 0.2264 (0.0677) 0.4683 (0.1803) 15.5046 (2.3557) 14.2239 (1.8959) 7.0322 (0.9985) 0.5123 (0.1273) 0.2412 (0.0917) 0.1337 (0.0190) 0.3828 (0.2221) 14.5111 (1.9255) 8.3647 (1.3729) 6.7768 (3.0065)	2.0921 (0.9338) 1.3059 (0.1362) 3.8559 (0.5422) 99.6555 (14.3327) 63.0351 (6.7923) 44.2601 (1.2242) 2.3708 (0.7181) 1.0941 (0.3829) 0.7691 (0.0814) 2.5303 (1.0144) 81.3095 (4.3530) 37.2592 (5.1038) 43.6364 (1.5572)	1.4182 (0.2737) 1.1909 (0.0607) 1.7800 (0.2434) 15.4490 (1.7975) 22.1338 (2.2755) 7.4614 (1.5881) 1.6514 (0.2616) 3.1580 (0.2162) 1.4892 (0.0454) 2.8544 (0.2497) 12.6523 (0.3463) 20.9121 (1.2496) 10.0926 (2.3311)	0.6936 (0.0346) 0.7310 (0.0126) 0.7299 (0.0272) 0.7372 (0.0764) 0.7331 (0.0463) 0.8838 (0.0351) 0.7387 (0.0198) 0.6226 (0.0071) 0.6170 (0.0063) 0.6624 (0.0135) 0.6942 (0.0217) 0.5090 (0.0686)
	MVE (NLL) MVE(Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD) Ours: VeBNN (pSGLD) MVE (NLL) MVE(Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD) Ours: VeBNN (pSGLD) MVE (NLL)	0.4545 (0.2802) 0.2264 (0.0677) 0.4683 (0.1803) 15.5046 (2.3557) 14.2239 (1.8959) 7.0322 (0.9985) 0.5123 (0.1273) 0.2412 (0.0917) 0.1337 (0.0190) 0.3828 (0.2221) 14.5111 (1.9255) 8.3647 (1.3729) 6.7768 (3.0065) 0.4755 (0.2692)	2.0921 (0.9338) 1.3059 (0.1362) 3.8559 (0.5422) 99.6555 (14.3327) 63.0351 (6.7923) 44.2601 (1.2242) 2.3708 (0.7181) 1.0941 (0.3829) 0.7691 (0.0814) 2.5303 (1.0144) 81.3095 (4.3530) 37.2592 (5.1038) 43.6364 (1.5572) 2.5667 (1.5456) 535.1390 (215.4460)	1.4182 (0.2737) 1.1909 (0.0607) 1.7800 (0.2434) 15.4490 (1.7975) 22.1338 (2.2755) 7.4614 (1.5881) 1.6514 (0.2616) 3.1580 (0.2162) 1.4892 (0.0454) 2.8544 (0.2497) 12.6523 (0.3463) 20.9121 (1.2496) 10.0926 (2.3311) 3.7509 (0.1950) 262.9090 (12.6078)	0.6936 (0.0346) 0.7310 (0.0126) 0.7299 (0.0272) 0.7372 (0.0764) 0.7331 (0.0463) 0.8838 (0.0351) 0.7387 (0.0198) 0.6226 (0.0071) 0.6170 (0.0063) 0.6624 (0.0135) 0.6942 (0.0217) 0.5090 (0.0686) 0.7665 (0.0737) 0.6684 (0.0446) 0.7973 (0.0270)
	MVE (NLL) MVE(Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD) Ours: VeBNN (pSGLD) MVE (NLL) MVE(Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD) Ours: VeBNN (pSGLD) MVE (NLL) MVE (Ensembles)	0.4545 (0.2802) 0.2264 (0.0677) 0.4683 (0.1803) 15.5046 (2.3557) 14.2239 (1.8959) 7.0322 (0.9985) 0.5123 (0.1273) 0.2412 (0.0917) 0.1337 (0.0190) 0.3828 (0.2221) 14.5111 (1.9255) 8.3647 (1.3729) 6.7768 (3.0065) 0.4755 (0.2692) 105.4170 (1.1518) 100.6390 (0.4642)	2.0921 (0.9338) 1.3059 (0.1362) 3.8559 (0.5422) 99.6555 (14.3327) 63.0351 (6.7923) 44.2601 (1.2242) 2.3708 (0.7181) 1.0941 (0.3829) 0.7691 (0.0814) 2.5303 (1.0144) 81.3095 (4.3530) 37.2592 (5.1038) 43.6364 (1.5572) 2.5667 (1.5456) 535.1390 (215.4460) 472.1400 (85.2654)	1.4182 (0.2737) 1.1909 (0.0607) 1.7800 (0.2434) 15.4490 (1.7975) 22.1338 (2.2755) 7.4614 (1.5881) 1.6514 (0.2616) 3.1580 (0.2162) 1.4892 (0.0454) 2.8544 (0.2497) 12.6523 (0.3463) 20.9121 (1.2496) 10.0926 (2.3311) 3.7509 (0.1950) 262.9090 (12.6078) 241.3947 (2.4395)	0.6936 (0.0346) 0.7310 (0.0126) 0.7299 (0.0272) 0.7372 (0.0764) 0.7331 (0.0463) 0.8838 (0.0351) 0.7387 (0.0198) 0.6226 (0.0071) 0.6170 (0.0063) 0.6624 (0.0135) 0.6942 (0.0217) 0.5090 (0.0686) 0.7665 (0.0737) 0.6684 (0.0446) 0.7973 (0.0270) 0.8229 (0.0134)
Chair-Tail	MVE (NLL) MVE(Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD) Ours: VeBNN (pSGLD) MVE (NLL) MVE(Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD) Ours: VeBNN (pSGLD) MVE (NLL) MVE (Ensembles) Evidential	0.4545 (0.2802) 0.2264 (0.0677) 0.4683 (0.1803) 15.5046 (2.3557) 14.2239 (1.8959) 7.0322 (0.9985) 0.5123 (0.1273) 0.2412 (0.0917) 0.1337 (0.0190) 0.3828 (0.2221) 14.5111 (1.9255) 8.3647 (1.3729) 6.7768 (3.0065) 0.4755 (0.2692) 105.4170 (1.1518) 100.6390 (0.4642)	2.0921 (0.9338) 1.3059 (0.1362) 3.8559 (0.5422) 99.6555 (14.3327) 63.0351 (6.7923) 44.2601 (1.2242) 2.3708 (0.7181) 1.0941 (0.3829) 0.7691 (0.0814) 2.5303 (1.0144) 81.3095 (4.3530) 37.2592 (5.1038) 43.6364 (1.5572) 2.5667 (1.5456) 535.1390 (215.4460) 472.1400 (85.2654)	1.4182 (0.2737) 1.1909 (0.0607) 1.7800 (0.2434) 15.4490 (1.7975) 22.1338 (2.2755) 7.4614 (1.5881) 1.6514 (0.2616) 3.1580 (0.2162) 1.4892 (0.0454) 2.8544 (0.2497) 12.6523 (0.3463) 20.9121 (1.2496) 10.0926 (2.3311) 3.7509 (0.1950) 262.9090 (12.6078) 241.3947 (2.4395)	0.6936 (0.0346) 0.7310 (0.0126) 0.7299 (0.0272) 0.7372 (0.0764) 0.7331 (0.0463) 0.8838 (0.0351) 0.7387 (0.0198) 0.6226 (0.0071) 0.6170 (0.0063) 0.6624 (0.0135) 0.6942 (0.0217) 0.5090 (0.0686) 0.7665 (0.0737) 0.6684 (0.0446) 0.7973 (0.0270)
