

---

# Lightning UQ Box: A Comprehensive Framework for Uncertainty Quantification in Deep Learning

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1       Uncertainty quantification (UQ) is an essential tool for applying deep neural  
2       networks (DNNs) to real world tasks, as it attaches a degree of confidence to  
3       DNN outputs. However, despite its benefits, UQ is often left out of the standard  
4       DNN workflow due to the additional technical knowledge required to apply and  
5       evaluate existing UQ procedures. Hence there is a need for a comprehensive  
6       toolbox that allows the user to integrate UQ into their modelling workflow, without  
7       significant overhead. We introduce *Lightning UQ Box*: a unified interface for  
8       applying and evaluating various approaches to UQ. In this paper, we provide a  
9       theoretical and quantitative comparison of the wide range of state-of-the-art UQ  
10      methods implemented in our toolbox. We focus on two challenging vision tasks:  
11      (i) estimating tropical cyclone wind speeds from infrared satellite imagery and  
12      (ii) estimating the power output of solar panels from RGB images of the sky. By  
13      highlighting the differences between methods our results demonstrate the need for  
14      a broad and approachable experimental framework for UQ, that can be used for  
15      benchmarking UQ methods. The toolbox, example implementations, and further  
16      information are available at: <https://github.com/lightning-uq-box/lightning-uq-box>.

## 17 **1 Introduction**

18    In real world applications, deep learning (DL) models are often deployed in safety-critical domains  
19    such as healthcare [45], robotics [56], and Earth observation [55, 60], with relevant areas includ-  
20    ing flood monitoring [7], wildfire mapping and forecasting [58], and weather forecasting [59]. In  
21    these fields, an incorrect prediction can cause significant damage and corresponding consequences.  
22    Uncertainty quantification (UQ) aims to provide a measure of confidence about a neural network's  
23    prediction and to support practitioners in identifying potentially false predictions to better guide anal-  
24    yses and decision-making processes [21]. Besides this, UQ can even improve predictive performance  
25    via regularization [16, 39].

26    The direct application of UQ to DL is often not straightforward for practitioners. Besides the  
27    implementation challenges associated with probabilistic modelling and stochastic training algorithms,  
28    the performance of UQ methods can fluctuate, depending on the data and the task [50]. Moreover,  
29    there is a lack of clear guidance on which methods are promising for specific tasks, given the ever-  
30    increasing zoo of UQ methods for DL [1, 21]. These challenges are particularly prominent for data  
31    modalities of higher dimensions, such as vision, where uncertainty modelling adds a further layer of  
32    complexity. Therefore, various approaches need to be considered, which usually involve different  
33    loss functions, training procedures, and model architectures. The need of accessible and open-source  
34    UQ frameworks is also called upon in a recent position paper on Bayesian Deep Learning (BDL)  
35    by leading experts in this field [53]: "Software development is key to encouraging DL practitioners  
36    to use Bayesian methods. More generally, there is a need for software that would make it easier for

37 practitioners to try BDL in their projects. The use of BDL must become competitive in human effort  
38 with standard deep learning." [53].

39 `Lightning UQ Box` provides users with all the tools needed to equip deep neural networks (DNNs)  
40 with UQ. We created `Lightning UQ Box` to tackle the gap between theoretical researchers and  
41 actual practitioners in the field of UQ in DL. The toolbox offers a comprehensive framework, building  
42 on top of `PyTorch` [54] and `Lightning` [18], as an accessible end-to-end solution. The toolbox  
43 is particularly suited for vision applications (see Section 3): it offers flexible layer configurations  
44 like Bayesian convolution layers that can be modularly placed in backbone architectures, which  
45 streamline UQ.

46 We underline the usefulness of the presented toolbox with two example applications: estimating the  
47 maximum sustained wind speed of tropical cyclones from satellite imagery and predicting the power  
48 voltage output of solar panels from a time series of sky images. These applications contain different  
49 sources and types of uncertainties in the input and target variables and illustrate the stochastic nature  
50 of real world phenomena and measurement systems practitioners are confronted with. Simultaneously,  
51 these applications carry an associated inherent risk that demands reliable predictive uncertainties.  
52 The central contributions of our work all aim to equip practitioners with the necessary tools to apply  
53 UQ methods for DL on their specific (real world) problem:

- 54 • **Comprehensive End-to-End UQ Toolbox:** `Lightning UQ Box` enables practitioners to ef-  
55 ficiently iterate over ideas without having to re-implement the provided UQ methods. To do  
56 so, it provides backbone architecture- and dataset-agnostic implementations of a wide array of  
57 UQ methods and corresponding evaluation schemes for DL, covering regression, classification,  
58 semantic segmentation, and pixel-wise regression tasks.
- 59 • **Adaptability and Expandability:** The modular implementation using `Lightning` encourages  
60 practitioners and the community to an individual adaptation and a continuous expansion and  
61 growth of the toolbox. Additionally, the implementation is adapted to vector or vision data.  
62 Specifically, partial stochasticity [65] is supported when applicable. This supports any larger  
63 architectures used for vision, and the "frozen" functionality enables retraining only a few layers.
- 64 • **Practical and Theoretical Introductions:** The toolbox contains comprehensive practical and  
65 theoretical introductions to the field of UQ and the application of the toolbox. UQ Tutorials and  
66 case studies on designing downstream tasks to compare various UQ methods are provided. A  
67 comprehensive theory guide provides methodological backgrounds on the implemented methods.

68 **Related Work** Frameworks for UQ in DL already exist in the `PyTorch` [54] ecosystem. However,  
69 they are limited to either a handful of UQ methods or a specific class of approaches, such as BDL.  
70 Several libraries exist for BDL, most notably `TorchBNN` [38], `BLiTz` [17], and `Bayesian-Torch` [35].  
71 Yet BNNs are only one approach to UQ and require choosing a prior distribution. When an abundance  
72 of data is available, frequentist procedures, such as conformal prediction, can be a more attractive  
73 alternative. The library `Fortuna` [14] supports several approaches to conformal prediction (CP), of  
74 which we currently support a subset (with plans to incorporate more). The primary difference between  
75 our work and `Fortuna` is that `Fortuna` is primarily compatible with `JAX` [9] and only supports post-  
76 hoc calibration of `PyTorch` models. `TorchCP` [71] is another framework that implements conformal  
77 prediction [4], but it does not support other UQ methods (such as BDL). The most closely related  
78 package to ours is `torch-uncertainty` [36], which implements both frequentist and Bayesian UQ  
79 methods in addition to common benchmarks. Yet our `Lightning UQ Box`, to date, implements the  
80 largest number of UQ methods across different theoretical frameworks, such as BDL and CP, while  
81 including cutting-edge techniques as partially stochastic networks [65], and additionally supports UQ  
82 methods for semantic segmentation tasks. Table 1 gives a comparison with previous libraries.

