

# HybridFlow: Quantification of Aleatoric and Epistemic Uncertainty with a Single Hybrid Model

Anonymous authors

Paper under double-blind review

## Abstract

Uncertainty quantification is critical for ensuring robustness in high-stakes machine learning applications. We introduce HybridFlow, a novel hybrid architecture that integrates normalizing flows for modeling aleatoric uncertainty with a probabilistic predictor model with the ability to quantify epistemic uncertainty, providing precise uncertainty estimates that are easily integrated into existing model architectures without sacrificing predictive performance. HybridFlow improves upon previous uncertainty quantification frameworks across a range of regression tasks, such as depth estimation, a collection of regression benchmarks, and a scientific case study of ice sheet emulation. We also provide an analysis of the quantified uncertainty, showing that the uncertainty quantified by HybridFlow is calibrated and better aligns with model error than existing methods for quantifying aleatoric and epistemic uncertainty. HybridFlow addresses a key challenge in Bayesian deep learning, unifying aleatoric and epistemic uncertainty modeling in a single robust framework.

## 1 Introduction

Uncertainty quantification plays a critical role in modern machine learning, particularly in high-stakes applications such as autonomous driving (Grewal et al., 2024), medical diagnosis (Jalal et al., 2024), and scientific modeling (Wang et al., 2023). The ability to estimate uncertainty allows models not only to make predictions but also to provide insights into the reliability of those predictions. This capability is especially valuable in scenarios where decision-making depends heavily on understanding the limitations and confidence of a model’s outputs (Thuy & Benoit, 2024; Marusich et al., 2024; Amodei et al., 2016). In recent years, researchers have increasingly focused on improving uncertainty estimation methods in deep learning to achieve more robust and trustworthy systems (Gal & Ghahramani, 2016; Kendall & Gal, 2017; Messuti et al., 2025; Yoon & Kim, 2024; Hwang & Shin, 2024). Despite significant progress, key challenges remain, particularly in achieving accurate predictions while effectively quantifying different types of uncertainties (Psaros et al., 2023).

In the context of machine learning, uncertainty is broadly categorized into two types: *aleatoric uncertainty*, which arises from the inherent noise and variability in data, and *epistemic uncertainty*, which reflects the model’s lack of knowledge, often due to limited or biased training data (Der Kiureghian & Ditlevsen, 2009; Gruber et al., 2023). While aleatoric uncertainty is irreducible, epistemic uncertainty can, in principle, be reduced by acquiring additional data or improving the model architecture. Techniques such as Monte Carlo (MC) Dropout (Gal & Ghahramani, 2016) and ensemble methods (Lakshminarayanan et al., 2017) have been developed to estimate epistemic uncertainty, while approaches like density estimation (Bishop, 1994) and heteroscedastic regression (Kendall & Gal, 2017) aim to model aleatoric uncertainty. Despite advancements in single-model frameworks, they often prioritize state-of-the-art uncertainty quantification at the expense of modularity and flexibility, which are critical for integrating uncertainty estimation into diverse existing architectures. For example, methods such as heteroscedastic regression (Kendall & Gal, 2017) are widely adopted due to their simplicity and ease of integration, even though they often suffer from calibration issues and suboptimal predictive accuracy (Seitzer et al., 2022). Heteroscedastic regression methods remain popular because they can be applied to diverse predictive models without significant architectural changes.

This paper introduces a novel hybrid flow-based architecture, HybridFlow, which addresses these challenges by quantifying both aleatoric and epistemic uncertainty without compromising predictive performance. This framework integrates a normalizing flow (NF) (Rezende & Mohamed, 2015) to model aleatoric uncertainty with any probabilistic predictor model for epistemic uncertainty estimation and model prediction. HybridFlow enables precise and reliable uncertainty quantification through leveraging both input data and the latent space generated by the NF as model inputs. By decoupling the aleatoric uncertainty quantification mechanism from the predictor model itself, the proposed HybridFlow framework is inherently modular, serving as a flexible method that can be adapted to a wide range of applications. While the current implementation demonstrates its efficacy using specific benchmarks and case studies, the design of the framework allows for the substitution of different predictors tailored to other tasks or domains.

We evaluate the HybridFlow model framework against three widely adopted frameworks (Kendall & Gal, 2017; Seitzer et al., 2022; Caprio et al., 2024) for quantifying both aleatoric and epistemic uncertainty in a single model, and demonstrate that our proposed framework provides accurate predictions and quantified uncertainty for a range of tasks: depth estimation in computer vision and a series of regression benchmarks. We also include a scientific case study, where we test the HybridFlow architecture as an emulator for future sea level rise. Another key contribution of this work is that we provide a comprehensive evaluation of quantified uncertainty, which is often overshadowed by the analysis of prediction accuracy alone (Wang et al., 2025). We employ a diverse set of metrics that collectively assess different aspects of uncertainty estimation to compare the quantified uncertainty from various methods. The proposed framework showcases a direct improvement in both predictive accuracy and uncertainty quantification over existing methods, providing a compelling step toward single, flexible models that can effectively model both aleatoric and epistemic uncertainties while maintaining high predictive accuracy.

## 2 Related Work

Uncertainty quantification has received significant attention in Bayesian deep learning due to its importance in improving model reliability and interpretability. Uncertainty quantification techniques can be model-agnostic, meaning that they can be incorporated into a variety of predictive architectures, or model-specific, where specific constraints are placed on the predictive architecture in order to quantify uncertainty. Kendall & Gal (2017) laid the foundation for model-agnostic uncertainty quantification by modeling both aleatoric and epistemic uncertainty within a single framework, specifically in the context of computer vision, by using a combination of existing Bayesian approaches and a heteroscedastic log likelihood loss for modeling aleatoric uncertainty. Their approach adds a prediction layer and variance layer at the end of a predictor model, such that the model predicts both the mean  $\mu(x)$  and variance  $\sigma^2(x)$  for a given input  $x$ , with the assumption that the uncertainty follows a Gaussian distribution. To find the optimal weights, maximum likelihood estimation is used, which is equivalent to minimizing the negative log-likelihood (NLL), or heteroscedastic loss, of the predictive distribution:

$$\mathcal{L} = \frac{1}{N} \sum_{i=1}^N \frac{1}{2\sigma^2(\mathbf{x}_i)} \|y_i - \mu(\mathbf{x}_i)\|^2 + \frac{1}{2} \log \sigma^2(\mathbf{x}_i) \quad (1)$$

This formulation balances the mean squared error (MSE) with the uncertainty estimates generated by the model. For the predictor model, Kendall & Gal (2017) integrate MC dropout (Gal & Ghahramani, 2016) to estimate epistemic uncertainty, and demonstrate the effectiveness of heteroscedastic regression in tasks such as semantic segmentation and depth regression. This approach enables flexibility, allowing for uncertainty quantification with minimal adjustment to existing predictor model architectures (Smith et al., 2024). However, this method is prone to calibration issues, particularly under distributional shifts or when noise patterns deviate from Gaussian assumptions. Additionally, reliance on a single combined loss function often leads to overestimation of one type of uncertainty at the expense of the other.

Recent research in model-agnostic uncertainty quantification techniques has highlighted the challenges and pitfalls of heteroscedastic uncertainty estimation. Specifically, Seitzer et al. (2022) critiqued this standard approach to modeling aleatoric uncertainty, emphasizing the tendency of NLL-trained probabilistic neural networks to underestimate uncertainties in the presence of out-of-distribution data or highly skewed noise distributions. They show that the use of the NLL loss is appropriate if the difference between the model

prediction and the true value is solely caused by noise, or aleatoric uncertainty. However, in practice, the inability of the model to perfectly approximate the phenomena may also contribute to error, which leads to an overestimation of aleatoric uncertainty when the data is poorly predicted by the model (Seitzer et al., 2022). To address this, they introduce the Beta-NLL loss, a refinement of the NLL loss that adjusts the balance between MSE and uncertainty estimates to address sub-optimal fits and training instability. This method represents an effort to mitigate calibration challenges while retaining the simplicity of the heteroscedastic regression framework. While in some tasks it does improve predictive performance, it can be unstable and fails to address the miscalibration of aleatoric uncertainty due to joint loss function training. Others have attempted to address miscalibration and the decrease in predictive accuracy using differing probabilistic predictors or post-hoc methods, but they similarly encounter comparable trade-offs between predictive accuracy and uncertainty quantification capabilities (Valdenegro-Toro & Mori, 2022; Yang & Li, 2023).

Model-specific uncertainty quantification methods, or those that require specific predictive architectures, have been successful in addressing miscalibration and reduced accuracy when quantifying both aleatoric and epistemic uncertainty, but at the cost of remaining generalizable to a wide range of predictive architectures. Credal Bayesian Deep Learning (CBDL) (Caprio et al., 2024) uses a credal set of Bayesian neural networks (BNNs) with a prescribed set of priors to estimate aleatoric and epistemic uncertainty. While BNNs are theoretically applicable to a wide range of tasks, in practice they are computationally intensive to implement and train, often requiring specialized inference procedures and long training times. Similarly, Evidential deep learning methods (Sensoy et al., 2018) learn the parameters of higher-order distributions directly from the data, but require specialized architectures and loss functions that restrict their usability. Berry & Meger (2023) propose NFlows Base and NFlows Out, two NF (see Section 3.1) ensemble methods that use fixed dropout masks to efficiently model uncertainty. By ensembling the NFs, they capture epistemic uncertainty, while the individual flows parameterize  $p(y|x)$  to model aleatoric uncertainty. This method is model-specific, as NFs are required for the prediction task, rather than solely distribution approximation. Similarly, FlowNet (Zhang et al., 2024) leverages NFs to parameterize  $p(y|x)$ , and uses predicted parameters to capture the aleatoric and epistemic uncertainties. Chan et al. (2024) introduce HyperDM, which employs a Bayesian hyper-network to generate an ensemble of model weights and combines it with a conditional diffusion model to estimate predictive distributions. These methods achieve state-of-the-art uncertainty quantification but require specific model architectures, limiting their practicality for users aiming to enhance existing models or create models for specific tasks and domains where the predictor proposed in these studies may be unfit.

Unlike the aforementioned methods, which integrate uncertainty quantification deeply into the custom architectures, methods such as our proposed HybridFlow framework and heteroscedastic regression (Kendall & Gal, 2017; Seitzer et al., 2022) are flexible by design, allowing integration of uncertainty quantification into any type of probabilistic predictive model. The importance of this capability is evident in the wide-spread adoption of heteroscedastic regression methods (Roohani et al., 2024; Huang et al., 2023; Yelleni et al., 2024), over custom models with complex, nonflexible model architectures. In this work, we introduce HybridFlow, a modular hybrid model architecture that surpasses prior flexible heteroscedastic regression frameworks by providing accurate, calibrated prediction and uncertainty estimates while avoiding combined loss function training and generalizing beyond Gaussian assumptions. This enables users to incorporate uncertainty modeling without being tied to specific predictor architectures or objectives. This approach aligns with the practical needs of many users who primarily aim to enhance existing predictors with uncertainty quantification, rather than adopting entirely new architectures that are not specific to applications or domains. HybridFlow is able to provide a balance between advanced uncertainty quantification and compatibility with diverse predictive frameworks.