	MVE (NLL) MVE(Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD) Ours: VeBNN (pSGLD) MVE (NLL) MVE(Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD) Ours: VeBNN (pSGLD) MVE (NLL) MVE (Ensembles) Evidential MVE (MC-Dropout) MVE (Ensembles) Evidential MVE (MC-Dropout)	0.4545 (0.2802) 0.2264 (0.0677) 0.4683 (0.1803) 15.5046 (2.3557) 14.2239 (1.8959) 7.0322 (0.9985) 0.5123 (0.1273) 0.2412 (0.0917) 0.1337 (0.0190) 0.3828 (0.2221) 14.5111 (1.9255) 8.3647 (1.3729) 6.7768 (3.0065) 0.4755 (0.2692) 105.4170 (1.1518) 100.6390 (0.4642) 99.6063 (6.9670) 258.0571 (17.4759)	2.0921 (0.9338) 1.3059 (0.1362) 3.8559 (0.5422) 99.6555 (14.3327) 63.0351 (6.7923) 44.2601 (1.2242) 2.3708 (0.7181) 1.0941 (0.3829) 0.7691 (0.0814) 2.5303 (1.0144) 81.3095 (4.3530) 37.2592 (5.1038) 43.6364 (1.5572) 2.5667 (1.5456) 535.1390 (215.4460) 472.1400 (85.2654) 4425.1200 (45.0808) 1356.7668 (118.0702)	1.4182 (0.2737) 1.1909 (0.0607) 1.7800 (0.2434) 15.4490 (1.7975) 22.1338 (2.2755) 7.4614 (1.5881) 1.6514 (0.2616) 3.1580 (0.2162) 1.4892 (0.0454) 2.8544 (0.2497) 12.6523 (0.3463) 20.9121 (1.2496) 10.0926 (2.3311) 3.7509 (0.1950) 262.9090 (12.6078) 241.3947 (2.4395) 260.6220 (11.7347) 645.0634 (5.2036)	0.6936 (0.0346) 0.7310 (0.0126) 0.7299 (0.0272) 0.7372 (0.0764) 0.7331 (0.0463) 0.8838 (0.0351) 0.7387 (0.0198) 0.6226 (0.0071) 0.6170 (0.0063) 0.6624 (0.0135) 0.6942 (0.0217) 0.5090 (0.0686) 0.7665 (0.0737) 0.6684 (0.0446) 0.7973 (0.0270) 0.8229 (0.0134) 0.5583 (0.3128) 0.7328 (0.0290)
Chair-Tail	MVE (NLL) MVE(Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD) Ours: VeBNN (pSGLD) MVE (NLL) MVE(Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD) Ours: VeBNN (pSGLD) MVE (NLL) MVE (Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (pSGLD) Ours: VeBNN (pSGLD) MVE (NLL) MVE (Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout)	0.4545 (0.2802) 0.2264 (0.0677) 0.4683 (0.1803) 15.5046 (2.3557) 14.2239 (1.8959) 7.0322 (0.9985) 0.5123 (0.1273) 0.2412 (0.0917) 0.1337 (0.0190) 0.3828 (0.2221) 14.5111 (1.9255) 8.3647 (1.3729) 6.7768 (3.0065) 0.4755 (0.2692) 105.4170 (1.1518) 100.6390 (0.4642) 99.6063 (6.9670) 258.0571 (17.4759) 305.9113 (35.1769)	2.0921 (0.9338) 1.3059 (0.1362) 3.8559 (0.5422) 99.6555 (14.3327) 63.0351 (6.7923) 44.2601 (1.2242) 2.3708 (0.7181) 1.0941 (0.3829) 0.7691 (0.0814) 2.5303 (1.0144) 81.3095 (4.3530) 37.2592 (5.1038) 43.6364 (1.5572) 2.5667 (1.5456) 535.1390 (215.4460) 472.1400 (85.2654) 425.1200 (45.0808) 1356.7668 (118.0702) 1350.9520 (188.2830)	1.4182 (0.2737) 1.1909 (0.0607) 1.7800 (0.2434) 15.4490 (1.7975) 22.1338 (2.2755) 7.4614 (1.5881) 1.6514 (0.2616) 3.1580 (0.2162) 1.4892 (0.0454) 2.8544 (0.2497) 12.6523 (0.3463) 20.9121 (1.2496) 10.0926 (2.3311) 3.7509 (0.1950) 262.9090 (12.6078) 241.3947 (2.4395) 260.6220 (11.7347) 645.0634 (5.2036) 464.4328 (25.0852)	0.6936 (0.0346) 0.7310 (0.0126) 0.7299 (0.0272) 0.7372 (0.0764) 0.7331 (0.0463) 0.8838 (0.0351) 0.7387 (0.0198) 0.6226 (0.0071) 0.6170 (0.0063) 0.6624 (0.0135) 0.6942 (0.0217) 0.5090 (0.0686) 0.7665 (0.0737) 0.6684 (0.0446) 0.7973 (0.0270) 0.8229 (0.0134) 0.5583 (0.3128) 0.7328 (0.0290) 0.7841 (0.0399)
Chair-Tail	MVE (NLL) MVE(Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD) Ours: VeBNN (pSGLD) MVE (NLL) MVE(Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD) Ours: VeBNN (pSGLD) MVE (NLL) MVE (Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD) Ours: VeBNN (MC-Dropout) Durs: VeBNN (MC-Dropout) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD)	0.4545 (0.2802) 0.2264 (0.0677) 0.4683 (0.1803) 15.5046 (2.3557) 14.2239 (1.8959) 7.0322 (0.9985) 0.5123 (0.1273) 0.2412 (0.0917) 0.1337 (0.0190) 0.3828 (0.2221) 14.5111 (1.9255) 8.3647 (1.3729) 6.7768 (3.0065) 0.4755 (0.2692) 105.4170 (1.1518) 100.6390 (0.4642) 99.6063 (6.9670) 258.0571 (17.4759) 305.9113 (35.1769) 127.9469 (5.2489)	2.0921 (0.9338) 1.3059 (0.1362) 3.8559 (0.5422) 99.6555 (14.3327) 63.0351 (6.7923) 44.2601 (1.2242) 2.3708 (0.7181) 1.0941 (0.3829) 0.7691 (0.0814) 2.5303 (1.0144) 81.3095 (4.3530) 37.2592 (5.1038) 43.6364 (1.5572) 2.5667 (1.5456) 535.1390 (215.4460) 472.1400 (85.2654) 425.1200 (45.0808) 1356.7668 (118.0702) 1350.9520 (188.2830) 550.7750 (25.6849)	1.4182 (0.2737) 1.1909 (0.0607) 1.7800 (0.2434) 15.4490 (1.7975) 22.1338 (2.2755) 7.4614 (1.5881) 1.6514 (0.2616) 3.1580 (0.2162) 1.4892 (0.0454) 2.8544 (0.2497) 12.6523 (0.3463) 20.9121 (1.2496) 10.0926 (2.3311) 3.7509 (0.1950) 262.9090 (12.6078) 241.3947 (2.4395) 260.6220 (11.7347) 645.0634 (5.2036) 464.4328 (25.0852) 267.6061 (4.5646)	0.6936 (0.0346) 0.7310 (0.0126) 0.7299 (0.0272) 0.7372 (0.0764) 0.7331 (0.0463) 0.8838 (0.0351) 0.7387 (0.0198) 0.6226 (0.0071) 0.6170 (0.0063) 0.6624 (0.0135) 0.6942 (0.0217) 0.5090 (0.0686) 0.7665 (0.0737) 0.6684 (0.0446) 0.7973 (0.0270) 0.8229 (0.0134) 0.5583 (0.3128) 0.7328 (0.0290) 0.7841 (0.0399) 0.8282 (0.0181)