## 83 2 Benchmarking UQ Methods: the `Lightning UQ Box`

84 The underlying design of `Lightning UQ Box` is based on three pillars:

- 85 • provide a comprehensive set of reference implementations of state-of-the-art UQ methods,
- 86 • optimally fit in the wide open-source landscape for DL based on `PyTorch`, and
- 87 • enhance automation, scalability, and reproducibility of experiments with `Lightning`.

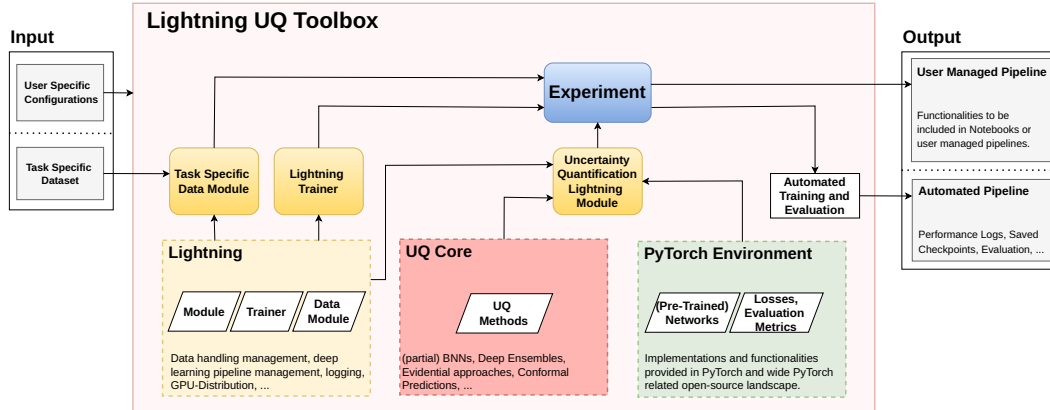


Figure 1: The structure of Lightning UQ Box. The experiments can be built and evaluated at scale or manually tailored to specific use cases. For large experiments at scale, only a dataset and a configuration file have to be provided.

88 These design goals are reflected in the structure of the toolbox, as visualized in Figure 1, and build  
 89 up on the three core components of the available DL functionalities provided within the Lightning  
 90 framework for structuring and pipeline managing, the UQ Core, which contains the UQ method  
 91 implementations, and the PyTorch ecosystem.

92 The UQ Core contains a comprehensive collection of UQ methods for DL with different theoretical  
 93 underpinnings consolidated and implemented for this toolbox. The theoretical backgrounds are very  
 94 diverse and cover, for example, mean-field estimation and various Bayesian-motivated approaches,  
 95 including kernel-based approaches and partially stochastic networks, ensemble methods, and evidence-  
 96 motivated approaches (cmp. Section 2.1). Besides the diversity in methodological approaches, the  
 97 toolbox provides unified interfaces and configuration patterns, thereby improving accessibility and,  
 98 importantly, enabling comparability between the methods.

99 The toolbox is compatible with common DL libraries and frameworks from the PyTorch ecosystem.  
 100 The provided UQ methods can be combined with user-specific architectures and implementations  
 101 provided in the PyTorch ecosystem, including pre-trained networks and foundation models. This is  
 102 especially useful as our framework can build upon or be included in existing code and pipelines based  
 103 on PyTorch-based libraries, such as `timm` [72]. In order to scale BDL to modern-sized architectures,  
 104 we offer functionality to convert existing deterministic architectures, or specified components thereof,  
 105 automatically to a Bayesian framework. As a result, the collection of UQ methods goes beyond  
 106 mere method compilation, offering not only comprehensiveness but also removing time-consuming  
 107 implementation overhead. This enables users to use the UQ toolbox as a simple extension of their  
 108 existing DL pipelines.

109 The toolbox utilizes the Lightning framework to enhance experiment automation, scalability, and re-  
 110 producibility. Lightning offers a flexible and user-friendly interface for the automated management  
 111 of complex pipelines. It is specifically designed to support practitioners in managing experiments  
 112 by providing functionalities to enhance their scalability and reproducibility. These include manag-  
 113 ing configurations, training loops, evaluation steps, and logging processes. To this end, each UQ  
 114 method is implemented as a `LightningModule` that can be used with a `LightningDataModule`  
 115 and a `Trainer` to execute training, evaluation, and inference for a desired task. The toolbox also  
 116 utilizes the Lightning command line interface (CLI) for better experiment reproducibility and for  
 117 setting up experiments at scale. This provides an end-to-end configuration, such that a full pipeline  
 118 of experiments can be built with minimal overhead. Many optional configurations and user-specific  
 119 objects, such as logging functionalities or models, can be included but are not mandatory. The general  
 120 concept of the toolbox is illustrated in Figure 1.

## 121 2.1 Provided Types of UQ Methods

122 Lightning UQ Box provides the most comprehensive collection of the extensive and versatile  
 123 landscape of UQ methods for DL. The following section gives an overview of these different UQ

Table 1: The methods provided with Lightning UQ Box and other available frameworks and reviews. The table represents the status at the time of publication and will be extended in the future. All currently available methods can be found in the provided repository.