HybridFlow remains modular and flexible by using a hybrid model architecture, as introduced in Nalisnick et al. (2019). Nalisnick et al. (2019) integrate a generative model with a predictive model, enabling exact computation of the joint distribution  $p(x, y)$ . Their work demonstrates a theoretically principled integration of tasks but does not explicitly disentangle aleatoric and epistemic uncertainty, and relies on a non-probabilistic generalized linear model for prediction, which limits the expressiveness of the predictive component. HybridFlow builds on Nalisnick et al. (2019) by explicitly using the generative component to model aleatoric uncertainty and by introducing flexibility in the choice of predictive architectures.

Unlike previous methods that rely on Gaussian assumptions or joint loss functions, HybridFlow uses conditional NFs to model  $p(y|x)$  and calculate aleatoric uncertainty through the variance of sampled predictions. Simultaneously, epistemic uncertainty is quantified using a probabilistic predictor model, decoupled from the aleatoric uncertainty module. This design avoids the trade-offs and calibration issues (Kendall & Gal, 2017; Seitzer et al., 2022; Valdenegro-Toro & Mori, 2022) inherent in joint loss functions while maintaining predictive accuracy. By employing a flow-based framework, HybridFlow enhances scalability to high-dimensional data and improves robustness to out-of-distribution samples. Its modularity enables straightforward adaptation to diverse tasks without the need for extensive domain-specific tuning, making it a versatile framework for robust uncertainty quantification.

### 3 Method

HybridFlow builds upon recent advancements in uncertainty quantification and generative modeling to address existing challenges in modeling both aleatoric and epistemic uncertainty in a single model. It leverages NFs (Rezende & Mohamed, 2015; Papamakarios et al., 2017) for modeling aleatoric uncertainty while integrating a probabilistic prediction model into a hybrid architecture (Nalisnick et al., 2019) for both accurate predictions and epistemic uncertainty quantification. This section describes the components of HybridFlow and their integration into a unified framework for uncertainty quantification.

#### 3.1 Normalizing Flows for Aleatoric Uncertainty

Normalizing Flows are a class of generative models that transform a simple base distribution  $p_z(z)$ , typically a standard Gaussian, into a more complex target distribution  $p(y|x)$  through a series of invertible transformations  $f$ . In the case of conditional NFs, the transformation is parameterized based on the conditioning variable  $x$ :

$$y = f(z|x), \quad z \sim p_z(z), \quad (2)$$

where  $f$  is designed to be invertible, ensuring that the probability density function of  $y$  given  $x$  can be computed via the change-of-variables formula:

$$p(y|x) = p_z(f^{-1}(y|x)) \left| \det \frac{\partial f^{-1}(y|x)}{\partial y} \right|, \quad (3)$$

where  $\frac{\partial f^{-1}(y|x)}{\partial y}$  represents the Jacobian of the transformation  $f^{-1}$  conditioned on  $x$ . The Jacobian matrix captures how the volume of the transformed space changes under the inverse mapping  $f^{-1}$ . The absolute value of its determinant accounts for this local volume change and ensures the resulting density  $p(y|x)$  remains normalized.

In HybridFlow, a *Conditional Masked Autoregressive Flow (CMAF)* (Papamakarios et al., 2017) is used. The CMAF decomposes the conditional joint distribution  $p(y|x)$  into a product of autoregressive conditionals:

$$p(y|x) = \prod_{i=1}^d p(y_i | y_{<i}, x), \quad (4)$$

where  $d$  is the dimensionality of  $y$ , and each conditional  $p(y_i | y_{<i}, x)$  is parameterized by a neural network. This conditional structure allows the flow to explicitly learn the distribution of  $y$  given  $x$ , capturing both the inherent noise in  $y$  and its dependence on  $x$ . Learning  $p(y|x)$  allows for the calculation of aleatoric uncertainty (Section 3.3), and the invertibility of the flow  $f$  enables us to produce a latent representation of the inputs,  $z = f^{-1}(y|x)$ , which serves as an additional feature for the predictor model (Section 3.2).

While normalizing flows are highly effective generative models for capturing complex conditional distributions, they are not typically used as predictors (Kobyzev et al., 2020). Their primary objective is to model the full data distribution  $p(y|x)$  rather than to optimize predictive accuracy for a specific target variable. In contrast, predictive models are explicitly trained to minimize prediction error and can incorporate mechanisms for estimating epistemic uncertainty. By introducing a dedicated predictor, HybridFlow leverages the

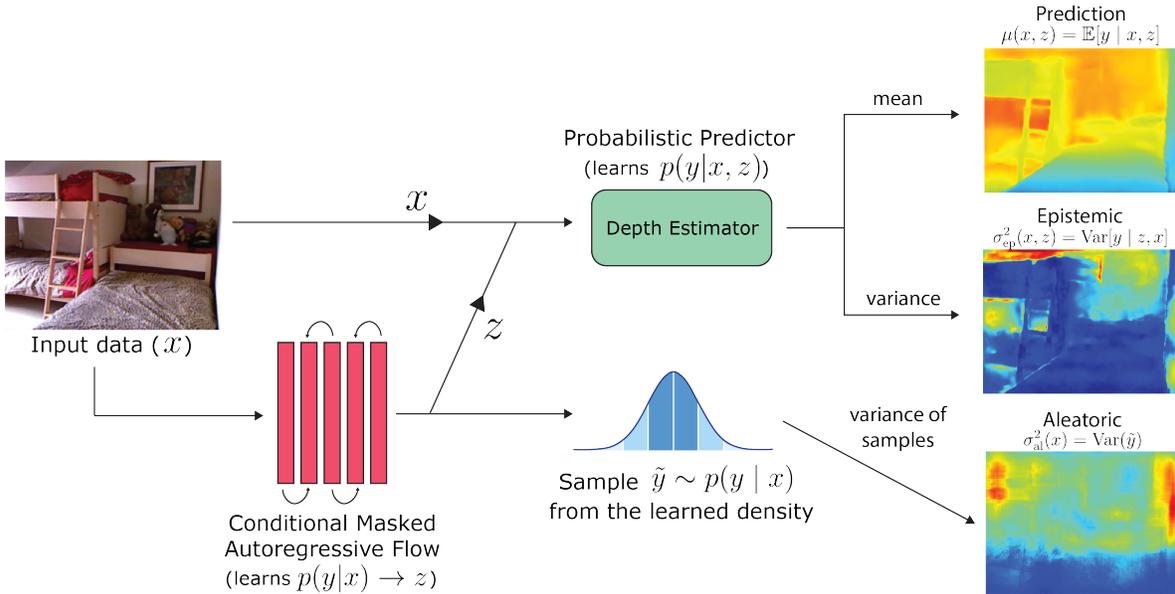


Figure 1: HybridFlow architecture, with depth estimation with the NYU Depth v2 dataset as the example task. A Conditional Masked Autoregressive Flow (CMAF) learns the density of input data ( $x$ ) through a latent representation ( $z$ ). The latent representation and the input data are concatenated and used as inputs to a probabilistic predictor model. From the predictor model, the prediction and the epistemic uncertainty can be calculated. The aleatoric uncertainty is estimated by calculating the variance of samples generated from the learned data distribution by the CMAF.

strengths of both components: the NF captures the uncertainty inherent in the data, while the predictor focuses on producing accurate point estimates and quantifying model-based uncertainty.

### 3.2 Hybrid Architecture

To achieve both predictive accuracy and robust uncertainty quantification, the *latent space* representation  $z$  generated by the conditional NF is passed, along with the original inputs  $x$ , into a probabilistic predictor.

The probabilistic predictor models the distribution  $p(y | z, x)$ , where  $y$  represents the prediction target. By combining  $z$  (the latent encoding capturing data-specific uncertainty) and  $x$  (the original inputs), the predictor jointly learns the relationship between the inputs and the outputs while retaining uncertainty information from the NF. The design of HybridFlow is flexible, such that any probabilistic predictor can be used. For example, methods such as MC Dropout (Gal & Ghahramani, 2016), Bayesian neural networks, deep ensembles (Lakshminarayanan et al., 2017), and other Bayesian methods are all viable predictor models that can be implemented within the hybrid architecture.

In the HybridFlow framework, the predictive mean  $\mu(x, z) = \mathbb{E}[y | x, z]$  represents the model’s prediction, while the epistemic uncertainty  $\sigma_{ep}^2(x, z) = \text{Var}[y | z, x]$  quantifies the uncertainty due to the model’s limited knowledge. The choice of the probabilistic predictor can be tailored based on the application’s requirements for computational efficiency, model complexity, or interpretability.

### 3.3 Aleatoric Uncertainty Estimation

The aleatoric uncertainty is captured directly from the learned data distribution modeled by the conditional NF. By sampling from the posterior distribution  $p(y | x)$  learned by the CMAF, we define the sampled predictions as  $\tilde{y} \sim p(y | x)$ . The aleatoric uncertainty is then computed as the variance of these samples,  $\sigma_{al}^2(x) = \text{Var}(\tilde{y})$ .

This variance can be considered aleatoric because it arises from the distribution  $p(y|x)$  that the flow explicitly learns from the data. The NF is trained to model the full conditional distribution of possible outcomes given each input, including any inherent noise or ambiguity in the data. When we sample multiple outputs for the same input, the spread of those outputs reflects the learned variability within the data itself. Thus, the variance of these samples naturally corresponds to aleatoric uncertainty, as it captures the irreducible variability in outcomes.

### 3.4 Summary of the Workflow

1. **Data Transformation:** Input data  $x$  is transformed into a latent representation  $z$  via a CMAF-based NF, conditioned on  $x$ .
2. **Prediction and Epistemic Uncertainty:** The latent  $z$  and input  $x$  are fed into a probabilistic predictor that outputs the mean prediction  $\mu(x, z)$  and epistemic uncertainty  $\sigma_{\text{ep}}^2(x, z)$ .
3. **Aleatoric Uncertainty:** Sampling from  $p(y|x)$ , the aleatoric uncertainty  $\sigma_{\text{al}}^2(x)$  is calculated as the variance of the samples.