Chair-Tail	MVE (NLL) MVE(Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD) Ours: VeBNN (pSGLD) MVE (NLL) MVE(Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD) Ours: VeBNN (pSGLD) MVE (NLL) MVE (Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (pSGLD) Ours: VeBNN (pSGLD) MVE (NLL) MVE (Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout)	0.4545 (0.2802) 0.2264 (0.0677) 0.4683 (0.1803) 15.5046 (2.3557) 14.2239 (1.8959) 7.0322 (0.9985) 0.5123 (0.1273) 0.2412 (0.0917) 0.1337 (0.0190) 0.3828 (0.2221) 14.5111 (1.9255) 8.3647 (1.3729) 6.7768 (3.0065) 0.4755 (0.2692) 105.4170 (1.1518) 100.6390 (0.4642) 99.6063 (6.9670) 258.0571 (17.4759) 305.9113 (35.1769)	2.0921 (0.9338) 1.3059 (0.1362) 3.8559 (0.5422) 99.6555 (14.3327) 63.0351 (6.7923) 44.2601 (1.2242) 2.3708 (0.7181) 1.0941 (0.3829) 0.7691 (0.0814) 2.5303 (1.0144) 81.3095 (4.3530) 37.2592 (5.1038) 43.6364 (1.5572) 2.5667 (1.5456) 535.1390 (215.4460) 472.1400 (85.2654) 425.1200 (45.0808) 1356.7668 (118.0702) 1350.9520 (188.2830)	1.4182 (0.2737) 1.1909 (0.0607) 1.7800 (0.2434) 15.4490 (1.7975) 22.1338 (2.2755) 7.4614 (1.5881) 1.6514 (0.2616) 3.1580 (0.2162) 1.4892 (0.0454) 2.8544 (0.2497) 12.6523 (0.3463) 20.9121 (1.2496) 10.0926 (2.3311) 3.7509 (0.1950) 262.9090 (12.6078) 241.3947 (2.4395) 260.6220 (11.7347) 645.0634 (5.2036) 464.4328 (25.0852) 267.6061 (4.5646)	0.6936 (0.0346) 0.7310 (0.0126) 0.7299 (0.0272) 0.7372 (0.0764) 0.7331 (0.0463) 0.8838 (0.0351) 0.7387 (0.0198) 0.6226 (0.0071) 0.6170 (0.0063) 0.6624 (0.0135) 0.6942 (0.0217) 0.5090 (0.0686) 0.7665 (0.0737) 0.6684 (0.0446) 0.7973 (0.0270) 0.8229 (0.0134) 0.5583 (0.3128) 0.7328 (0.0290) 0.7841 (0.0399)
Chair-Tail	MVE (NLL) MVE(Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD) Ours: VeBNN (pSGLD) MVE (NLL) MVE(Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD) Ours: VeBNN (pSGLD) MVE (NLL) MVE (Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD) Ours: VeBNN (pSGLD) MVE (NLL) MVE (Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) Ours: VeBNN (MC-Dropout) Ours: VeBNN (pSGLD) MVE (NLL)	0.4545 (0.2802) 0.2264 (0.0677) 0.4683 (0.1803) 15.5046 (2.3557) 14.2239 (1.8959) 7.0322 (0.9985) 0.5123 (0.1273) 0.2412 (0.0917) 0.1337 (0.0190) 0.3828 (0.2221) 14.5111 (1.9255) 8.3647 (1.3729) 6.7768 (3.0065) 0.4755 (0.2692) 105.4170 (1.1518) 100.6390 (0.4642) 99.6063 (6.9670) 258.0871 (17.4759) 305.9113 (35.1769) 127.9469 (5.2489) 105.5894 (4.1459) 217.8770 (1.7249)	2.0921 (0.9338) 1.3059 (0.1362) 3.8559 (0.5422) 99.6555 (14.3327) 63.0351 (6.7923) 44.2601 (1.2242) 2.3708 (0.7181) 1.0941 (0.3829) 0.7691 (0.0814) 2.5303 (1.0144) 81.3095 (4.3530) 37.2592 (5.1038) 43.6364 (1.5572) 2.5667 (1.5456) 535.1390 (215.4460) 472.1400 (85.2654) 425.1200 (45.0808) 1356.7668 (118.0702) 1350.9520 (188.2830) 550.7750 (25.6849) 446.7374 (9.2086)	1.4182 (0.2737) 1.1909 (0.0607) 1.7800 (0.2434) 15.4490 (1.7975) 22.1338 (2.2755) 7.4614 (1.5881) 1.6514 (0.2616) 3.1580 (0.2162) 1.4892 (0.0454) 2.8544 (0.2497) 12.6523 (0.3463) 20.9121 (1.2496) 10.0926 (2.3311) 3.7509 (0.1950) 262.9090 (12.6078) 241.3947 (2.4395) 260.6220 (11.7347) 645.0634 (5.2036) 464.4328 (25.0852) 267.6061 (4.5646) 260.3773 (9.0958) 330.8905 (32.7377)	0.6936 (0.0346) 0.7310 (0.0126) 0.7299 (0.0272) 0.7372 (0.0764) 0.7331 (0.0463) 0.8838 (0.0351) 0.7387 (0.0198) 0.6226 (0.0071) 0.6170 (0.0063) 0.6624 (0.0135) 0.6942 (0.0217) 0.5090 (0.0686) 0.7665 (0.0737) 0.6684 (0.0446) 0.7973 (0.0270) 0.8229 (0.0134) 0.5583 (0.3128) 0.7328 (0.0290) 0.7841 (0.0399) 0.8282 (0.0181) 0.8036 (0.0038) 0.6988 (0.0662)
Chair-Tail	MVE (NLL) MVE(Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD) Ours: VeBNN (pSGLD) MVE (NLL) MVE(Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD) MVE (NLL) MVE (Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) Ours: VeBNN (pSGLD) MVE (Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD) Ours: VeBNN (pSGLD) MVE (NLL) MVE (NLLL) MVE (Ensembles)	0.4545 (0.2802) 0.2264 (0.0677) 0.4683 (0.1803) 15.5046 (2.3557) 14.2239 (1.8959) 7.0322 (0.9985) 0.5123 (0.1273) 0.2412 (0.0917) 0.1337 (0.0190) 0.3828 (0.2221) 14.5111 (1.9255) 8.3647 (1.3729) 6.7768 (3.0065) 0.4755 (0.2692) 105.4170 (1.1518) 100.6390 (0.4642) 99.6063 (6.9670) 258.0571 (17.4759) 305.9113 (35.1769) 127.9469 (5.2489) 105.5894 (4.1459) 208.4870 (1.7249) 208.4870 (1.7249)	2.0921 (0.9338) 1.3059 (0.1362) 3.8559 (0.5422) 99.6555 (14.3327) 63.0351 (6.7923) 44.2601 (1.2242) 2.3708 (0.7181) 1.0941 (0.3829) 0.7691 (0.0814) 2.5303 (1.0144) 81.3095 (4.3530) 37.2592 (5.1038) 43.6364 (1.5572) 2.5667 (1.5456) 535.1390 (215.4460) 472.1400 (85.2654) 425.1200 (45.0808) 1356.7668 (118.0702) 1350.9520 (188.2830) 550.7750 (25.6849) 446.7374 (9.2086) 929.5620 (47.6606) 885.3490 (27.8584)	