	Publication	[26]	[63]	[15]	[31]	[64]	[50]	[46]	[36]	Lightning UQ Box
<b>Deterministic Methods</b>										
Gaussian (MVE)		✓							✓	✓
Deep Evidential Networks (DER)									✓	✓
<b>Neural Network Ensembles</b>										
	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
<b>Bayesian Neural Networks</b>										
MC Dropout (GMM)		✓	✓	✓	✓	✓	✓	✓		✓
BNN with VI ELBO				✓	✓	✓	✓			✓
BNN with VI (alpha divergence)										✓
VBLL										✓
Laplace Approximation					✓					✓
SWAG				✓	✓					✓
DVI, SI				✓						
HMC				✓						
Radial BNN								✓		
Rank-1 BNN								✓		
<b>Gaussian Process based</b>										
Deep Kernel Learning (DKL)										✓
Det. Unc. Estimation (DUE)										✓
Spectral Normalized GPs (SNGP)					✓			✓		✓
<b>Quantile based</b>										
Quantile Regression (QR)	✓		✓							✓
Conformal Prediction (CQR)	✓		✓							✓
<b>Diffusion Model</b>										
CARD										✓
<b>Post-hoc Calibration</b>										
RAPS										✓
TempScaling							✓		✓	✓

124 methods, which are listed in Table 1. For comprehensive explanations, we refer to the theory guide in  
 125 the supplement and to existing reviews [1, 21]. For **regression tasks** NNs predict a continuous target  
 126  $y^*$ . Currently, the toolbox supports six classes of UQ methods for regression: deterministic, quantile,  
 127 ensemble, Bayesian, Gaussian Process, and diffusion-based methods.

- 128 1. Deterministic methods: use a DNN,  $f_\theta : X \rightarrow \mathcal{P}(Y)$ , that map inputs  $x$  to the parameters of  
 129 a probability distribution  $f_\theta(x^*) = p_\theta(x^*) \in \mathcal{P}(Y)$ , and include methods like Deep Evidential  
 130 Regression (**DER**) [2] and Mean Variance Estimation (**MVE**) [49]. The latter, for example, outputs  
 131 the mean and standard deviation of a Gaussian distribution  $f_\theta^{\text{MVE}}(x^*) = (\mu_\theta(x^*), \sigma_\theta(x^*))$ .
- 132 2. Quantile based models: use a DNN,  $f_\theta : X \rightarrow Y^n$ , that map to  $n$  quantiles,  $f_\theta(x^*) =$   
 133  $(q_1(x^*), \dots, q_n(x^*)) \in Y^n$ , and include Quantile Regression [33] (**Quantile Regression**) and the  
 134 conformalized version thereof (**ConformalQR**) [62].
- 135 3. Ensembles: Deep Ensembles [37], which utilize an ensemble over MVE networks.
- 136 4. Bayesian methods: model the network parameters as random variables. Multiple principles and  
 137 techniques to approximate BNNs have been introduced. We include BNNs with Variational  
 138 Inference (VI) (**BNN VI ELBO**) [8], BNNs with VI and alpha divergence (**BNN VI**) [13],  
 139 Variational Bayesian Last Layers (**VBLL**) [28], MC-Dropout (**MCDropout**) [20], the Laplace  
 140 Approximation (**Laplace**) [61][12] and **SWAG** [43] with partially stochastic variants [65].
- 141 5. Gaussian Process-based methods: these model a joint distribution over a set of functions in a  
 142 data-driven manner that approximates the first and second moment of the marginalized distribution.  
 143 These include Deep Kernel Learning (**DKL**) [73], an extension thereof Deterministic Uncertainty  
 144 Estimation (**DUE**) [69, 70], and Spectral Normalized Gaussian Process (**SNGP**) [40].
- 145 6. Conditional Generative model: Classification and Regression Diffusion (**CARD**) [27].

146 For **classification**, the toolbox currently supports six classes of UQ methods. Vanilla softmax proba-  
 147 bilities can be directly used to obtain predictive uncertainties. However, they are often miscalibrated  
 148 and have lead to post-hoc recalibration methods being proposed [25].

- 149 1. Deep Ensembles (**DeepEnsembles**) [37]: utilize an ensemble over independent standard classification networks.
- 150
- 151 2. Bayesian methods: **BNN VI ELBO** [8], **VLL** [28], **MCDropout** [20], **Laplace** [61][12],
- 152 **SWAG** [43].
- 153 3. Gaussian Process based methods: **DKL** [73], **DUE** [69] and Spectral-normalized Neural Gaussian
- 154 Processes (**SNGP**) [40].
- 155 4. Conformal Prediction: [62], Regularized Adaptive Prediction Sets (**RAPS**) [3].
- 156 5. Other: Test-time Augmentation (**TTA**) [41], Temperature Scaling [25].

157 Additionally to the general purpose tasks of regression and classification, **Lightning UQ Box**  
 158 supports UQ methods for vision-specific tasks. These include segmentation and pixel-wise regression,  
 159 where an extensive overview of supported UQ methods can be found on our documentation page.

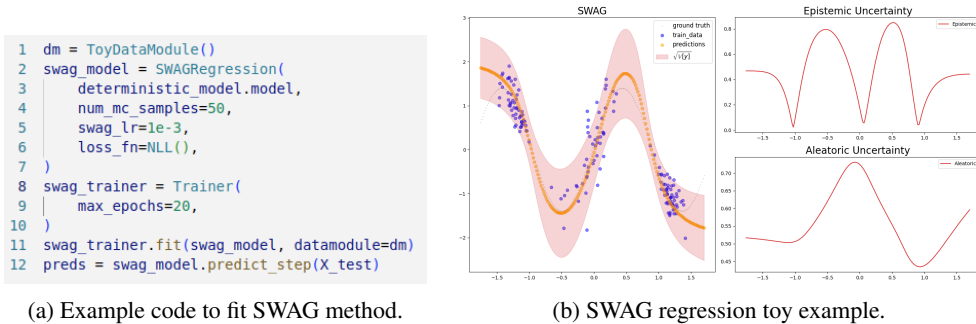


Figure 2: Example code and visualization on toy regression dataset.