### 3.5 Training Protocol and Implementation Considerations

HybridFlow is trained in a modular two-stage process that reflects its hybrid architecture. First, the normalizing flow (NF) component is trained independently to learn the conditional distribution  $p(y|x)$ , which is used to quantify aleatoric uncertainty. The NF is optimized using the NLL loss of the conditional distribution (Equation 3), a standard and widely adopted objective for training normalizing flows to model the full data distribution. After training the flow, the latent representation  $z = f^{-1}(y|x)$  is extracted and used as an additional input to the predictor model, along with the original input  $x$ . During this second stage, the predictor model is trained using a task-appropriate loss function, such as mean squared error (MSE), and optionally paired with a Bayesian approximation method to estimate epistemic uncertainty.

The NF can be optionally fine-tuned during the training of the predictor model to allow improved joint representations, though in most applications the flow converges independently and does not require additional updates. Similarly, the predictor model can be fine-tuned after integration with the NF outputs. HybridFlow supports any probabilistic predictor architecture for epistemic uncertainty estimation; in our experiments, we use MC Dropout for simplicity and compatibility with prior work, but the framework can accommodate deep ensembles, variational inference, or other Bayesian or Bayesian-approximation techniques.

A key consideration when implementing HybridFlow is the dimensionality of the target variable  $y$ , which constrains the use of normalizing flows. High-dimensional output spaces increase computational and memory requirements due to the need for tractable Jacobian determinants in the NF. In such cases, the NF can be paired with an autoencoder to first compress the output space into a lower-dimensional latent representation, as we demonstrate in Section 4.1. This preserves the flow’s modeling capacity while ensuring that the method remains tractable for large-scale or high-resolution outputs. This design choice reflects HybridFlow’s flexibility in balancing expressiveness, efficiency, and integration with domain-specific architectures.

## 4 Experiments

This section presents the evaluation of the HybridFlow architecture on various distinct datasets, including the NYU Depth v2 dataset for depth estimation and 12 UCI regression benchmarks. The experiments focus on assessing both predictive accuracy and uncertainty quantification capabilities. We benchmark HybridFlow against a baseline model based on the framework of Kendall & Gal (2017), a similar framework from Seitzer et al. (2022) using the Beta-NLL (BNLL) loss, as well as a CBDL model, a state-of-the-art model-specific uncertainty quantification framework (Caprio et al., 2024).

### 4.1 Depth Estimation

We use the NYU Depth v2 dataset (Silberman et al., 2012), which is widely used for evaluating models in computer vision tasks, and the primary benchmark in Kendall & Gal (2017). It consists of 400,000 high-

dimensional RGB and depth image pairs, covering a variety of environments such as offices, homes, and classrooms. For our experiments, we used a preprocessed subset comprising 795 training images and 654 test images, following the standard splits used in previous studies. Each RGB image is paired with a depth map representing the ground truth distances of each pixel from the camera. The diversity and complexity of the scenes, combined with the inherent measurement noise in depth sensors, make this dataset ideal for evaluating both aleatoric and epistemic uncertainty.

We train three models for quantifying aleatoric and epistemic uncertainty, two baseline heteroscedastic regression models, one with a Gaussian NLL loss (Kendall & Gal, 2017), and another with a Beta-NLL loss (Seitzer et al., 2022), and a hybrid regression model with the proposed HybridFlow architecture. For both heteroscedastic regression models, we include a 2D Convolutional layer at the end of the predictor model to output the mean (prediction) and variance (aleatoric uncertainty) and train with their respective NLL-based losses. For the hybrid regression model, we use a CMAF paired with a UNet-based autoencoder to quantify the aleatoric uncertainty. While the autoencoder is not a necessary component of HybridFlow, we use an autoencoder to compress the input images (which are  $384 \times 512 \times 3$ ) into a latent space of size 256, thus demonstrating the proposed approach is computationally manageable when working with high-dimensional data (for more information and performance of the autoencoder, see Appendix B). We then use the generated latent representation from the CMAF and the RGB image as inputs to the predictor model.

For all experiments, we use a predictor model presented in Ganj et al. (2025), which achieves state-of-the-art performance on the NYU Depth v2 dataset across key metrics such as Absolute Relative Error (AbsRel) and Root Mean Squared Error (RMSE). Following Kendall & Gal (2017), we implement MC Dropout (Gal & Ghahramani, 2016) to make the Ganj et al. (2025) predictor probabilistic for the quantification of epistemic uncertainty. While other epistemic uncertainty quantification methods have been shown to produce more calibrated uncertainty (e.g., Deep Ensembles (Lakshminarayanan et al., 2017), Variational Inference (Blundell et al., 2015)), we choose MC dropout due to the simplicity of incorporating into existing models and to be consistent with the implementations of the baseline model in Kendall & Gal (2017). However, we highlight that any other probabilistic method can be used due to the modularity of the HybridFlow framework. For reference, we have also included a non-probabilistic version of the Ganj et al. (2025) model (without the ability to quantify uncertainty) in Table 1 in order to compare the effect that quantifying uncertainty has on the predictive performance. We also highlight that our implementation results show improvements over the Kendall & Gal (2017), solely because we use a predictor model (Ganj et al., 2025) that showcases current state-of-the-art performance rather than any difference in training or uncertainty quantification protocol.

Despite its strong performance on a variety of uncertainty quantification tasks (Caprio et al., 2024), CBDL cannot be applied in this depth estimation task because the underlying UNet-based predictor for the depth estimation task is not a Bayesian model. As a result, CBDL’s ensemble of variational guides cannot be constructed for a purely deterministic architecture, which limits its flexibility. This incompatibility represents a significant drawback of model-specific uncertainty frameworks when extending to non-Bayesian deep learning modules.

Training was conducted for 100 epochs using a batch size of 16, with the Adam optimizer (Kingma & Ba, 2015), a learning rate of  $1e-4$ , and a weight decay of  $1e-3$ . For the NF components, we used a CMAF with 4 flow layers using `nflows` (Durkan et al., 2020). We train the CMAF flow with a NLL loss prior to training the prediction model. Then, during the training stage of the predictor, we fine-tune the flow with a lower learning rate ( $1e-6$ ) to integrate the latent representation  $z$  into the hybrid architecture.

The probabilistic predictor for HybridFlow was configured with MC Dropout at a rate of 0.2 for epistemic uncertainty quantification, and we performed 30 stochastic forward passes at test time to approximate the predictive mean and epistemic variance. The heteroscedastic regression models minimized a Gaussian NLL and Beta-NLL losses, while the HybridFlow model minimized a combined loss consisting of the flow’s negative log-likelihood, and a scale-invariant log loss (Eigen et al., 2014) and multi-scale gradient loss (Ganj et al., 2025) for the predictor.

All training was conducted on a single Tesla V100 GPU, with the heteroscedastic models requiring approximately 25 hours to converge and HybridFlow requiring 36 hours. The models were evaluated using the

standard metrics for depth estimation, including MSE, AbsRel, and thresholds ( $\delta_{1-3}$ ), which indicate the proportion of predictions within 25%, 56%, and 95% of the ground truth values (Eigen et al., 2014).

In this study, we include a robust analysis of the quantified uncertainty from the approaches tested. Following Wang et al. (2025), we evaluate calibration quality using the Expected Calibration Error (ECE), which measures the average discrepancy between the confidence intervals and predicted accuracy, and the Prediction Interval Coverage Probability (PICP) which reports the fraction of true values contained within the computed uncertainty intervals. To assess the sharpness of our uncertainty estimates, we compute the Mean Prediction Interval Width (MPIW), capturing the typical size of the predictive intervals, and the Winkler score, which penalizes both overly wide and under-covering intervals. Additionally, we apply strictly proper scoring rules (Gneiting & Raftery, 2007), such as the Negative Log Likelihood score (NLL) and the Continuous Ranked Probability Score (CRPS), which penalize miscalibration in the full predictive distribution. Together, these metrics furnish a concise yet comprehensive characterization of aleatoric and epistemic uncertainty across predictive accuracy, calibration, sharpness, and reliability.

## 4.2 UCI Regression Benchmarks

To evaluate the robustness and accuracy of HybridFlow, we compared it against the same heteroscedastic regression methods on the UCI regression datasets, as detailed in Seitzer et al. (2022), and a CBDL model (Caprio et al., 2024). The UCI datasets are a widely recognized benchmark suite for regression tasks, offering diverse data distributions, varying noise levels, and differing scales that are ideal for evaluating models’ predictive accuracy and uncertainty estimation capabilities.

We follow a standard experimental protocol for Bayesian deep learning frameworks (Seitzer et al., 2022; Hernández-Lobato & Adams, 2015; Gal & Ghahramani, 2016), using the same datasets and splits. Each dataset was evaluated 20 times with random 80/20 splits for training and testing, normalizing input features and target outputs to zero mean and unit variance based on the training split. For more experimental details, see Seitzer et al. (2022).

We train a flow using the HybridFlow framework with a CMAF with 5 flow layers. We first train the flow alone using a NLL loss. Then, the latent representations generated by the flow were combined with the original inputs and fed into the predictor model, which incorporates MC Dropout for modeling epistemic uncertainty. The predictor model used for all three methods is a simple MLP with one hidden layer of 50 ReLU-activated units that is trained using an MSE loss and early stopping, as in Seitzer et al. (2022).

HybridFlow is compared to heteroscedastic regression models trained with a Gaussian NLL loss (Kendall & Gal, 2017) and Beta-NLL loss (Seitzer et al., 2022), as well as a non-probabilistic version of the model trained with an MSE loss. We conduct the experiments using a Tesla V100 GPU with an Adam optimizer and perform a grid search of possible learning rates ranging from  $1e-3$  to  $1e-5$ . To quantify epistemic uncertainty, we applied MC Dropout with a dropout probability of 0.2 and performed 50 stochastic forward passes during evaluation.

In addition to the heteroscedastic regression baselines, we also benchmark against CBDL (Caprio et al., 2024), a model-specific uncertainty framework consisting of a credal set of Bayesian neural networks to jointly estimate aleatoric and epistemic uncertainty. CBDL often achieves strong calibration by averaging over the credal set of network, but this approach entails significant computational overhead and requires specialized inference protocols that can be difficult to integrate into existing pipelines. By including CBDL in our comparisons, we aim to illustrate that HybridFlow not only meets or exceeds CBDL’s calibration and predictive performance, but does so with a single, modular framework that is simpler to train and can be used with existing predictive models.

We implement CBDL as a credal set of Bayesian models in Pyro (Bingham et al., 2019) that closely resemble the MLP models used for HybridFlow and heteroscedastic regression baselines. We initialize four BNNs, each with one hidden layer of 50 Tanh activated units, by varying the Gaussian prior scale and likelihood over the parameters and train each of them using stochastic variational inference (SVI). Aleatoric and epistemic uncertainties are computed using individual model entropies and variability within the credal set as outlined in Caprio et al. (2024).