1.4182 (0.2737) 1.1909 (0.0607) 1.7800 (0.2434) 15.4490 (1.7975) 22.1338 (2.2755) 7.4614 (1.5881) 1.6514 (0.2616) 3.1580 (0.2162) 1.4892 (0.0454) 2.8544 (0.2497) 12.6523 (0.3463) 20.9121 (1.2496) 10.0926 (2.3311) 3.7509 (0.1950) 262.9090 (12.6078) 241.3947 (2.4395) 260.6220 (11.7347) 645.0634 (5.2036) 464.4328 (25.0852) 267.6061 (4.5646) 260.3773 (9.0958) 330.8905 (32.7377) 295.1917 (9.6539)	0.6936 (0.0346) 0.7310 (0.0126) 0.7299 (0.0272) 0.7372 (0.0764) 0.7331 (0.0463) 0.8838 (0.0351) 0.7387 (0.0198) 0.6226 (0.0071) 0.6170 (0.0063) 0.6624 (0.0135) 0.5990 (0.0686) 0.7665 (0.0737) 0.6684 (0.0446) 0.7973 (0.0270) 0.8229 (0.0134) 0.5583 (0.3128) 0.7328 (0.0290) 0.7841 (0.0399) 0.8282 (0.0181) 0.8036 (0.0038) 0.6988 (0.0662) 0.8123 (0.0667)
Chair-Tail Skin	MVE (NLL) MVE(Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD) Ours: VeBNN (pSGLD) MVE (NLL) MVE(Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD) Ours: VeBNN (pSGLD) MVE (NLL) MVE (Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (pSGLD) MVE (Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD) Ours: VeBNN (pSGLD) MVE (NLL) MVE (Ensembles) Evidential	0.4545 (0.2802) 0.2264 (0.0677) 0.4683 (0.1803) 15.5046 (2.3557) 14.2239 (1.8959) 7.0322 (0.9985) 0.5123 (0.1273) 0.2412 (0.0917) 0.1337 (0.0190) 0.3828 (0.2221) 14.5111 (1.9255) 8.3647 (1.3729) 6.7768 (3.0065) 0.4755 (0.2692) 105.4170 (1.1518) 100.6390 (0.4642) 99.6063 (6.9670) 258.0571 (17.4759) 305.9113 (35.1769) 127.9469 (5.2489) 105.5894 (4.1459) 217.8770 (1.7249) 208.4870 (1.0358) 180.5486 (2.4021)	2.0921 (0.9338) 1.3059 (0.1362) 3.8559 (0.5422) 99.6555 (14.3327) 63.0351 (6.7923) 44.2601 (1.2242) 2.3708 (0.7181) 1.0941 (0.3829) 0.7691 (0.0814) 2.5303 (1.0144) 81.3095 (4.3530) 37.2592 (5.1038) 43.6364 (1.5572) 2.5667 (1.5456) 535.1390 (215.4460) 472.1400 (85.2654) 4425.1200 (45.0808) 1356.7668 (118.0702) 1350.9520 (188.2830) 550.7750 (25.6849) 446.7374 (9.2086) 929.5620 (47.6606) 885.3490 (27.8584) 747.9998 (6.5119)	1.4182 (0.2737) 1.1909 (0.0607) 1.7800 (0.2434) 15.4490 (1.7975) 22.1338 (2.2755) 7.4614 (1.5881) 1.6514 (0.2616) 3.1580 (0.2162) 1.4892 (0.0454) 2.8544 (0.2497) 12.6523 (0.3463) 20.9121 (1.2496) 10.0926 (2.3311) 3.7509 (0.1950) 262.9090 (12.6078) 241.3947 (2.4395) 260.6220 (11.7347) 645.0634 (5.2036) 464.4328 (25.0852) 267.6061 (4.5646) 260.3773 (9.0958) 330.8905 (32.7377) 295.1917 (9.6539) 354.8147 (12.1219)	0.6936 (0.0346) 0.7310 (0.0126) 0.7329 (0.0272) 0.7372 (0.0764) 0.7331 (0.0463) 0.8838 (0.0351) 0.7387 (0.0198) 0.6226 (0.0071) 0.6170 (0.0063) 0.6624 (0.0135) 0.6942 (0.0217) 0.5090 (0.0686) 0.7665 (0.0737) 0.6684 (0.0446) 0.7973 (0.0270) 0.8229 (0.0134) 0.5583 (0.3128) 0.7328 (0.0290) 0.7841 (0.0399) 0.8282 (0.0181) 0.8036 (0.0038) 0.6988 (0.0662) 0.8123 (0.0617) 0.6857 (0.0174)
Chair-Tail	MVE (NLL) MVE(Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD) Ours: VeBNN (pSGLD) MVE (NLL) MVE(Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD) Ours: VeBNN (pSGLD) MVE (NLL) MVE (Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (pSGLD) Ours: VeBNN (pSGLD) MVE (NLL) MVE (Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (pSGLD) Ours: VeBNN (pSGLD) MVE (NLL) MVE (Ensembles) Evidential MVE (MC-Dropout)	0.4545 (0.2802) 0.2264 (0.0677) 0.4683 (0.1803) 15.5046 (2.3557) 14.2239 (1.8959) 7.0322 (0.9985) 0.5123 (0.1273) 0.2412 (0.0917) 0.1337 (0.0190) 0.3828 (0.2221) 14.5111 (1.9255) 8.3647 (1.3729) 6.7768 (3.0065) 0.4755 (0.2692) 105.4170 (1.1518) 100.6390 (0.4642) 99.6063 (6.9670) 258.0571 (17.4759) 305.9113 (35.1769) 127.9469 (5.2489) 105.5894 (4.1459) 217.8770 (1.7249) 208.4870 (1.0358) 180.5486 (2.4021) 302.0488 (6.7523)	2.0921 (0.9338) 1.3059 (0.1362) 3.8559 (0.5422) 99.6555 (14.3327) 63.0351 (6.7923) 44.2601 (1.2242) 2.3708 (0.7181) 1.0941 (0.3829) 0.7691 (0.0814) 2.5303 (1.0144) 81.3095 (4.3530) 37.2592 (5.1038) 43.6364 (1.5572) 2.5667 (1.5456) 535.1390 (215.4460) 472.1400 (85.2654) 425.1200 (45.0808) 1356.7668 (118.0702) 1350.9520 (188.2830) 550.7750 (25.6849) 446.7374 (9.2086) 929.5620 (47.6606) 885.3490 (27.8584) 747.9998 (6.5119) 1455.8063 (51.6707)	1.4182 (0.2737) 1.1909 (0.0607) 1.7800 (0.2434) 15.4490 (1.7975) 22.1338 (2.2755) 7.4614 (1.5881) 1.6514 (0.2616) 3.1580 (0.2162) 1.4892 (0.0454) 2.8544 (0.2497) 12.6523 (0.3463) 20.9121 (1.2496) 10.0926 (2.3311) 3.7509 (0.1950) 262.9090 (12.6078) 241.3947 (2.4395) 260.6220 (11.7347) 645.0634 (5.2036) 464.4328 (25.0852) 267.6061 (4.5646) 260.3773 (9.0958) 330.8905 (32.7377) 295.1917 (9.6539) 354.8147 (12.1219) 497.3040 (8.2818)	0.6936 (0.0346) 0.7310 (0.0126) 0.7299 (0.0272) 0.7372 (0.0764) 0.7331 (0.0463) 0.8838 (0.0351) 0.7387 (0.0198) 0.6226 (0.0071) 0.6170 (0.0063) 0.6624 (0.0135) 0.6942 (0.0217) 0.5090 (0.0686) 0.7665 (0.0737) 0.6684 (0.0446) 0.7973 (0.0270) 0.8229 (0.0134) 0.5583 (0.3128) 0.7328 (0.0290) 0.7841 (0.0390) 0.7841 (0.0399) 0.8282 (0.0181) 0.8036 (0.0038) 0.6988 (0.0662) 0.8123 (0.0617) 0.6887 (0.0174) 0.7156 (0.0346)