160 **Quantifying Predictive Uncertainty:** By default, we quantify predictive uncertainty via the standard  
 161 deviation for regression and via the entropy of the predictive distribution for classification. In general,  
 162 for UQ in DL, two main types of uncertainties can be considered: aleatoric and epistemic [13, 21].  
 163 Aleatoric uncertainty originates from random, or partially observable, effects in the data itself and is  
 164 not reducible: for instance, the Earth covered with clouds does not contain enough information to  
 165 surely assign the land cover type to one of multiple given options. In contrast, epistemic uncertainty  
 166 quantifies the model’s predictive uncertainty originating from uncertainty over its parameters: it will  
 167 typically shrink as more data becomes available [30]. See Figure 2b for an example decomposition.  
 168 Depending on the underlying theoretical assumptions, UQ methods model these types of uncertainties  
 169 individually or mutually [30]. From a statistical perspective, Gruber et al. [24] allude that such a  
 170 distinction is often not possible. Thus, in the examples given here, we focus on the approximate  
 171 predictive distributions of the UQ methods  $p_{\theta}(y_{*}|x_{*})$ , from which we derive the aforementioned  
 172 uncertainty measures. However, where applicable, the toolbox also enables researchers to decompose  
 173 these two types of uncertainties.

174 **Limitations:** Despite the robustness and versatility of the **Lightning UQ Box**, it is tightly inte-  
 175 grated within the PyTorch ecosystem, limiting its applicability to other existing DL frameworks  
 176 like Tensorflow and JAX. Furthermore, merely using UQ methods does not guarantee complete  
 177 reliability, and applications nevertheless require proper experimental design and evaluation.

### 178 3 Experimental Setup for Validation

179 We now showcase **Lightning UQ Box** as a valuable tool for conducting experimental studies  
 180 including benchmarking. We exemplify this by comparing UQ methods on three challenging computer  
 181 vision datasets from two different domains. More concretely, we evaluate the methods on selected  
 182 downstream tasks that highlight the efficacy of UQ and the usefulness of a unified framework<sup>1</sup>. Each  
 183 experiment was completed using the UQ toolbox in less than 10 hours (6 hours on average) on a  
 184 single A100 40GB GPU.

<sup>1</sup>Code for all experiments available at this link: Github Repo.

185 **3.1 Datasets**

186 For our experiments, we consider three datasets: the Tropical Cyclone Driven Data Challenge dataset  
 187 (TC) [44], the Digital Typhoon (DT) dataset [32], and the SKy Images and Photovoltaic Power  
 188 Generation Dataset (SKIPP'D) [47]. An overview of the datasets is given in Table 2. For a detailed  
 explanations of the datasets see supplementary section 1.

Table 2: Dataset Overview.

Dataset	Image source/Satellite	Spatial Res.	Temporal Res.	Train Samples	Val. Samples	Test Samples
Tropical Cyclone	GOES	2 km	15 min	53k	11k	43k
Digital Typhoon	Himawari	5 km	60 min	64.5k	14k	20k
SKIPP'D	Fisheye camera	-	1 min	280k	63k	14k

189

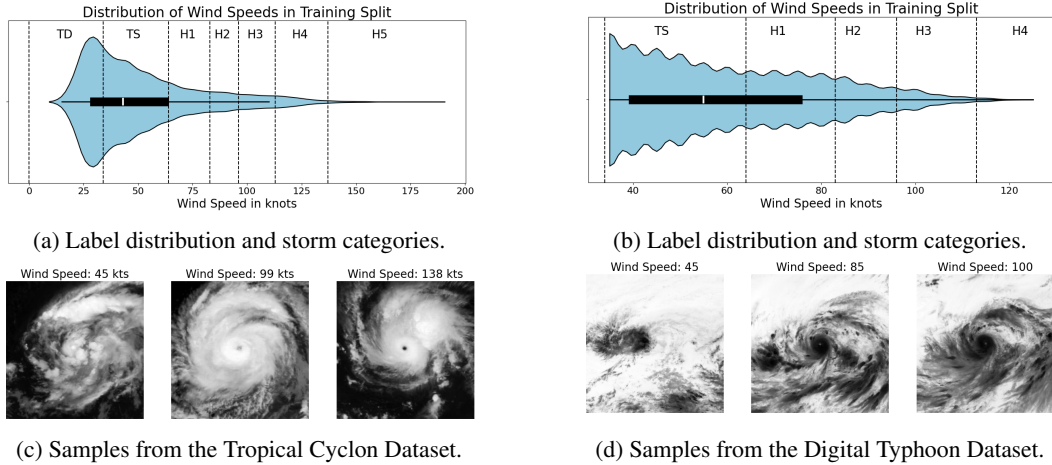


Figure 3: Visualization of the Tropical Cyclone (left) and the Digital Typhoon Dataset (right).

190 **Cyclone and Typhoon Dataset:** The TC and DT datasets consist of infrared measurements that  
 191 capture the spatial structure of storms. Corresponding wind speed targets are matched based on  
 192 hurricane databases. There are varying sources of uncertainty in the inputs, such as missing pixels due  
 193 to the swath of the satellites, and in the targets, such as measurement uncertainties and interpolations  
 194 over time with respect to non-uniform time steps. As such, these datasets exemplify real world  
 195 stochastic phenomena, where predictive uncertainties are essential for decision-making due to the  
 196 inherent risk associated with these potentially extreme events. The magnitude of rapid intensification  
 197 events has been increasing [6], thus causing more damage if not properly detected and predicted. One  
 198 such recent example is Hurricane Otis in October 2023, where existing models had to disproportionately  
 199 rely on satellite data, due to limited in-situ data, which lead to erroneous forecasts [34]. Given the  
 200 extensive availability of satellite imagery, research efforts using this modality are a promising avenue  
 201 to enhance existing forecasts.

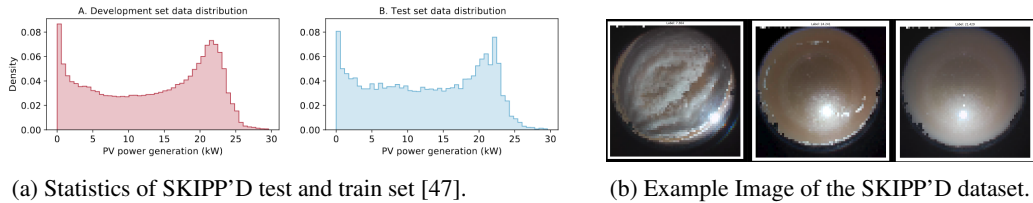


Figure 4: Visualization of SKIPP'D Dataset.