Table 1: Comparison of hybrid depth estimation model with the heteroscedastic depth estimation models using NYU Depth v2. The Ganj et al. (2025) model does not quantify uncertainty, but is included to compare the predictive accuracy changes when adding uncertainty quantification capabilities.

TRAINING METHOD	MSE	ABSREL	$\delta_1$	$\delta_2$	$\delta_3$
NLL (KENDALL & GAL, 2017)	0.129	0.083	0.947	0.987	0.995
B-NLL (SEITZER ET AL., 2022)	0.134	0.087	0.947	0.987	0.994
HYBRIDFLOW	<b>0.058</b>	<b>0.041</b>	<b>0.989</b>	<b>0.998</b>	<b>0.999</b>
GANJ ET AL. (2025) (NO UQ)	0.057	0.039	0.989	0.998	0.999

Table 2: Uncertainty quantification metrics across models on the NYU Depth v2 datasets. Metrics include NLL (Negative Log-Likelihood), calibration error (ECE), sharpness, total PICP (Prediction Interval Coverage Probability), CRPS (Continuous Ranked Probability Score), and MPIW (Mean Prediction Interval Width).

METHOD	NLL			ECE			WINKLER			MPIW			PICP	CRPS
	ALE.	EPI.	TOT.											
NLL	5.79	10.88	7.39	4.42	<b>2.98</b>	3.07	14.53	11.51	10.16	5.38	0.035	4.04	0.97	2.71
BNLL	2.95	3.37	2.73	2.32	10.96	4.63	6.18	10.52	8.23	<b>3.70</b>	0.78	<b>3.60</b>	0.92	1.66
HYBRIDFLOW	<b>1.96</b>	<b>0.09</b>	<b>0.73</b>	<b>0.68</b>	4.14	<b>0.46</b>	<b>4.82</b>	<b>7.46</b>	<b>4.83</b>	4.70	<b>0.64</b>	4.76	<b>0.99</b>	<b>0.37</b>

We evaluate the predictive accuracy of each model using the mean RMSE value and the standard deviation of the 20 variations of the model based on the random splits. We then evaluate the uncertainty quantification performance of each model using a robust selection of metrics that collectively capture critical aspects of uncertainty estimation including the NLL, ECE, CRPS, PICP, and MPIW. NLL and CRPS quantify the accuracy of the probabilistic predictions, with CRPS offering insights into the entire predictive distribution. ECE specifically assesses calibration quality, measuring the alignment between predicted uncertainty intervals and observed accuracy. Additionally, PICP evaluates whether prediction intervals reliably encompass true values, and MPIW assesses the sharpness or informativeness of these intervals. Together, these metrics provide a comprehensive and nuanced evaluation of both predictive performance and the reliability of the quantified uncertainty. We evaluate aleatoric, epistemic, and total uncertainty with all of the aforementioned metrics. Total uncertainty, within the scope of this work, is the sum of the aleatoric and epistemic uncertainty, which is common in uncertainty quantification literature (Depeweg et al., 2018; Hüllermeier & Waegeman, 2021; Wimmer et al., 2023).

## 5 Results

### 5.1 Depth Estimation Results

HybridFlow demonstrates superior predictive accuracy on the NYU Depth v2 dataset, outperforming baseline models that use heteroscedastic regression methods and the CBDL model. As seen in Table 1, the HybridFlow model achieves improved results across several evaluation metrics, including MSE, AbsRel, and  $\delta$  thresholds ( $\delta_1$ - $\delta_3$ ). The Hybrid model achieves an MSE of 0.058 and an AbsRel of 0.041, which are improvements over the Gaussian NLL framework (MSE: 0.129, AbsRel: 0.083) and the Beta-NLL approach (MSE: 0.134, AbsRel: 0.087). Notably, the performance is on par with the Ganj et al. (2025) predictor model when uncertainty quantification is excluded (MSE: 0.057, AbsRel: 0.039), illustrating that the integration of uncertainty estimation does not compromise predictive accuracy.

HybridFlow also exhibits robust uncertainty quantification capabilities, addressing both aleatoric and epistemic uncertainties effectively. The aleatoric uncertainty, modeled using a CMAF in the hybrid architecture, captures inherent noise in the dataset, while epistemic uncertainty is quantified through MC Dropout in the predictor model. As shown in Table 2, HybridFlow achieves consistently better aleatoric ECE and PICP

Table 3: RMSE values for each of the UCI regression datasets from a standard MLP model compared to heteroscedastic training methods (NLL, BNLL) and a non-probabilistic predictor.  $N$ ,  $d$ , and  $k$  represent the dataset length, input dimensions, and output dimensions respectively.

DATASET	$N$	$d$	$k$	NLL	BNLL	CBDL	HYBRIDFLOW	MLP (NO UQ)
BOSTON HOUSING	506	13	1	3.56 ± 1.07	3.42 ± 1.04	3.40 ± 0.74	<b>3.12 ± 0.85</b>	3.24 ± 1.08
CARBON	10721	5	3	0.0068 ± 0.0029	0.0068 ± 0.0029	0.0076 ± 0.0019	<b>0.0068 ± 0.0018</b>	0.0068 ± 0.0028
CONCRETE STRENGTH	1030	8	1	6.08 ± 0.65	5.61 ± 0.65	5.52 ± 0.22	<b>5.07 ± 0.52</b>	4.96 ± 0.64
CYCLE POWER PLANT	9568	4	1	4.06 ± 0.18	4.04 ± 0.15	4.26 ± 0.20	<b>4.02 ± 0.18</b>	4.01 ± 0.19
ENERGY EFFICIENCY	768	8	2	2.25 ± 0.34	1.12 ± 0.25	<b>0.97 ± 0.08</b>	1.00 ± 0.14	0.92 ± 0.11
KIN8M	8192	8	1	0.087 ± 0.004	0.081 ± 0.003	0.84 ± 0.002	<b>0.078 ± 0.003</b>	0.081 ± 0.003
NAVAL PROPULSION	11934	16	2	0.0021 ± 0.0006	<b>0.0004 ± 0.0001</b>	0.026 ± 0.0043	<b>0.0004 ± 0.0001</b>	0.0004 ± 0.0001
PROTEIN STRUCTURE	45730	8	1	4.49 ± 0.11	4.28 ± 0.02	4.91 ± 0.11	<b>4.26 ± 0.05</b>	4.28 ± 0.07
SUPERCONDUCTIVITY	21263	81	1	13.87 ± 0.50	13.02 ± 0.47	14.98 ± 2.57	<b>11.87 ± 0.45</b>	12.48 ± 0.40
WINE QUALITY (RED)	1599	11	1	0.636 ± 0.038	<b>0.635 ± 0.037</b>	0.665 ± 0.040	0.639 ± 0.036	0.633 ± 0.036
WINE QUALITY (WHITE)	4898	11	1	0.691 ± 0.032	0.685 ± 0.035	0.737 ± 0.037	<b>0.678 ± 0.030</b>	0.684 ± 0.038
YACHT HYDRODYNAMICS	308	6	1	<b>1.22 ± 0.47</b>	1.73 ± 1.00	2.303 ± 1.94	1.55 ± 0.79	0.78 ± 0.25

Table 4: Training run-time (in minutes wall-time) for each of the UCI regression datasets from a standard MLP model compared to heteroscedastic training methods (NLL, BNLL) and a non-probabilistic predictor.

DATASET	NLL	BNLL	CBDL	HYBRIDFLOW	MLP (NO UQ)
BOSTON HOUSING	<b>0.19</b>	0.21	10.91	1.11	0.37
CARBON	4.49	3.38	24.50	16.62	<b>2.95</b>
CONCRETE STRENGTH	<b>0.22</b>	0.30	12.27	2.04	0.80
CYCLE POWER PLANT	5.97	4.69	43.84	21.95	<b>4.30</b>
ENERGY EFFICIENCY	<b>0.75</b>	0.89	17.15	2.63	1.19
KIN8M	<b>2.94</b>	3.28	29.51	18.43	2.99
NAVAL PROPULSION	6.60	6.30	24.79	16.47	<b>4.93</b>
PROTEIN STRUCTURE	6.05	6.10	44.83	28.59	<b>4.75</b>
SUPERCONDUCTIVITY	4.70	<b>4.39</b>	37.28	43.35	4.95
WINE QUALITY (RED)	<b>0.30</b>	<b>0.30</b>	7.12	2.73	0.32
WINE QUALITY (WHITE)	0.83	0.82	24.03	8.61	<b>0.66</b>
YACHT HYDRODYNAMICS	<b>0.25</b>	0.28	14.48	1.69	1.17

with epistemic ECE on par with the other models, indicating that its uncertainty estimates are better calibrated than those of existing methods. It also shows improved performance across other uncertainty metrics, including lower CRPS and NLL, and more reliable predictive intervals as measured by the Winkler score. While HybridFlow has slightly wider uncertainty intervals on average (as reflected by MPIW), this trade-off supports better coverage and reliability.

Qualitative analyses of depth maps further support the quantitative findings, the HybridFlow model predicts sharper depth edges and maintains consistency across scenes with varied lighting and texture conditions (Appendix A). HybridFlow’s uncertainty maps reveal elevated aleatoric uncertainty for areas with abnormal lighting or occlusion boundaries, consistent with regions expected to pose challenges due to aleatoric noise. For epistemic uncertainty, higher uncertainty is more localized to objects or textures that are underrepresented in the training dataset. We also see that the aleatoric uncertainty predicted by the heteroscedastic methods closely resembles the predicted depth, which may indicate that the model is unable to effectively learn noise distributions at variable depths.

## 5.2 UCI Regression Benchmarks Results

Table 3 shows that the HybridFlow models consistently are more accurate across the suite of UCI regression benchmarks. Based on RMSE, out of the 12 UCI regression benchmarks we tested, the HybridFlow models outperform the Gaussian NLL, BNLL, and CBDL models on 9 of the datasets. This highlights HybridFlow’s ability to excel in diverse conditions, outperforming the other methods in both high-dimensional settings and tasks characterized by challenging noise distributions.

The NF not only allows the model to quantify aleatoric uncertainty, but in 6 of the benchmarks the HybridFlow model achieves even greater accuracy than the non-probabilistic models. Of the remaining 6 benchmarks, 5 show RMSE values within 5% of the non-probabilistic MLP baseline. This suggests that the

latent feature representations learned by the NF contribute meaningful additional structure to the data, improving predictive performance beyond what is achievable with MSE optimization alone. While the model’s performance surpasses others in most scenarios, a few cases where it matches the performance of baseline methods suggest the influence of dataset-specific complexities on results. Table 4 showcases the training run-times of each of the 4 models tested. Our results show HybridFlow does add to total training time, but less so than the equivalent Bayesian method, CBDL. However, we consider that a small trade-off when comparing with HybridFlow accuracy and argue that the extra training time provides enough practical benefit to justify its use.