Chair-Tail Skin	MVE (NLL) MVE(Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD) Ours: VeBNN (pSGLD) MVE (NLL) MVE(Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD) Ours: VeBNN (pSGLD) MVE (NLL) MVE (Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) Durs: VeBNN (pSGLD) MVE (NLL) MVE (Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD) Ours: VeBNN (pSGLD) MVE (NLL) MVE (Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (pSGLD)	0.4545 (0.2802) 0.2264 (0.0677) 0.4683 (0.1803) 15.5046 (2.3557) 14.2239 (1.8959) 7.0322 (0.9985) 0.5123 (0.1273) 0.2412 (0.0917) 0.1337 (0.0190) 0.3828 (0.2221) 14.5111 (1.9255) 8.3647 (1.3729) 6.7768 (3.0065) 0.4755 (0.2692) 105.4170 (1.1518) 100.6390 (0.4642) 99.6063 (6.9670) 258.0571 (17.4759) 305.9113 (35.1769) 127.9469 (5.2489) 105.5894 (4.1459) 217.8770 (1.7249) 208.4870 (1.0358) 180.5486 (2.4021) 302.0488 (6.7523) 296.4448 (9.5291)	2.0921 (0.9338) 1.3059 (0.1362) 3.8559 (0.5422) 99.6555 (14.3327) 63.0351 (6.7923) 44.2601 (1.2242) 2.3708 (0.7181) 1.0941 (0.3829) 0.7691 (0.0814) 2.5303 (1.0144) 81.3095 (4.3530) 37.2592 (5.1038) 43.6364 (1.5572) 2.5667 (1.5456) 535.1390 (215.4460) 472.1400 (85.2654) 425.1200 (45.0808) 1356.7668 (118.0702) 1350.9520 (188.2830) 550.7750 (25.6849) 446.7374 (9.2086) 929.5620 (47.6606) 885.3490 (27.8584) 747.9998 (6.5119) 1455.8003 (51.6707) 1295.8913 (57.9997)	1.4182 (0.2737) 1.1909 (0.0607) 1.7800 (0.2434) 15.4490 (1.7975) 22.1338 (2.2755) 7.4614 (1.5881) 1.6514 (0.2616) 3.1580 (0.2162) 1.4892 (0.0454) 2.8544 (0.2497) 12.6523 (0.3463) 20.9121 (1.2496) 10.0926 (2.3311) 3.7509 (0.1950) 262.9090 (12.6078) 241.3947 (2.4395) 260.6220 (11.7347) 645.0634 (5.2036) 464.4328 (25.0852) 267.6061 (4.5646) 260.3773 (9.0958) 330.8905 (32.7377) 295.1917 (9.6539) 354.8147 (12.1219) 497.3040 (8.2818) 537.5555 (38.6030)	0.6936 (0.0346) 0.7310 (0.0126) 0.7299 (0.0272) 0.7372 (0.0764) 0.7331 (0.0463) 0.8838 (0.0351) 0.7387 (0.0198) 0.6226 (0.0071) 0.6170 (0.0063) 0.6624 (0.0135) 0.6942 (0.0217) 0.5090 (0.0686) 0.7665 (0.0737) 0.6684 (0.0446) 0.7973 (0.0270) 0.8229 (0.0134) 0.5583 (0.3128) 0.7328 (0.0290) 0.7841 (0.0399) 0.8282 (0.0181) 0.8036 (0.0038) 0.6988 (0.0662) 0.8123 (0.0617) 0.6857 (0.0174) 0.7156 (0.0346) 0.7150 (0.0346) 0.7150 (0.0346)
Chair-Tail Skin	MVE (NLL) MVE(Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD) Ours: VeBNN (pSGLD) MVE (NLL) MVE(Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (MC-Dropout) BNN-End-to-End (pSGLD) Ours: VeBNN (pSGLD) MVE (NLL) MVE (Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (pSGLD) Ours: VeBNN (pSGLD) MVE (NLL) MVE (Ensembles) Evidential MVE (MC-Dropout) Ours: VeBNN (pSGLD) Ours: VeBNN (pSGLD) MVE (NLL) MVE (Ensembles) Evidential MVE (MC-Dropout)	0.4545 (0.2802) 0.2264 (0.0677) 0.4683 (0.1803) 15.5046 (2.3557) 14.2239 (1.8959) 7.0322 (0.9985) 0.5123 (0.1273) 0.2412 (0.0917) 0.1337 (0.0190) 0.3828 (0.2221) 14.5111 (1.9255) 8.3647 (1.3729) 6.7768 (3.0065) 0.4755 (0.2692) 105.4170 (1.1518) 100.6390 (0.4642) 99.6063 (6.9670) 258.0571 (17.4759) 305.9113 (35.1769) 127.9469 (5.2489) 105.5894 (4.1459) 217.8770 (1.7249) 208.4870 (1.0358) 180.5486 (2.4021) 302.0488 (6.7523)	2.0921 (0.9338) 1.3059 (0.1362) 3.8559 (0.5422) 99.6555 (14.3327) 63.0351 (6.7923) 44.2601 (1.2242) 2.3708 (0.7181) 1.0941 (0.3829) 0.7691 (0.0814) 2.5303 (1.0144) 81.3095 (4.3530) 37.2592 (5.1038) 43.6364 (1.5572) 2.5667 (1.5456) 535.1390 (215.4460) 472.1400 (85.2654) 425.1200 (45.0808) 1356.7668 (118.0702) 1350.9520 (188.2830) 550.7750 (25.6849) 446.7374 (9.2086) 929.5620 (47.6606) 885.3490 (27.8584) 747.9998 (6.5119) 1455.8063 (51.6707)	1.4182 (0.2737) 1.1909 (0.0607) 1.7800 (0.2434) 15.4490 (1.7975) 22.1338 (2.2755) 7.4614 (1.5881) 1.6514 (0.2616) 3.1580 (0.2162) 1.4892 (0.0454) 2.8544 (0.2497) 12.6523 (0.3463) 20.9121 (1.2496) 10.0926 (2.3311) 3.7509 (0.1950) 262.9090 (12.6078) 241.3947 (2.4395) 260.6220 (11.7347) 645.0634 (5.2036) 464.4328 (25.0852) 267.6061 (4.5646) 260.3773 (9.0958) 330.8905 (32.7377) 295.1917 (9.6539) 354.8147 (12.1219) 497.3040 (8.2818)	0.6936 (0.0346) 0.7310 (0.0126) 0.7299 (0.0272) 0.7372 (0.0764) 0.7331 (0.0463) 0.8838 (0.0351) 0.7387 (0.0198) 0.6226 (0.0071) 0.6170 (0.0063) 0.6624 (0.0135) 0.6942 (0.0217) 0.5090 (0.0686) 0.7665 (0.0737) 0.6684 (0.0446) 0.7973 (0.0270) 0.8229 (0.0134) 0.5583 (0.3128) 0.7328 (0.0290) 0.7841 (0.0390) 0.7841 (0.0399) 0.8282 (0.0181) 0.8036 (0.0038) 0.6988 (0.0662) 0.8123 (0.0617) 0.6887 (0.0174) 0.7156 (0.0346)