202 **Photovoltaic Dataset:** The SKIPP'D dataset consists of ground-based fish-eye RGB images over  
 203 a 3-year period (2017–2019), where associated targets are power output measurements from a  
 204 30 kilowatt (kW) rooftop photovoltaic array [47]. Given the urgent necessity to transform the world’s  
 205 energy sector to more sustainable solutions [5], this dataset aims to support research efforts of  
 206 large-scale integration of power voltage into electricity grids, where the main problem is to manage  
 207 the non-constant and intermittent power source [47].

208 **3.2 Methodological Setup**

209 **Cyclone and Typhoon Dataset:** Various works have framed the wind speed estimation of tropical  
 210 cyclones from a satellite image as both a regression [10, 42, 75] and classification [57, 74] task.  
 211 We apply all UQ methods provided by the toolbox to the regression and classification task. For  
 212 all wind speed experiments, we use the same ResNet-18 [29] pre-trained on ImageNet<sup>2</sup> as the  
 213 backbone architecture of compared UQ methods. For the TC and DT datasets, the chosen task is  
 214 selective prediction, as introduced by Geifman et al. [22]. Here, samples with a predictive uncertainty  
 215 above a given threshold are omitted and can be referred to domain experts or further decision-making  
 216 pipelines. If the corresponding UQ method has higher uncertainties for inaccurate predictions, leaving  
 217 out the predictions for these samples should increase the overall accuracy, indicating a correlation  
 218 between predictive uncertainty and model error. This can be beneficial in a deployment setting where  
 219 automated analysis systems are paired with human expertise. Examples are visualized in Figure 6.

220 **Photovoltaic Dataset:** Previous work have demonstrated promising results of such image data for  
 221 photovoltaic power generation estimation modeled as a regression task [67, 76, 68, 48, 19, 51, 52].  
 222 We apply all UQ methods provided by the toolbox (see Table 1) to this regression task. Here, we use  
 223 the proposed CNN architecture of Nie et al. [47], which requires only a single line code change in  
 224 experiment configuration for each respective UQ method.<sup>3</sup> Given the central problem of photovoltaics  
 225 being a non-constant power source, we analyze the additional benefits of UQ by evaluating predictive  
 226 uncertainty on annotated sunny and cloudy days. From a reliable model, we expect that both the  
 227 predictive error as well as the predictive uncertainty is larger on the cloudy samples because the  
 228 partial occlusions make it more difficult to estimate the corresponding power voltage output.

229 **Evaluation Metrics:** As evaluation metrics, we use the root mean squared error (RMSE), as well as  
 230 proper scoring rules such as the negative log-likelihood (NLL) [23]. Furthermore, we also consider  
 231 the mean absolute calibration error (MACE) and correlation between the predictive uncertainties and  
 232 mean absolute error (MAE).<sup>4</sup> A detailed description of the employed metrics is in the supplementary.

233 **4 Results: Examples of UQ Method Analysis**

234 The following section provides a quantitative performance comparison of different UQ methods  
 235 under a possible benchmark setting, easily enabled by our proposed framework.

236 **4.1 Selective Prediction for Wind Speed Estimation**

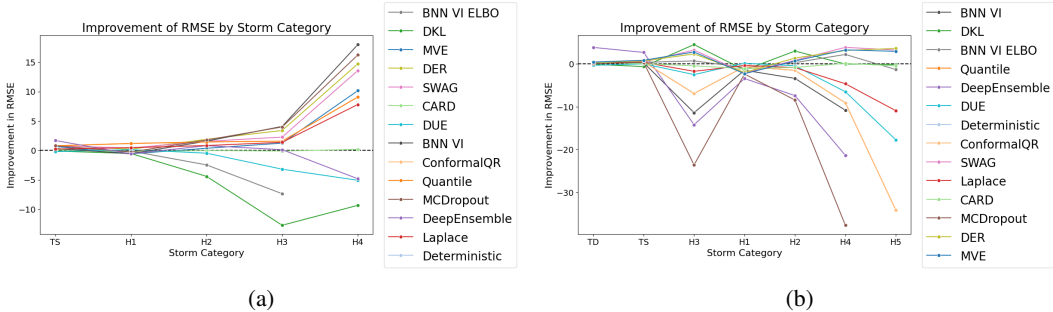


Figure 5: Selective Prediction RMSE improvement per category on the Digital Typhoon Dataset (left) and Tropical Cyclone Dataset (right).

237 Table 3 shows that most UQ methods improve model accuracy when applying selective prediction  
 238 with respect to a deterministic baseline, which cannot express any uncertainty. Compared to Table  
 239 3, Figure 5 demonstrates a different ranking of the UQ methods, with respect to the accuracy  
 240 improvement due to selective prediction, when evaluated per category, according to the Saffir-  
 241 Simpson scale [66]. This ranking also differs on the DT and TC dataset, as Figure 5 shows. The  
 242 skewed data distribution of both datasets, 3a and 3b, and the different uncertainty sources in the

<sup>2</sup>As available in the `timm` library [72]  
<sup>3</sup>More examples are shown in the Github Repo for these experiments.  
<sup>4</sup>Metrics computed with the library provided by [11]



243 TC and DT datasets 3.1 may contribute to these observations of aggregation pathologies. For the  
 244 classification task the ranking of methods varies with Gaussian Process based methods performing  
 245 better, see supplementary section 2.

Table 3: Evaluation of Regression Results on the test set. Note that [64] observe a similar ranking in terms of accuracy, also with respect to Deep Ensembles. RMSE  $\Delta$  shows the improvement after selective prediction, while Coverage denotes the fraction of remaining samples that were not omitted. Selective prediction is based on the 0.8 quantile of predictive uncertainties on a validation set.