Table 5: Uncertainty quantification metrics across models on the UCI datasets. Metrics include NLL (Negative Log-Likelihood), calibration error (ECE), Winkler score, MPIW (Mean Prediction Interval Width), total PICP (Prediction Interval Coverage Probability), and CRPS (Continuous Ranked Probability Score). Metrics are reported for each quantified aleatoric, epistemic, and total (additive) uncertainty.

DATASET	METHOD	NLL			ECE			WINKLER			MPIW			PICP	CRPS
		ALE.	EPI.	TOT.	ALE.	EPI.	TOT.	ALE.	EPI.	TOT.	ALE.	EPI.	TOT.		
BOSTON HOUSING	NLL	2.80	4.11	2.88	2.39	1.79	2.40	80.82	35.17	86.35	80.59	5.73	86.32	<b>1.00</b>	1.98
	BNLL	2.39	<b>3.15</b>	1.98	1.83	<b>1.38</b>	<b>1.88</b>	29.96	<b>28.28</b>	33.43	27.94	5.00	32.94	<b>1.00</b>	<b>1.65</b>
	CBDL	491.23	6.28	5.97	188.00	199.59	200.15	5516.66	991.58	769.99	142.70	520.21	662.91	0.94	1.83
	HYBRIDFLOW	<b>1.64</b>	4.42	<b>1.63</b>	<b>1.78</b>	1.43	1.90	<b>18.77</b>	33.18	<b>19.65</b>	<b>11.81</b>	<b>4.82</b>	<b>16.63</b>	0.98	1.76
CARBON	NLL	43.71	<b>-2.96</b>	<b>-2.92</b>	39.28	22.44	21.06	0.61	<b>0.19</b>	<b>0.20</b>	<b>0.01</b>	<b>0.18</b>	<b>0.19</b>	<b>1.00</b>	<b>0.01</b>
	BNLL	<b>23.53</b>	-2.86	-2.79	27.32	19.40	17.81	<b>0.46</b>	0.21	0.23	0.02	0.21	0.23	<b>1.00</b>	0.02
	CBDL	370.94	317.27	345.93	<b>19.62</b>	<b>4.29</b>	<b>4.39</b>	7.92	37.46	36.61	0.02	0.96	0.98	0.89	0.03
	HYBRIDFLOW	24.63	7.68	52.15	35.24	22.21	23.90	<b>0.46</b>	0.24	0.26	0.02	0.19	0.21	0.99	0.02
CONCRETE STRENGTH	NLL	3.56	4.03	3.65	4.51	4.12	4.51	165.80	60.02	177.13	165.80	11.33	177.13	<b>1.00</b>	3.58
	BNLL	3.36	<b>3.80</b>	3.46	<b>4.15</b>	<b>3.76</b>	<b>4.16</b>	126.25	<b>54.34</b>	137.57	126.25	11.32	137.57	<b>1.00</b>	3.29
	CBDL	148.05	7.07	6.97	588.89	587.23	588.65	1567.22	254.89	257.75	138.99	209.12	348.11	0.98	<b>3.22</b>
	HYBRIDFLOW	<b>2.29</b>	7.85	<b>2.40</b>	4.40	4.06	4.44	<b>35.74</b>	78.59	<b>40.83</b>	<b>30.28</b>	<b>9.23</b>	<b>39.51</b>	0.99	3.76
CYCLE POWER	NLL	3.03	5.45	3.10	3.22	2.56	3.23	82.94	49.66	89.61	82.94	6.68	89.62	<b>1.00</b>	<b>2.64</b>
	BNLL	3.00	<b>5.30</b>	3.07	3.22	2.57	3.22	79.40	<b>49.41</b>	86.03	79.39	6.64	86.03	<b>1.00</b>	<b>2.64</b>
	CBDL	135.34	60.87	55.74	577.36	770.71	770.32	306.80	253.72	251.01	142.80	293.63	436.43	0.72	2.91
	HYBRIDFLOW	<b>1.93</b>	6.53	<b>2.02</b>	<b>3.03</b>	<b>2.51</b>	<b>3.10</b>	<b>20.53</b>	51.85	<b>24.61</b>	<b>16.62</b>	<b>6.35</b>	<b>22.97</b>	0.99	2.66
ENERGY EFFICIENCY	NLL	2.79	3.49	2.90	2.25	1.70	2.26	76.70	33.73	83.42	76.70	6.72	83.42	<b>1.00</b>	1.86
	BNLL	1.64	2.76	1.89	1.57	0.98	1.61	31.78	27.65	36.60	31.78	<b>4.82</b>	36.60	<b>1.00</b>	1.46
	CBDL	171.68	7.87	6.69	225.24	117.98	189.01	709.48	240.01	156.14	104.48	416.82	521.30	0.65	1.48
	HYBRIDFLOW	<b>1.17</b>	<b>2.08</b>	<b>1.43</b>	<b>1.20</b>	<b>0.66</b>	<b>1.18</b>	<b>14.70</b>	<b>18.50</b>	<b>19.28</b>	<b>14.19</b>	5.09	<b>19.28</b>	<b>1.00</b>	<b>1.16</b>
KIN8M	NLL	20.67	0.29	-1.25	69.27	<b>20.00</b>	14.50	2.49	1.28	0.84	0.09	0.25	0.34	0.79	0.08
	BNLL	23.49	<b>-0.60</b>	-1.37	76.36	20.01	15.47	2.32	<b>0.91</b>	0.69	<b>0.07</b>	0.25	<b>0.32</b>	0.81	0.07
	CBDL	2991.21	400.10	363.54	67.91	72.79	146.19	8.87	11.16	10.64	0.25	<b>0.08</b>	0.33	0.84	0.07
	HYBRIDFLOW	<b>-1.94</b>	1.76	<b>-1.86</b>	<b>11.61</b>	35.99	<b>8.42</b>	<b>0.43</b>	1.13	<b>0.51</b>	0.37	0.13	0.50	<b>0.99</b>	<b>0.06</b>
NAVAL PROPULSION	NLL	-3.24	-0.60	-2.89	0.26	0.18	0.29	0.39	1.03	0.59	0.05	0.21	0.26	0.98	0.02
	BNLL	-2.61	-0.54	-2.33	0.38	0.27	0.41	0.46	0.95	0.62	0.06	0.22	0.28	0.97	0.03
	CBDL	201.38	132.27	155.82	778.69	182.37	261.39	0.50	<b>0.51</b>	0.51	0.05	0.35	0.40	0.88	<b>0.01</b>
	HYBRIDFLOW	<b>-3.79</b>	<b>-0.47</b>	<b>-3.32</b>	<b>0.17</b>	<b>0.12</b>	<b>0.19</b>	<b>0.31</b>	0.86	<b>0.48</b>	<b>0.03</b>	<b>0.19</b>	<b>0.22</b>	<b>0.99</b>	0.02
PROTEIN STRUCTURE	NLL	3.34	63.83	3.37	4.40	2.61	4.40	112.45	133.39	115.13	112.45	<b>2.68</b>	115.13	<b>1.00</b>	4.11
	BNLL	2.97	33.82	3.02	3.66	2.28	3.66	84.57	<b>96.59</b>	88.05	84.45	3.51	87.96	<b>1.00</b>	3.32
	CBDL	92.66	<b>4.88</b>	4.90	40.91	47.15	48.17	162.60	473.69	487.74	<b>14.05</b>	473.69	487.74	<b>1.00</b>	<b>3.00</b>
	HYBRIDFLOW	<b>1.97</b>	38.88	<b>1.95</b>	<b>3.26</b>	<b>1.84</b>	<b>3.30</b>	<b>21.54</b>	97.91	<b>22.74</b>	17.17	2.87	<b>20.04</b>	0.97	3.20
SUPERCONDUCTIVITY	NLL	6.04	10.57	6.06	14.55	14.29	14.55	2013.14	313.96	2033.23	2013.14	20.09	2033.23	<b>1.00</b>	12.49
	BNLL	<b>4.28</b>	<b>6.01</b>	4.33	8.12	7.73	8.12	629.83	<b>121.04</b>	647.08	629.24	17.58	646.82	<b>1.00</b>	6.51
	CBDL	707.73	8.45	9.47	580.43	157.97	157.99	402.27	173.89	1770.73	3730.98	1730.89	5461.87	<b>1.00</b>	7.60
	HYBRIDFLOW	8.63	6.67	<b>2.74</b>	<b>7.65</b>	<b>7.54</b>	<b>7.77</b>	<b>57.98</b>	128.32	<b>63.02</b>	<b>41.91</b>	<b>16.89</b>	<b>58.80</b>	0.99	<b>6.40</b>
WINE QUALITY (RED)	NLL	1.03	23.57	2.09	1.91	10.09	1.38	4.12	13.52	3.41	1.76	<b>0.43</b>	2.19	0.91	0.46
	BNLL	1.57	<b>17.24</b>	0.25	3.18	8.98	1.95	5.10	<b>11.94</b>	3.70	<b>1.35</b>	0.51	<b>1.86</b>	0.86	<b>0.43</b>
	CBDL	75.83	106.97	47.26	8.65	<b>2.44</b>	2.65	172.82	144.63	137.57	0.37	1.85	2.22	0.76	<b>0.43</b>
	HYBRIDFLOW	<b>0.08</b>	19.97	<b>0.05</b>	<b>1.26</b>	9.48	<b>0.95</b>	<b>3.31</b>	12.52	<b>3.27</b>	2.38	0.47	2.85	<b>0.97</b>	0.44
WINE QUALITY (WHITE)	NLL	0.40	20.24	0.16	1.70	8.59	1.20	4.19	13.78	3.54	1.92	0.50	2.42	0.92	0.48
	BNLL	0.49	<b>18.62</b>	0.16	1.88	8.14	1.28	4.36	<b>12.99</b>	3.53	<b>1.81</b>	0.56	<b>2.37</b>	0.91	<b>0.47</b>
	CBDL	26.46	33.25	25.18	6.36	<b>3.60</b>	3.78	198.85	154.03	147.42	0.34	2.70	3.04	0.64	0.49
	HYBRIDFLOW	<b>0.22</b>	22.47	<b>0.13</b>	<b>0.99</b>	9.08	<b>0.75</b>	<b>3.41</b>	13.85	<b>3.48</b>	2.65	<b>0.49</b>	3.14	<b>0.97</b>	0.48
YACHT HYDRODYNAMICS	NLL	1.21	4.33	1.31	4.07	3.53	4.29	6.03	5.57	6.58	6.30	3.31	9.61	<b>1.00</b>	0.28
	BNLL	1.35	<b>3.98</b>	1.25	3.10	2.54	3.11	5.36	<b>4.03</b>	5.52	5.66	2.50	8.16	<b>1.00</b>	0.25
	CBDL	29.11	32.17	30.11	10.99	18.83	15.50	138.35	168.73	183.31	25.39	15.44	40.83	0.71	0.30
	HYBRIDFLOW	<b>0.63</b>	5.00	<b>0.62</b>	<b>2.01</b>	<b>1.51</b>	<b>2.19</b>	<b>2.59</b>	5.80	<b>3.17</b>	<b>2.26</b>	<b>1.21</b>	<b>3.47</b>	0.98	<b>0.22</b>

Table 6: Comparison of performance metrics for three ice sheet emulators: HybridFlow, a Gaussian Process, and a non-probabilistic predictor (Deep Ensemble of LSTMs without uncertainty quantification), evaluated on the Antarctic Ice Sheet (AIS) and Greenland Ice Sheet (GrIS).