Experiments on another problem with larger aleatoric uncertainty Here we conduct the same experiment as in Section 4.2 for the second problem of Table 4, and show the results in Figure 17. The conclusions are similar to what is observed in Figure 2, demonstrating that the proposed cooperative learning strategy is also effective for problems with large data uncertainty. Unsurprisingly, all approaches are less accurate in this problem for the same amount of training data when comparing

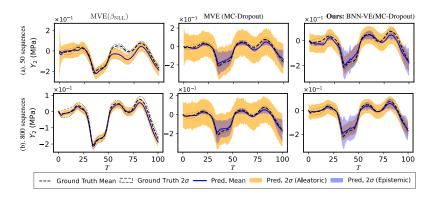


Figure 16: Predictions of MVE (β -NLL), MVE (MC-Dropout), and VeBNN (MC-Dropout) methods on plasticity law discovery dataset. The upper and bottom rows correspond to 50 and 800 training sequences, respectively.

with Figure 2 because this problem has larger data uncertainty. We also see that the proposed cooperative learning strategy consistently has the best performance across all accuracy metrics. Note, however, that the proposed cooperative learning strategy is particularly effective in estimating the mean when compared to all other methods, especially ME (MSE). This demonstrates that the performance of other methods degrades with increasing data uncertainty, unlike the proposed method.

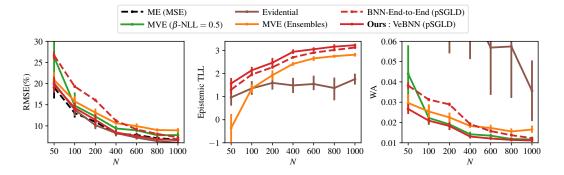


Figure 17: Accuracy metrics (RMSE \downarrow , TLL \uparrow , and WA \downarrow) obtained for the second problem of the plasticity law discovery dataset. The Wasserstein distance (WA) represents the closeness of the estimated aleatoric uncertainty distribution to the ground truth distribution. All metrics result from repeating the training of each method 5 times by randomly resampling points from the training datasets.

D.5 ABLATION STUDY ON ITERATION K

We claim in Section 5 that the iteration K is not a hyperparameter but a parameter, where we present an experiment conducted based on the plasticity law discovery dataset with 800 training sequences. The results are shown in Figure 18.

According to Figure 18, the proposed method exhibits rapid convergence, typically requiring only a single iteration to achieve stable performance, highlighting the effectiveness of the cooperative training strategy. Specifically, the initial warm-up of the mean network provides a favorable starting point for both **Step 2** (aleatoric modeling) and **Step 3** (Bayesian training), facilitating accurate uncertainty decomposition in subsequent stages. However, as observed in the case of seed 1, an unfavorable initialization may lead the warm-up stage to overfit noise. In such cases, additional iterations can effectively mitigate the issue and restore model performance. Based on empirical evidence, convergence occurs within 2 iterations.

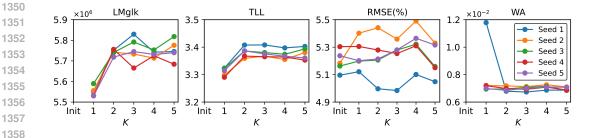


Figure 18: Trajectory of all essential components with respect to the iteration K for 800 training sequences in the plasticity law discovery problem. In the figures, "Init" represents the initialization of training the mean network only as described in Section 3.2, and different curves are realizations under different seeds to restart the procedure, as well as using new samples of training sequences in the dataset.

D.6 ABLATION STUDY ON THE SIZE OF VARIANCE ESTIMATION NETWORK

In Section 5, we also comment on the impact of the size of the variance network, and we provide detailed results of varying the variance of the architecture configurations in Figure 19. As observed, the predictions do not change significantly when choosing different network architectures.

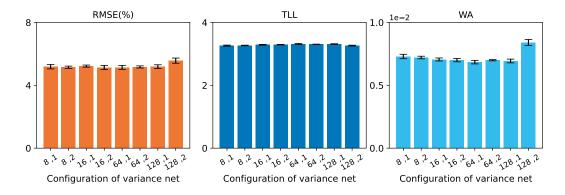


Figure 19: Barplot of all performance metrics with different variance network configurations under 800 training points for plasticity law discovery problem.

E HYPERPARAMETER SETTINGS

E.1 ILLUSTRATIVE DATASETS

Heteroscedastic noise In Appendix D.1, we consider the following one-dimensional example (Skafte et al., 2019):

$$y = x\sin x + 0.3 \cdot x \cdot \varepsilon_1 + 0.3 \cdot \varepsilon_2 \tag{36}$$

where $\varepsilon_1, \varepsilon_2 \sim \mathcal{N}(0, 1)$. In the experiments of Appendix D.1 and Appendix D.1, we sample points uniformly from [0, 10] for training. We generate 1000 points in [0, 10] and 1000 points $[-4, 0] \cup [10, 14]$ to test the performance of different methods.

We depict the results of the illustrative example in Figure 1, Figure 9, and Figure 6, for which the hyperparameters are summarized as follows.