UQ group	Method	RMSE $\downarrow$	RMSE $\Delta$ $\uparrow$	NLL $\downarrow$	MACE $\downarrow$
None	Deterministic	9.64	0.00	NaN	NaN
Deterministic	MVE	10.10	0.64	3.74	0.06
	DER	9.59	<b>1.07</b>	4.32	0.30
Quantile	QR	9.54	<b>1.03</b>	<b>3.64</b>	0.05
	CQR	9.54	<b>1.03</b>	3.72	0.10
Ensemble	Deep Ensemble	14.37	0.77	4.05	<b>0.01</b>
	MC Dropout	9.77	<b>1.03</b>	3.75	0.10
	SWAG	<b>9.10</b>	0.97	3.67	0.12
	Laplace	9.64	0.44	3.69	0.03
	BNN VI ELBO	<b>9.15</b>	0.17	15.82	0.35
	BNN VI	10.74	0.94	3.76	0.03
Bayesian	SNGP	9.33	-0.05	14.00	0.36
	VBLL	9.72	0.06	3.70	0.03
	DKL	10.35	-0.31	3.77	<b>0.01</b>
	DUE	9.46	-0.10	3.68	<b>0.01</b>
	Diffusion	CARD	9.57	0.09	9.35

(a) Digital Typhoon Dataset.

UQ group	Method	RMSE $\downarrow$	RMSE $\Delta$ $\uparrow$	NLL $\downarrow$	MACE $\downarrow$
None	Deterministic	10.50	0.00	NaN	NaN
Deterministic	MVE	9.95	1.15	<b>3.64</b>	0.04
	DER	10.14	1.17	4.60	0.35
Quantile	QR	10.95	<b>1.05</b>	3.73	<b>0.01</b>
	CQR	10.95	<b>1.05</b>	3.79	0.10
Ensemble	Deep Ensemble	16.19	<b>3.30</b>	4.15	0.05
	MC Dropout	10.23	0.87	3.81	0.16
	SWAG	9.78	1.13	3.71	0.13
	Laplace	10.53	0.60	4.31	0.28
	BNN VI ELBO	11.82	1.56	5.57	0.23
	BNN VI	11.20	1.45	3.74	0.02
Bayesian	SNGP	12.01	0.28	5.53	0.18
	VBLL	10.79	0.51	3.80	0.07
	DKL	12.59	0.21	3.95	0.06
	DUE	9.95	-0.21	3.73	0.08
	Diffusion	CARD	10.86	0.45	3.92

(b) Tropical Cyclone Dataset.

246 Figure 6 gives a visual intuition of the selective prediction scheme. If the predictive uncertainty  
 247 (red-shaded region) exceeds the established threshold (blue-shaded region), individual predictions are  
 248 deferred to an expert. The models provide a reasonable mean estimate of a storm track, even though  
 249 predictions are made for single image instances and the regression task is modeled by ResNet-18  
 250 without a notion of time. Figure 6 additionally showcases the effect of conformalizing the predictive  
 251 uncertainty of an underlying model. Conformal prediction aims to calibrate prediction sets while  
 252 providing theoretical coverage guarantees; it can be particularly interesting in the case of a high-risk  
 253 task such as wind speed estimation. Figure 6b demonstrates the effectiveness of the procedure, as the  
 254 coverage has increased from 0.73 to 0.97, which is also reflected in the wider prediction intervals that  
 255 cover the targets without sacrificing any accuracy.

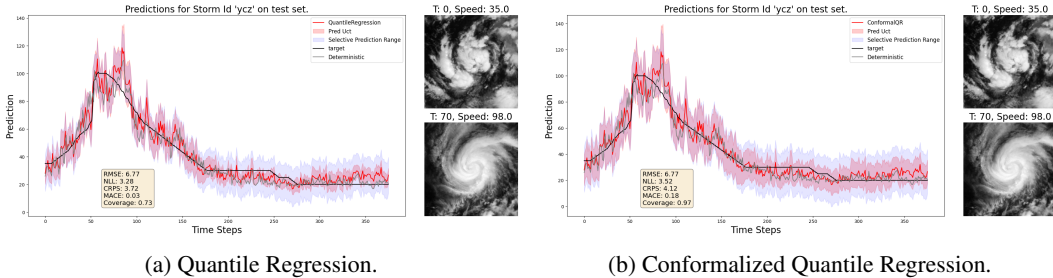


Figure 6: Individual nowcasting predictions are stitched together to recover a time series. Areas where the red-shaded regions exceed the blue denote samples that *would* be omitted during selective prediction. The example showcases the effect of the conformal procedure, where conformalized prediction intervals increase the desired empirical coverage.

## 256 4.2 Photovoltaic Power Output Estimation Under Cloudy and Sunny Conditions

257 Figure 7 demonstrates that model performance differs under cloudy or sunny conditions. Across  
 258 methods the NLL demonstrates differences in the model performance and related calibration between  
 259 cloudy and sunny days. The consideration of uncertainty improved the accuracy of models compared  
 260 to the deterministic baseline, as shown in the supplementary material. The correlation between the  
 261 model error (in terms of MAE) and the predictive uncertainty shows a clear positive correlation  
 262 ( $>0.45$ ) across all methods. However, there are differences in the magnitude between methods and  
 263 cloud conditions. Stakeholders might prefer good UQ estimates on more complex days, i.e., the  
 264 cloudy ones, than for sunny days, where the output is much easier to predict. Exhaustive results can  
 265 be found in the supplementary material.



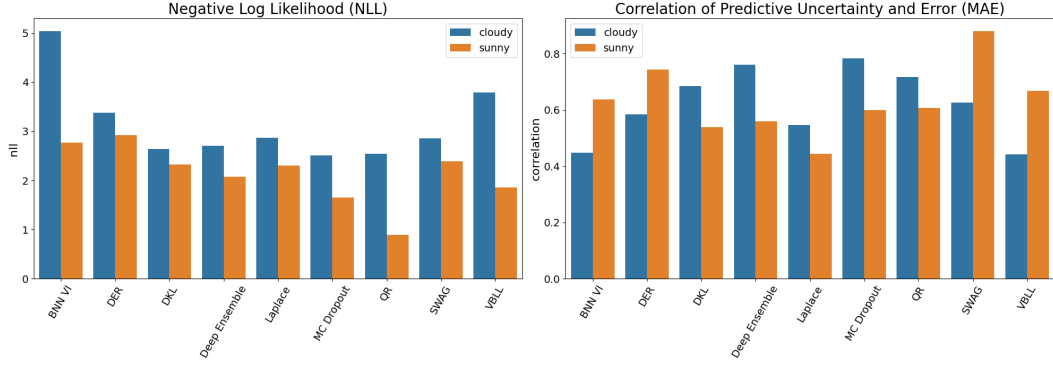
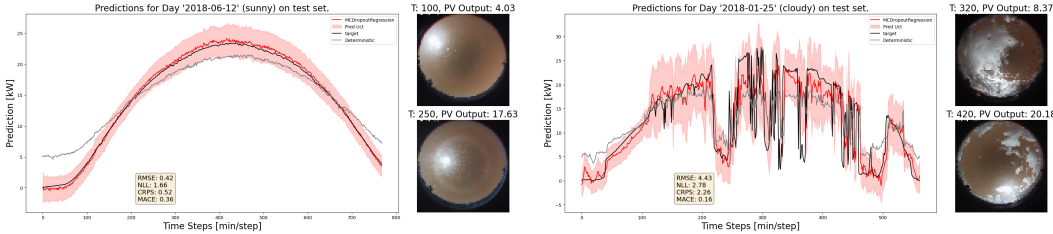