	EMULATOR	MSE	ECE	PICP
AIS	HYBRIDFLOW	<b>1.20</b>	<b>0.02</b>	<b>0.95</b>
	GAUSSIAN PROCESS	3.05	0.08	0.90
	DEEP ENSEMBLE (NO UQ)	1.09	–	–
GRIS	HYBRIDFLOW	<b>1.02</b>	<b>0.01</b>	<b>0.96</b>
	GAUSSIAN PROCESS	9.87	0.12	0.91
	DEEP ENSEMBLE (NO UQ)	0.87	–	–

Across the UCI suite HybridFlow showcases superior performance on nearly all uncertainty quantification metrics. It attains the lowest total NLL on 11 of 12 datasets, the lowest Winkler score on 11 of 12, and the smallest total ECE on 9 of 12 (Table 5). Despite these sharper and better-calibrated distributions, it still produces the narrowest intervals (lowest MPIW on 8 of 12 datasets) while preserving reliability (PICP remains near 0.97 on every benchmark). These proportions confirm that decoupling aleatoric and epistemic components with HybridFlow yields predictive uncertainties that are simultaneously sharp, well-calibrated, and trustworthy across nearly all of the UCI benchmarks.

## 6 Case Study: Hybrid Ice Sheet Emulator

To test the efficacy of HybridFlow in a real-world scenario and demonstrate how it could be implemented with an existing model structure, we use the HybridFlow framework to create an emulator for projecting ice sheet dynamics and their contributions to sea level rise. Ice sheet evolution is governed by complex, nonlinear processes, such as melting, snow accumulation, and ice shelf instability. These processes involve large amounts of uncertainty, which originate both from the physical modeling of the processes and from the chaotic nature of the processes themselves. Ice sheet emulators (Van Katwyk et al., 2023; Edwards et al., 2021) provide efficient and reliable approximation of the complex physics-based models, which allows ice sheet scientists to perform sensitivity testing to understand the effect that climate variables have on sea level projections, and enables more experimentation to understand future sea level.

Recent work has demonstrated the effectiveness of LSTM-based ice sheet emulators (Van Katwyk et al., 2023) in improving predictive accuracy over Gaussian Process-based approaches (Edwards et al., 2021) on the ISMIP6 ice sheet model projection dataset (Nowicki et al., 2016). Building on these advancements, we use a CMAF and a Deep Ensemble (Lakshminarayanan et al., 2017) of LSTM models to effectively model the temporal structure of ice sheet projections from 2015 to 2100. We compare the HybridFlow emulator to a Gaussian Process emulator based on the individual projection performance as well as the ability to approximate the full range of ensemble predictions. Consistent with the approach in Section 4, we train a non-probabilistic emulator using the LSTM model based on Van Katwyk et al. (2023).

Table 6 shows that the HybridFlow framework achieves superior performance in both test accuracy and the ability to reproduce the full range of ensemble projections compared to the Gaussian Process-based ice sheet emulator. HybridFlow performs similarly to the non-probabilistic predictor in test accuracy (MSE), highlighting the framework’s ability to deliver reliable and actionable projections in the context of real data. Furthermore, the framework’s ability to generate calibrated (PICP) and accurate (ECE) uncertainty estimates ensures that it can provide confidence intervals suitable for policy and planning decisions.

## 7 Conclusion

In this work we introduce HybridFlow, a novel hybrid architecture that combines the strength of NFs and probabilistic prediction models to quantify both aleatoric and epistemic uncertainty within a single unified model. By leveraging the flexibility of NFs to model data-specific aleatoric uncertainty and incorporating

probabilistic models for epistemic uncertainty estimation, HybridFlow achieves robust and calibrated predictions without compromising predictive accuracy. HybridFlow’s design addresses two gaps in uncertainty quantification methods: the decrease in model predictive accuracy to quantify uncertainty and the miscalibration of current state-of-the-art uncertainty quantification methods (Seitzer et al., 2022). HybridFlow achieves this by decoupling the loss functions used for aleatoric and epistemic uncertainty estimation, avoiding the pitfalls of joint NLL-based training methods. Decoupling the loss functions allows for task-specific loss functions for predictors, which results in predictive accuracy levels comparable to non-probabilistic models while maintaining the ability to quantify and separate sources of uncertainty. We demonstrate the ability to consistently achieve better metrics for accuracy (RMSE) and uncertainty quantification (ECE) on a variety of regression tasks, including depth estimation and a series of standard benchmarks. Furthermore, we present a case study that demonstrates the ability of the HybridFlow framework to be effective in real-world applications, offering accurate emulation of complex systems such as continental-scale ice sheets.

HybridFlow offers a particularly practical solution for domain scientists and researchers who seek to incorporate robust uncertainty quantification into existing modeling workflows without changing their entire system. In many applied settings, such as Earth system modeling or ice sheet forecasting (Section 6), predictive frameworks are already well established and finely tuned (Irrgang et al., 2021); replacing them with task-specific uncertainty models is often infeasible. HybridFlow avoids this disruption by providing a modular framework that can be integrated with existing probabilistic predictors, enabling calibrated and interpretable estimates of both aleatoric and epistemic uncertainty with minimal architectural changes. This makes HybridFlow especially valuable for scientific applications where improving the transparency and trustworthiness of predictions is critical, but where flexibility and compatibility with legacy models remain essential.

HybridFlow’s modular design ensures adaptability to a wide range of tasks, allowing for future advancements in robust and trustworthy AI systems. However, a key drawback of this method is the added computation required for implementation. In the test cases presented in this study, the HybridFlow model for depth estimation took 44% longer to train than the NLL-based models. Similarly, Table 4 shows that HybridFlow was also slower to train for the UCI datasets, as the NF must be trained before the training of the predictor models. However, HybridFlow requires less computational overhead than the CBDL method (as can be seen in Table 4), which took one or two orders of magnitude longer than both the NLL methods and HybridFlow to train, but still remains a potential drawback of HybridFlow. The NLL and BNLL methods, on the other hand, are simple to implement and the added ability to quantify aleatoric uncertainty requires a negligible increase in computation to add to any model. By separating the quantification of aleatoric uncertainty into a separate module, we quantify uncertainty more accurately but at the cost of added computational complexity. Therefore, future work may focus both on expanding the framework’s scalability and exploring alternative generative modeling techniques, including Bayesian NFs, to further enhance its versatility, precision, and reliability in high-dimensional or complex data environments.

While HybridFlow provides a modular approach to separately estimate aleatoric and epistemic uncertainty, it is important to acknowledge that disentangling these two sources of uncertainty is inherently difficult in practice. Despite conceptual distinctions, they are often entangled in complex models and real-world data, and there is currently no definitive or universally accepted solution to completely separate them (Smith et al., 2024; de Jong et al., 2024; Valdenegro-Toro & Mori, 2022). Attempts to decompose uncertainty are still valuable, as demonstrated in our ice sheet case study, where understanding the source of uncertainty informs scientific and policy decisions. Even imperfect separation can provide useful insight into model limitations and the nature of data variability.

Future work should also include the investigation of epistemic uncertainty within the NF itself, as it is not explicitly modeled in the current implementation. Although this practice aligns with conventional flow-based modeling approaches, incorporating methods such as flow ensembles (Berry & Meger, 2023) could enhance the robustness of aleatoric uncertainty estimates by capturing model-level epistemic uncertainty, even though due to the size of the NF model, the epistemic uncertainty is likely small. Future work may also include the rigorous testing of this framework on a variety of classification tasks, with a comparison to Sale et al. (2024), including a detailed analysis of quantified uncertainty.

## Broader Impact Statement

HybridFlow provides a practical and accessible framework for uncertainty quantification, enabling users to enhance the reliability of their models without significant architectural changes or specialized expertise. By prioritizing modularity and ease of integration, HybridFlow lowers the barrier to adopting robust uncertainty estimation, making it especially valuable for applications where trust and interpretability are critical.

Rather than focusing solely on novelty, this work emphasizes usability and flexibility—characteristics that support broader adoption across disciplines, including scientific modeling, environmental forecasting, and healthcare. While users should remain mindful of data limitations, HybridFlow equips them with a straightforward tool for improving model transparency and decision-making confidence, contributing to more trustworthy AI systems.

## References

- Dario Amodei et al. Concrete problems in AI safety. *arXiv preprint arXiv:1606.06565*, 2016.
- Lucas Berry and David Meger. Normalizing flow ensembles for rich aleatoric and epistemic uncertainty modeling. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 6806–6814, 2023. doi: 10.1609/aaai.v37i6.2583.
- Eli Bingham, Jonathan P Chen, Martin Jankowiak, Fritz Obermeyer, Neeraj Pradhan, Theofanis Karaletsos, Rohit Singh, Paul Szerlip, Paul Horsfall, and Noah D Goodman. Pyro: Deep universal probabilistic programming. *Journal of machine learning research*, 20(28):1–6, 2019.
- Christopher M Bishop. Mixture density networks. Technical Report NCRG/94/004, Aston University, Birmingham, UK, 1994.
- Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in neural networks. In *Proceedings of the International Conference on Machine Learning (ICML)*, pp. 1613–1622, 2015. doi: 10.48550/arXiv.1505.05424.
- Michele Caprio, Souradeep Dutta, Kuk Jin Jang, Vivian Lin, Radoslav Ivanov, Oleg Sokolsky, and Insup Lee. Credal Bayesian deep learning. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856.
- Matthew Albert Chan, Maria J Molina, and Christopher Metzler. Estimating epistemic and aleatoric uncertainty with a single model. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. doi: 10.48550/arXiv.2402.03478.
- Ivo Pascal de Jong, Andreea Ioana Sburlea, and Matias Valdenegro-Toro. How disentangled are your classification uncertainties? *arXiv preprint arXiv:2408.12175*, 2024.
- Stefan Depeweg, Jose-Miguel Hernandez-Lobato, Finale Doshi-Velez, and Steffen Udluft. Decomposition of uncertainty in Bayesian deep learning for efficient and risk-sensitive learning. In *International Conference on Machine Learning*, pp. 1184–1193. PMLR, 2018.
- Armen Der Kiureghian and Ove Ditlevsen. Aleatoric or epistemic? Does it matter? *Structural Safety*, 31(2):105–112, 2009. doi: 10.1016/j.strusafe.2008.06.020.
- Conor Durkan, Artur Bekasov, Iain Murray, and George Papamakarios. nflows: normalizing flows in PyTorch, November 2020.
- Tamsin L Edwards, Sophie Nowicki, Ben Marzeion, Regine Hock, Heiko Goelzer, Hélène Seroussi, et al. Projected land ice contributions to twenty-first-century sea level rise. *Nature*, 593(7857):74–82, 2021. doi: 10.1038/s41586-021-03302-y.
- David Eigen, Christian Puhersch, and Rob Fergus. Depth map prediction from a single image using a multi-scale deep network. In *Advances in Neural Information Processing Systems*, 2014. doi: 10.48550/arXiv.1406.2283.