- Architectures: We use a two-layer multi-layer perception (MLP) with 256 neurons for the methods that only require one neural network, except BNN (BBB), where a one-layer with 50 neurons is adopted ³. *Tanh* function is used as the activation function. We note that only the mean is outputted for the ME network, but there are two outputs, mean and aleatoric variance, for the MVE network. For the proposed cooperative learning strategy, we employ the same architecture for the BNN as the ME network. An additional one-layer MLP with 5 neurons is employed as the variance neural network to learn data uncertainty.
- Optimizer/Inference: We use *Adam* optimizer with a learning rate of 0.001 for 20000 epochs to optimize Equation (2) for all deterministic methods and MC-Dropout (Dropout rate is 0.1 for each layer). *Adam* is also used to optimize the ELBO of BBB for 10000 epochs with a learning rate of 0.01. As for pSGLD, we set the burn-in epoch to be 10000 and collect 100 posterior samples every 100 epochs. We also optimize the variance network for 5000 epochs with a learning rate of 0.001 and early stopping of 100 epochs for the proposed cooperative learning strategy.
- Hyperparameter selection: We select 70% data points for training, and the remaining data is used for finding the best hyperparameters, namely number of training epochs and β in the case based on the NLL loss value. After identifying the best number of epochs, we use all data points to re-train the model under the best hyperparameters. For the proposed cooperative learning strategy, we feed all the data to Algorithm 1 and set the iteration K=5.

Homoscedastic noise The homoscedastic noise case has the same ground truth, but instead of input-dependent noise, we consider constant noise, expressed as:

$$y = x\sin x + 0.5 \cdot \varepsilon_2 \tag{37}$$

where $\varepsilon_2 \sim \mathcal{N}(0,1)$.

E.2 UCI REGRESSION DATASETS

We carefully reviewed the experiments and hyperparameters that were found in other studies using UCI regression datasets (Skafte et al., 2019; Immer et al., 2023; Seitzer et al., 2022). We used similar configurations to these studies. Our dataset splits and randomizations can be found in our code for reproducibility because the UCI dataset results can be sensitive to data splits (Seitzer et al., 2022). We also considered multiple experiments per dataset to minimize the impact of randomization. The dataset size and input/output dimensions of each problem are listed in Table 2 and Table 1 under each column within parentheses. We report the metrics at the original scale and averaged over all outputs.

• **Architectures:** We use a one-layer MLP with 50 neurons followed by *ReLU* activation function for all ME and MVE methods. For the proposed cooperative learning strategy, we employ the same architecture for mean net and BNN, and we additionally use an MLP with 5 neurons followed by *ReLU* activation function as the variance network.

³We also try the same architecture with other approaches; however, BBB has difficulty with such a large MLP architecture

- Optimizer/Inference: The methods, including ME, MVE, BBB, as well as the warm-up of the proposed cooperative learning strategy, are trained via *Adam* for 20000 epochs with a learning rate of 0.001 (0.01 for training ELBO), in which the MC-Dropout has a Dropout rate of 0.1 for each layer. As for the pSGLD, we set the burn-in epoch to be 5000 and collected 150 posterior samples every 100 epochs (In total 20000 epochs, which is the same as that of *Adam*). For the additional variance network of the proposed cooperative learning strategy, we use *Adam* to train for 10000 epochs with an early stopping of 100 epochs. The batch size is set to be 256 for all approaches.
- Hyperparameter selection: We split each dataset into train-test by 80%-20% randomly 20 times. In addition, the train set is further divided into 80%-20% for training and validation. We search for the best learning rate within $\left\{10^{-4}, 3 \cdot 10^{-4}, 10^{-3}, 3 \cdot 10^{-3}, 7 \cdot 10^{-4}\right\}$ and as well as record the best epoch utilizing the validation NLL loss. After finding the best learning rate and epoch, the final model is trained using all training data points, and metrics are calculated for the test set. It is noted that the presented results of MVE (β -NLL) in Table 2 and Table 1 are the overall best results among $\{\beta=0.0, \beta=0.25, \beta=0.5, \beta=0.75, \beta=1.0\}$. We also conducted a grid search for BNN-Homo methods for suitable constant noise, where the space is set between zero and one, with 10 different values given the consideration of similar computational resources. For the proposed cooperative learning strategy, we feed all the data to Algorithm 1 and set the iteration K=2.

E.3 IMAGE REGRESSION DATASETS

We take the image regression dataset introduced in (Gustafsson et al., 2023), which consists of relatively large-scale problems, with each containing between 6,592 and 20,614 training images depending on the specific task. Each image has a resolution of 64×64 and the target is a one-dimensional output y. Full details can be found in (Gustafsson et al., 2023); a summary of the relevant information is provided below:

- Cells The dataset contains 10000, 2000, and 10000 images for training, validation, and testing, respectively. The labels have a range of [0, 200], and there is no distribution shift among all the datasets.
- **Cells-Tail** The training and validation datasets contain images with labels in the range of [50, 150]; the test dataset has labels with a range of [0, 200]
- Cells-Gap The training and validation datasets contain images with labels in the range of $[0, 50] \cup [150, 200]$; the test dataset has labels with a range of [0, 200]
- ChairAngle The dataset 17640, 4410, and 11225 images for training, validation, and testing, respectively. The labels have a range of [0.1°, 89.9°], there is no distribution shift among all the datasets.
- Chair Angle-Tail The training/validation datasets contain images whose labels have a range of [15°, 75°]; and the test labels have a range of [0.1°, 89.9°].
- ChairAngle-Gap The labels have a range of [0.1°, 30°] ∪ [60°, 89.9°], and the test labels are in [0.1°, 89.9°].
- SkinLesion The dataset contains 6592, 1164, and 2259 images for training, validation, and testing, respectively. The dataset contains four different sub-datasets, in which the first three are split into train/val with 85%/15%; and the fourth sub-dataset is used as the test dataset.
- AreaBuilding The dataset contains 180 large aerial images with corresponding building segmentation masks. Specifically, the train/val is obtained from two densely populated American cities cities while the test dataset is from rural European cities. Overall, it contains 11184, 2797, and 3890 images for training, validation, and testing, respectively.

We leverage the hyperparameters in (Gustafsson et al., 2023) and define ours as follows:

 Architectures: ResNet34 backbone (He et al., 2016) is employed for this problem. We use a two-layer MLP to decode the prediction into a Gaussian distribution for MVE (β-NLL), MVE (Ensembles), and MVE (MC-Droout). It is noted that a dropout layer with a dropout rate of 0.1 is followed for each MLP layer. As for the variance net and deep evidential regression, the decoding layer is set to be the corresponding outputs after the ResNet34 backbone.

• Optimizer/Inference We use *Adam* with a learning rate of 0.001 for MVE, as well as the warm-up step of the proposed cooperative learning strategy, and employ *Adam* to run for 75 epochs for the above methods with a batch size of 32. As for the pSGLD, we set the burn-in epoch to be 20 and we sample 10 posterior samples every 2 epochs (in total 40 epochs). As for the additional variance network, which is trained with *Adam* for 20 epochs.