Figure 7: Negative Log Likelihood (left) and correlation between model error (measured by MAE) and predictive uncertainty for different methods on cloudy and sunny test examples.

266 Figure 8 showcases concrete examples with power voltage estimates plotted over the duration of a  
 267 cloudy and a non-cloudy day. Compared to the smooth and consistent power output on a sunny day  
 268 8a, the predictive uncertainty is larger under cloudy conditions. This may reflect the uncertainty in  
 269 the input images due to cloudiness changing faster than the time step resolution.



(a) MC Dropout prediction: sunny day example. (b) MC Dropout prediction: cloudy day example.

Figure 8: Individual nowcasting predictions stitched together to recover a time series. The plot shows qualitative and quantitative differences between the two methods for the same set of predictions.

## 270 5 Conclusion

271 We have introduced Lightning UQ Box, a comprehensive framework for enhancing neural networks  
 272 with uncertainty estimates. Additionally, we have showcased its usefulness for comparing a broad  
 273 range of methods from different theoretical foundations on three relevant tasks with various sources  
 274 of uncertainty. Our framework not only makes it easier for practitioners to use Bayesian methods for  
 275 DL as demanded by [53] but goes beyond this by supporting UQ methods stemming from various  
 276 theoretical frameworks and assumptions. Our experimental results demonstrate the differences and  
 277 variability between UQ methods and, therefore, the benefit of this benchmarking framework. In  
 278 conclusion, our open-source framework and the accompanying resources can be both an entry point  
 279 for researchers to the field of UQ and also aid the development of new methods that address the  
 280 shortcomings of existing ones [50].

## 281 6 Ethics and Broader Impact Statement

282 Including UQ in DL applied to real world and safety critical applications is of significant importance.  
 283 UQ can provide the means to reduce risks, yet practitioners should not succumb to a false sense  
 284 of security provided by such methods. The performance and reliability of UQ methods may be  
 285 dataset and task dependent. Exactly for that reason we provide our framework under the open-source  
 286 Apache-2.0 license to support open science, transparency, and collaborative research efforts.

287 **References**

- 288 [1] Moloud Abdar, Farhad Pourpanah, Sadiq Hussain, Dana Rezazadegan, Li Liu, Mohammad  
289 Ghavamzadeh, Paul Fieguth, Xiaochun Cao, Abbas Khosravi, U Rajendra Acharya, et al. A  
290 review of uncertainty quantification in deep learning: Techniques, applications and challenges.  
291 *Information fusion*, 76:243–297, 2021.
- 292 [2] Alexander Amini, Wilko Schwarting, Ava Soleimany, and Daniela Rus. Deep evidential  
293 regression. *Advances in Neural Information Processing Systems*, 33:14927–14937, 2020.
- 294 [3] Anastasios Angelopoulos, Stephen Bates, Jitendra Malik, and Michael I Jordan. Uncertainty  
295 sets for image classifiers using conformal prediction. *arXiv preprint arXiv:2009.14193*, 2020.
- 296 [4] Anastasios N Angelopoulos and Stephen Bates. A gentle introduction to conformal prediction  
297 and distribution-free uncertainty quantification. *arXiv preprint arXiv:2107.07511*, 2021.
- 298 [5] Svante Arrhenius. Xxxi. on the influence of carbonic acid in the air upon the temperature of the  
299 ground. *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science*,  
300 41(251):237–276, 1896.
- 301 [6] Karthik Balaguru, Gregory R Foltz, and L Ruby Leung. Increasing magnitude of hurricane  
302 rapid intensification in the central and eastern tropical atlantic. *Geophysical Research Letters*,  
303 45(9):4238–4247, 2018.
- 304 [7] Roberto Bentivoglio, Elvin Isufi, Sebastian Nicolaas Jonkman, and Riccardo Taormina. Deep  
305 learning methods for flood mapping: a review of existing applications and future research  
306 directions. *Hydrology and earth system sciences*, 26(16):4345–4378, 2022.
- 307 [8] Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty  
308 in neural network. In *International conference on machine learning*, pages 1613–1622. PMLR,  
309 2015.
- 310 [9] James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal  
311 Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao  
312 Zhang. JAX: Composable transformations of Python+NumPy programs, 2018.
- 313 [10] Buo-Fu Chen, Boyo Chen, Hsuan-Tien Lin, and Russell L Elsberry. Estimating tropical cyclone  
314 intensity by satellite imagery utilizing convolutional neural networks. *Weather and Forecasting*,  
315 34(2):447–465, 2019.
- 316 [11] Youngseog Chung, Ian Char, Han Guo, Jeff Schneider, and Willie Neiswanger. Uncertainty tool-  
317 box: an open-source library for assessing, visualizing, and improving uncertainty quantification.  
318 *arXiv preprint arXiv:2109.10254*, 2021.
- 319 [12] Erik Daxberger, Agustinus Kristiadi, Alexander Immer, Runa Eschenhagen, Matthias Bauer,  
320 and Philipp Hennig. Laplace redux-effortless bayesian deep learning. *Advances in Neural  
321 Information Processing Systems*, 34:20089–20103, 2021.
- 322 [13] Stefan Depeweg, Jose-Miguel Hernandez-Lobato, Finale Doshi-Velez, and Steffen Udluft.  
323 Decomposition of uncertainty in bayesian deep learning for efficient and risk-sensitive learning.  
324 In *International Conference on Machine Learning*, pages 1184–1193. PMLR, 2018.
- 325 [14] Gianluca Detommaso, Alberto Gasparin, Michele Donini, Matthias Seeger, Andrew Gordon  
326 Wilson, and Cedric Archambeau. Fortuna: A library for uncertainty quantification in deep  
327 learning. *arXiv preprint arXiv:2302.04019*, 2023.
- 328 [15] Nicolas Dewolf, Bernard De Baets, and Willem Waegeman. Valid prediction intervals for  
329 regression problems. *Artificial Intelligence Review*, pages 1–37, 2022.
- 330 [16] Codrut-Andrei Diaconu and Nina Maria Gottschling. Uncertainty-aware learning with label  
331 noise for glacier mass balance modelling. *IEEE Geoscience and Remote Sensing Letters*, 2024.
- 332 [17] Piero Esposito. BLiTz - Bayesian Layers in Torch Zoo (a Bayesian deep learning library for  
333 Torch). <https://github.com/piEsposito/blitz-bayesian-deep-learning/>, 2020.