- Yarin Gal and Zoubin Ghahramani. Dropout as a Bayesian approximation: Representing model uncertainty in deep learning. In *International Conference on Machine Learning*, pp. 1050–1059. PMLR, 2016. doi: 10.48550/arXiv.1506.02142.
- Ashkan Ganj, Hang Su, and Tian Guo. HybridDepth: Robust metric depth fusion by leveraging depth from focus and single-image priors. In *2025 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pp. 973–982, 2025. doi: 10.1109/WACV61041.2025.00104.
- Tilmann Gneiting and Adrian E Raftery. Strictly proper scoring rules, prediction, and estimation. *Journal of the American Statistical Association*, 102(477):359–378, 2007. doi: 10.1198/016214506000001437.
- Ruben Grewal, Paolo Tonella, and Andrea Stocco. Predicting Safety Misbehaviours in Autonomous Driving Systems Using Uncertainty Quantification . In *2024 IEEE Conference on Software Testing, Verification and Validation (ICST)*, pp. 70–81, 2024. doi: 10.1109/ICST60714.2024.00016.
- Cornelia Gruber, Patrick Oliver Schenk, Malte Schierholz, Frauke Kreuter, and Göran Kauermann. Sources of uncertainty in machine learning—a statisticians’ view. *arXiv preprint arXiv:2305.16703*, 2023.
- José Miguel Hernández-Lobato and Ryan Adams. Probabilistic backpropagation for scalable learning of Bayesian neural networks. In *International Conference on Machine Learning*, pp. 1861–1869. PMLR, 2015.
- Kexin Huang, Ying Jin, Emmanuel Candes, and Jure Leskovec. Uncertainty quantification over graph with conformalized graph neural networks. *Advances in Neural Information Processing Systems*, 36:26699–26721, 2023. doi: 10.48550/arXiv.2305.14535.
- Hyekyoung Hwang and Jitae Shin. Uncertainty Measurement of Deep Learning System Based on the Convex Hull of Training Sets. *arXiv preprint arXiv:2405.16082*, 2024.
- Eyke Hüllermeier and Willem Waegeman. Aleatoric and epistemic uncertainty in machine learning: an introduction to concepts and methods. *Machine Learning*, 110(3):457–506, March 2021. ISSN 1573-0565. doi: 10.1007/s10994-021-05946-3.
- Christopher Irrgang, Niklas Boers, Maike Sonnwald, Elizabeth A Barnes, Christopher Kadow, Joanna Staneva, and Jan Saynisch-Wagner. Towards neural earth system modelling by integrating artificial intelligence in earth system science. *Nature Machine Intelligence*, 3(8):667–674, 2021.
- Nyaz Jalal, Małgorzata Śliwińska, Wadim Wojciechowski, Iwona Kucybała, Miłosz Rozynek, Kamil Krupa, Patrycja Matusik, Jarosław Jarczewski, and Zbysław Tabor. Evaluating Uncertainty Quantification in Medical Image Segmentation: A Multi-Dataset, Multi-Algorithm Study. *Applied Sciences (2076-3417)*, 14(21), 2024. doi: 10.3390/app142110020.
- Alex Kendall and Yarin Gal. What uncertainties do we need in Bayesian deep learning for computer vision? *Advances in Neural Information Processing Systems*, 30, 2017.
- Diederik P Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. In *3rd International Conference for Learning Representations (ICLR)*, 2015. doi: 10.48550/arXiv.1412.6980.
- Ivan Kobyzev, Simon JD Prince, and Marcus A Brubaker. Normalizing flows: An introduction and review of current methods. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43(11):3964–3979, 2020. doi: 10.1109/TPAMI.2020.2992934.
- Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In *Advances in Neural Information Processing Systems (NeurIPS)*, pp. 6402–6413, 2017.
- Laura Marusich, Jonathan Bakdash, Yan Zhou, and Murat Kantarcioglu. Using ai uncertainty quantification to improve human decision-making. In *International Conference on Machine Learning*, pp. 34949–34960. PMLR, 2024. doi: 10.48550/arXiv.2309.10852.

- Giovanni Messuti, Ortensia Amoroso, Ferdinando Napolitano, Mariarosaria Falanga, Paolo Capuano, and Silvia Scarpetta. Uncertainty estimation via ensembles of deep learning models and dropout layers for seismic traces. In *Advanced Neural Artificial Intelligence: Theories and Applications*, pp. 107–117. Springer, 2025. doi: 10.1007/978-981-96-0994-9\_10.
- Eric Nalisnick, Akihiro Matsukawa, Yee Whye Teh, Dilan Gorur, and Balaji Lakshminarayanan. Hybrid models with deep and invertible features. In *International Conference on Machine Learning*, pp. 4723–4732. PMLR, 2019.
- Sophie Nowicki, Anthony Payne, Eric Larour, Helene Seroussi, Heiko Goelzer, William Lipscomb, Jonathan Gregory, Ayako Abe-Ouchi, and Andrew Shepherd. Ice sheet model intercomparison project (ISMIP6) contribution to CMIP6. *Geoscientific Model Development*, 9(12):4521–4545, 2016. doi: 10.5194/gmd-9-4521-2016.
- George Papamakarios, Theo Pavlakou, and Iain Murray. Masked autoregressive flow for density estimation. In *Advances in Neural Information Processing Systems*, pp. 2338–2347, 2017.
- Apostolos F Psaros, Xuhui Meng, Zongren Zou, Ling Guo, and George Em Karniadakis. Uncertainty quantification in scientific machine learning: Methods, metrics, and comparisons. *Journal of Computational Physics*, 477:111902, 2023. doi: 10.1016/j.jcp.2022.111902.
- Danilo Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *International Conference on Machine Learning*, pp. 1530–1538. PMLR, 2015.
- Yusuf Roohani, Kexin Huang, and Jure Leskovec. Predicting transcriptional outcomes of novel multigene perturbations with GEARS. *Nature Biotechnology*, 42(6):927–935, 2024. doi: 10.1038/s41587-023-01905-6.
- Yusuf Sale, Paul Hofman, Timo Löhr, Lisa Wimmer, Thomas Nagler, and Eyke Hüllermeier. Label-wise aleatoric and epistemic uncertainty quantification. In *Uncertainty in Artificial Intelligence*, pp. 3159–3179. PMLR, 2024. doi: 10.48550/arXiv.2406.02354.
- Maximilian Seitzer, Arash Tavakoli, Dimitrije Antic, and Georg Martius. On the pitfalls of heteroscedastic uncertainty estimation with probabilistic neural networks. In *Tenth International Conference on Learning Representations (ICLR)*, 2022. doi: 10.48550/arXiv.2203.09168.
- Murat Sensoy, Lance Kaplan, and Melih Kandemir. Evidential deep learning to quantify classification uncertainty. *Advances in Neural Information Processing Systems*, 31, 2018. doi: 10.48550/arXiv.1806.01768.
- Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob Fergus. Indoor segmentation and support inference from RGBD images. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 746–760. Springer, 2012. doi: 10.1007/978-3-642-33715-4\_54.
- Freddie Bickford Smith, Jannik Kossen, Eleanor Trollope, Mark van der Wilk, Adam Foster, and Tom Rainforth. Rethinking aleatoric and epistemic uncertainty. In *NeurIPS 2024 Workshop on Bayesian Decision-making and Uncertainty*, 2024. doi: 10.48550/arXiv.2412.20892.
- Arthur Thuy and Dries F Benoit. Explainability through uncertainty: Trustworthy decision-making with neural networks. *European Journal of Operational Research*, 317(2):330–340, 2024. doi: 10.1016/j.ejor.2023.09.009.
- 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)*, pp. 1508–1516. IEEE, 2022. doi: 10.1109/CVPRW56347.2022.00157.
- Peter Van Katwyk, Baylor Fox-Kemper, Helene Seroussi, Sophie Nowicki, and Karianne J Bergen. A variational LSTM emulator of sea level contribution from the Antarctic ice sheet. *Journal of Advances in Modeling Earth Systems*, 15(12):e2023MS003899, 2023. doi: 10.1029/2023MS003899.

- Hanchen Wang, Tianfan Fu, Yuanqi Du, Wenhao Gao, Kexin Huang, Ziming Liu, Payal Chandak, Shengchao Liu, Peter Van Katwyk, Andreea Deac, et al. Scientific discovery in the age of artificial intelligence. *Nature*, 620(7972):47–60, 2023. doi: 10.1038/s41586-023-06221-2.
- Tianyang Wang, Yunze Wang, Jun Zhou, Benji Peng, Xinyuan Song, Charles Zhang, Xintian Sun, Qian Niu, Junyu Liu, Silin Chen, et al. From aleatoric to epistemic: Exploring uncertainty quantification techniques in artificial intelligence. *arXiv preprint arXiv:2501.03282*, 2025.
- Lisa Wimmer, Yusuf Sale, Paul Hofman, Bernd Bischl, and Eyke Hüllermeier. Quantifying aleatoric and epistemic uncertainty in machine learning: Are conditional entropy and mutual information appropriate measures? In Robin J. Evans and Ilya Shpitser (eds.), *Proceedings of the Thirty-Ninth Conference on Uncertainty in Artificial Intelligence*, volume 216 of *Proceedings of Machine Learning Research*, pp. 2282–2292. PMLR, 31 Jul–04 Aug 2023.
- Chu-I Yang and Yi-Pei Li. Explainable uncertainty quantifications for deep learning-based molecular property prediction. *Journal of Cheminformatics*, 15(1):13, 2023. doi: 10.1186/s13321-023-00682-3.
- Sai Harsha Yelleni, Deepshikha Kumari, PK Srijith, et al. Monte Carlo DropBlock for modeling uncertainty in object detection. *Pattern Recognition*, 146:110003, 2024. doi: 10.1016/j.patcog.2023.110003.
- Tae Sung Yoon and Heeyoung Kim. Uncertainty Estimation by Density Aware Evidential Deep Learning. In *Proceedings of the 41st International Conference on Machine Learning*, pp. 57217–57243. PMLR, 2024. doi: 10.48550/arXiv.2409.08754.
- Baobing Zhang, Wanxin Sui, Zhengwen Huang, Maozhen Li, and Man Qi. Normalizing flow based uncertainty estimation for deep regression analysis. *Neurocomputing*, 585:127645, 2024. doi: 10.1016/j.neucom.2024.127645.