E.4 PLASTICITY LAW DATASETS

 Hyperparameter setting for Section 4 In the literature of data-driven constitutive laws, several studies address similar problems without considering noise (Anonymous, 2025). We leverage their setups and define the hyperparameter setting as follows:

- Architectures: For the mean network and BNN, we adopt a two-layer GRU architecture with 128 hidden neurons. It is noted that we only apply the Dropout operation to the hidden-to-decoding layer with a Dropout rate of 0.02. The reason is that this inference method does not show compatible performance when making all weights and biases the Dropout layer. For the proposed cooperative learning strategy, a smaller GRU network with two layers and 8 hidden neurons is employed for the variance network.
- Optimizer/Inference We use *Adam* with a learning rate of 0.001 for ME, MVE, as well as the warm-up step of the proposed cooperative learning strategy and employ *Adam* to run for 2000 epochs for the above methods. As for the pSGLD, we set the burn-in epoch to be 500 and we sample 150 posterior samples every 10 epochs (in total 2000 epochs). As for the additional variance network, which is trained with *Adam* for 4000 epochs with an early stopping patience of 50 epochs.
- Hyperparameter selection: For every experiment of different training points we reserve 100 validation data points from the training dataset to determine the best epoch for ME and MVE. Subsequently, we combine the validation points with the training set and retrain the final model using the best epoch configuration. For the same reason, the results depicted of MVE (β -NLL) is the overall best among $\{\beta=0.0,\beta=0.25,\beta=0.5,\beta=0.75,\beta=1.0\}$. As for the BNN-Homo, we follow the same strategy to execute a grid search for the best homoscedastic noise, where the noise variance is set to be $\{0.04,0.06,0.08,0.10,0.12\}$ empirically. For the proposed cooperative learning strategy, we set the iteration K=2

Hyperparameter setting for Appendix D.6 To investigate the influence of the variance network architecture, we consider eight configurations where the number of layers varies from 1 to 2, and the number of hidden neurons is set to $\{8, 16, 64, 128\}$. The largest configuration, $\{128, 2\}$, is identical to the mean network. All other hyperparameters are consistent with those used in the previous section.

E.5 COMPUTATIONAL RESOURCES

Each experiment, except the image regression datasets, is conducted on a node of an HPC cluster platform with an Intel® Xeon(R) E5-2643v3 CPU with 6 cores of 3.40GHz and 128 GB of RAM. Regarding the image regression datasets, the experiment is conducted based on a platform with an H100 NVL GPU.

F DESCRIPTION OF PLASTICITY LAW DATASET

The fundamental mechanical law of materials is called a constitutive law. It relates average material deformations to average material stresses at any point in a structure. Constitutive laws can model different Physics behaviors, such as elasticity, hyperelasticity, plasticity, and damage. In this paper, we focus on generating datasets for plastically deforming composite materials, coming from prior work (Anonymous, 2025). Without loss of generality, the constitutive law of such path-dependent materials can be formulated as follows:

$$\mathbf{y} = f(\mathbf{x}, \dot{\mathbf{x}}, \tau, \dot{\tau}, \mathbf{h}) \tag{38}$$

where y, x, τ are stress, strain, and temperature respectively, h is a set of internal variables. The constitutive law can be predicted by micro-scale simulations of material domains that are called stochastic volume elements (SVEs) – see Figure 20. These SVEs are simulated by rigorous Physics simulators based on the Finite Element Method (FEM). Each material SVE is utilized as the basic simulation unit. Many factors bring data uncertainty into the data generation process; we focus on data uncertainties from two aspects: (1) SVE size; and (2) particle distribution. As we randomize particle distribution, the stress obtained for an input deformation exhibits stochasticity (aleatoric uncertainty). Therefore, two datasets are generated from simulations according to Table 4.

Table 4: Parameter configuration material plasticity law simulations (Units: SI(mm)).

Name	Microstructure Parameters		Parameters	Hardening Law	$E_{ m fiber}$	Size	E_{matrix}	$\nu_{ m matrix}$	$ u_{ m fiber}$
	v_f	r	$r_{ m std}$	C					
Material 1	0.30	0.003	0.0	$\sigma_y = 0.5 + 0.5(\bar{\epsilon})^{0.4}$	1	0.048	100	0.30	0.19
Material 2	0.30	0.003	0.0	$\sigma_y = 0.5 + 0.5(\bar{\epsilon})^{0.4}$	1	0.030	100	0.30	0.19

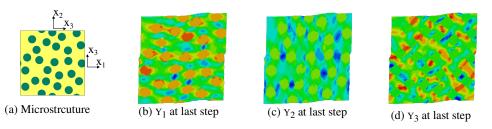


Figure 20: **Material plasticity law simulation illustration.** Figure. (a) shows an arbitrary realization of material microstructure, and the following figures show the contour plot of this material simulation at the final step.

According to Table 4, we have two materials that have different SVE sizes that control the noise sources of the data. Materials with a smaller SVE size have larger data noise. Figure 20 and Figure 21 illustrate details of the simulations and how the uncertainty in the data originates. Specifically, according to the microstructure configuration in Table 4, we generate SVEs using the Monte Carlo Sampling strategy (Melro et al., 2008) and simulate the stress responses through the commercial software ABAQUS (Dassault Systèmes, 2024) with the input of the strain sequence shown in the first row of Figure 21. After simulation, we get a series of contours of stress components in Figure 20 and average the field (color contours) for each input sequence point, leading to the output sequence. An example of 3 realizations of SVEs and corresponding 3 output response sequences, shown as a dashed line in the second row of Figure 21. Each realization corresponds to a randomization of the microstructure of the material (particle distribution). By running multiple realizations, we can obtain the statistics of those strain inputs. By definition, the variation in the stress output is the uncertainty of the data.

Each input deformation sequence and output stress sequence used in training contains 100 points. In total, we use 50 different SVEs to ensure that we have enough realizations to calculate the ground truth aleatoric uncertainty. Overall, we simulate 1000 sequences for training and 100 sequences for testing, respectively, for each problem shown in Table 4.

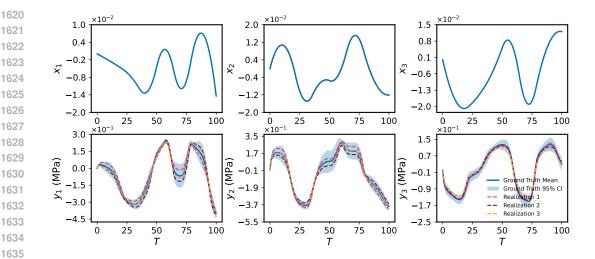


Figure 21: Plasticity law data illustration. The first column is the strain inputs for the material law simulation, and the second row is the stress outputs. Each dashed line represents one particle material microstructure realization in Figure 20; the mean and the confidence interval are obtained via multiple realizations.

The new dataset is made available as open-source in the hope of creating a more interesting problem for assessing future methods because we had difficulties in finding more challenging heteroscedastic problems with ground truth aleatoric uncertainty to assess our method. This dataset is threedimensional and history-dependent, i.e., $\mathcal{D} = \{\mathbf{x}_{n,t}, \mathbf{y}_{n,t}\}$ with features $\mathbf{x}_{n,t} \in \mathbb{R}^3$ and targets $\mathbf{y}_{n,t} \in \mathbb{R}^3$, where n = 1, ..., N are the training sequences (deformation paths) and t = 1, ..., T are the points in each sequence. We highlight two aspects about this dataset. First, the targets y are history-dependent, so estimating a new state y' requires to know the sequence of states needed to reach that state, i.e., regression requires recurrent neural network architectures (Anonymous, 2025), specifically adopting a Gated Recurrent Unit (GRU) architecture (Cho, 2014; Gan et al., 2017) in this work. Furthermore, the dataset was created synthetically by physically-accurate computer simulations of materials under mechanical deformation, so it was possible to generate enough data to determine the ground-truth aleatoric uncertainty (arising from variations within the material). In other words, we have a good estimate of the heteroscedastic noise in the data.