## Appendix

### A Depth Estimation Visualizations

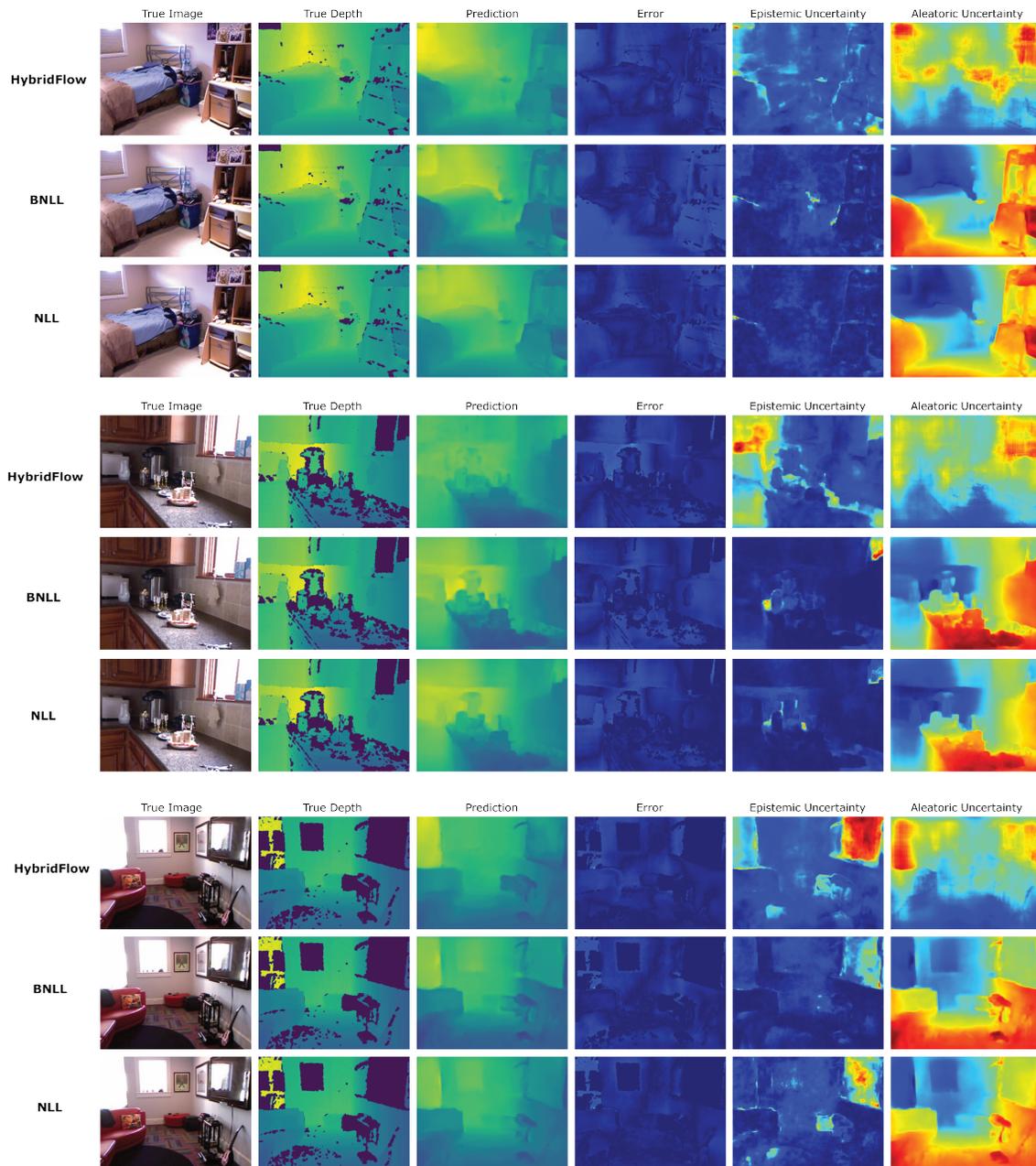


Figure 2: Visualization of depth estimation results using the HybridFlow framework compared to the BNLL (Seitzer et al., 2022) and NLL (Kendall & Gal, 2017) on the NYU Depth v2 dataset. Each row showcases a sample scene, with the columns representing (from left to right): the input RGB image, ground truth depth map, predicted depth map, model error, aleatoric uncertainty map, and epistemic uncertainty map. The predicted depth maps closely align with the ground truth, demonstrating accurate scene reconstruction. The aleatoric uncertainty maps highlight areas with inherent measurement noise, such as changes in image lighting (glare) or occlusion boundaries, while the epistemic uncertainty maps identify regions of the scene where the predictive model is uncertain. These visualizations illustrate the HybridFlow model’s ability to provide both accurate predictions and reliable uncertainty quantification.

## B Autoencoder for Depth Estimation

For the depth estimation experiments detailed in Section 4.1, we utilized an autoencoder to compress the dimensionality of RGB images from the NYU Depth v2 dataset. The autoencoder employs a UNet-inspired architecture with an encoder-decoder structure enhanced by skip connections. The encoder gradually reduces the spatial dimensions of the input images through six convolutional layers, each followed by batch normalization and ReLU activation. The bottleneck compresses the high-dimensional input ( $384 \times 512 \times 3$ ) into a compact 256-dimensional latent representation. The decoder reconstructs the images using transposed convolutions, with skip connections from the encoder aiding in preserving spatial details. A sigmoid activation in the final layer ensures that pixel values are mapped to the range  $[0, 1]$ .

The autoencoder was trained using a composite loss function designed to ensure high-quality reconstruction and edge preservation. This loss combines Mean Squared Error (MSE) with the Structural Similarity Index Measure (SSIM) to focus on perceptual similarity, alongside an edge preservation term that penalizes differences in gradients between the input and reconstructed images. Training was conducted with the Adam optimizer at a learning rate of  $1 \times 10^{-4}$ , decayed by a factor of 0.5 based on validation loss, for a maximum of 50 epochs. Early stopping with a patience of 5 epochs was applied to prevent overfitting. Data augmentation, including random cropping, flipping, and brightness adjustments, was employed to improve generalization.

Extensive testing was conducted with latent dimensions of 512, 1024, 128, and 64, but it was observed that performance improvements plateaued beyond 256 dimensions, so we selected the 256-dimensional latent space to balance computational efficiency and reconstruction quality. On the NYU Depth v2 validation set, the autoencoder achieved an MSE of 0.012, SSIM of 0.943, and an edge loss of 0.005.

The latent representations generated by the autoencoder were used as inputs to the Conditional Masked Autoregressive Flow (CMAF) within the HybridFlow framework. This approach allowed for accurate modeling of aleatoric uncertainty while maintaining computational efficiency. The decoupling of representation learning from flow modeling streamlined training and facilitated convergence, showcasing the flexibility and modularity of the HybridFlow architecture.

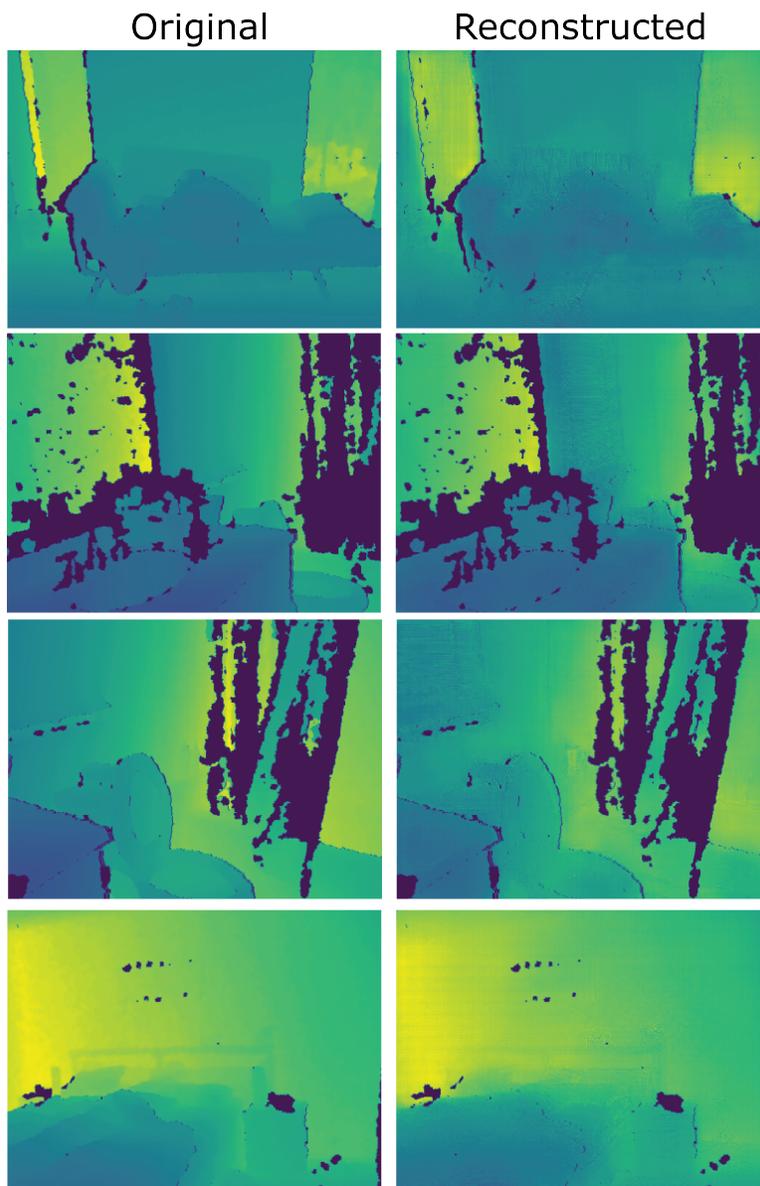


Figure 3: Reconstruction results from the autoencoder trained on the NYU Depth v2 dataset. Each column shows the true depth image (left) and the corresponding reconstructed depth map (right). The autoencoder effectively compresses high-dimensional input data into a 256-dimensional latent representation, enabling efficient dimensionality reduction while preserving fine-grained spatial details.