- 334 [18] William A. Falcon. PyTorch Lightning. [https://github.com/Lightning-AI/](https://github.com/Lightning-AI/pytorch-lightning)  
335 [pytorch-lightning](https://github.com/Lightning-AI/pytorch-lightning), 2019.
- 336 [19] Cong Feng and Jie Zhang. Solarnet: A sky image-based deep convolutional neural network for  
337 intra-hour solar forecasting. *Solar Energy*, 204:71–78, 2020.
- 338 [20] Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model  
339 uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059.  
340 PMLR, 2016.
- 341 [21] Jakob Gawlikowski, Cedrique Rovile Njiteucheu Tassi, Mohsin Ali, Jongseok Lee, Matthias  
342 Humt, Jianxiang Feng, Anna Kruspe, Rudolph Triebel, Peter Jung, Ribana Roscher, et al.  
343 A survey of uncertainty in deep neural networks. *Artificial Intelligence Review*, 56(Suppl  
344 1):1513–1589, 2023.
- 345 [22] Yonatan Geifman and Ran El-Yaniv. Selective classification for deep neural networks. *Advances*  
346 *in neural information processing systems*, 30, 2017.
- 347 [23] Tilmann Gneiting and Adrian E Raftery. Strictly proper scoring rules, prediction, and estimation.  
348 *Journal of the American statistical Association*, 102(477):359–378, 2007.
- 349 [24] Cornelia Gruber, Patrick Oliver Schenk, Malte Schierholz, Frauke Kreuter, and Göran Kauer-  
350 mann. Sources of uncertainty in machine learning—a statisticians’ view. *arXiv preprint*  
351 *arXiv:2305.16703*, 2023.
- 352 [25] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural  
353 networks. In *International conference on machine learning*, pages 1321–1330. PMLR, 2017.
- 354 [26] Fredrik K. Gustafsson, Martin Danelljan, and Thomas B. Schön. How reliable is your regression  
355 model’s uncertainty under real-world distribution shifts?, 2023.
- 356 [27] Xizewen Han, Huangjie Zheng, and Mingyuan Zhou. Card: Classification and regression  
357 diffusion models. *Advances in Neural Information Processing Systems*, 35:18100–18115, 2022.
- 358 [28] James Harrison, John Willes, and Jasper Snoek. Variational bayesian last layers. In *International*  
359 *Conference on Learning Representations (ICLR)*, 2024.
- 360 [29] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image  
361 recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*,  
362 pages 770–778, 2016.
- 363 [30] Eyke Hüllermeier and Willem Waegeman. Aleatoric and epistemic uncertainty in machine  
364 learning: An introduction to concepts and methods. *Machine Learning*, 110:457–506, 2021.
- 365 [31] Pavel Izmailov, Sharad Vikram, Matthew D Hoffman, and Andrew Gordon Gordon Wilson.  
366 What are bayesian neural network posteriors really like? In *International conference on machine*  
367 *learning*, pages 4629–4640. PMLR, 2021.
- 368 [32] Asanobu Kitamoto, Jared Hwang, Bastien Vuillod, Lucas Gautier, Yingtao Tian, and Tarin  
369 Clanuwat. Digital typhoon: Long-term satellite image dataset for the spatio-temporal modeling  
370 of tropical cyclones. *Advances in Neural Information Processing Systems*, 36, 2024.
- 371 [33] Roger Koenker and Gilbert Bassett Jr. Regression quantiles. *Econometrica: Journal of the*  
372 *Econometric Society*, pages 33–50, 1978.
- 373 [34] Katrina Krämer. Daily briefing: Why forecasters failed to predict hurricane otis. *Nature*, 2023.
- 374 [35] Ranganath Krishnan, Pi Esposito, and Mahesh Subedar. Bayesian-torch: Bayesian neural  
375 network layers for uncertainty estimation, January 2022.
- 376 [36] Adrian Lafage and Olivier Laurent. Torch Uncertainty. [https://github.com/](https://github.com/ENSTA-U2IS-AI/torch-uncertainty)  
377 [ENSTA-U2IS-AI/torch-uncertainty](https://github.com/ENSTA-U2IS-AI/torch-uncertainty), 2024.

- 378 [37] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable  
379 predictive uncertainty estimation using deep ensembles. *Advances in neural information*  
380 *processing systems*, 30, 2017.
- 381 [38] Sungyoon Lee, Hoki Kim, and Jaewook Lee. GradDiv: Adversarial robustness of randomized  
382 neural networks via gradient diversity regularization. *IEEE Transactions on Pattern Analysis*  
383 *and Machine Intelligence*, 2022.
- 384 [39] Nils Lehmann, Nina Maria Gottschling, Stefan Depeweg, and Eric Nalisnick. Uncertainty aware  
385 tropical cyclone wind speed estimation from satellite data. *arXiv preprint arXiv:2404.08325*,  
386 2024.
- 387 [40] Jeremiah Liu, Zi Lin, Shreyas Padhy, Dustin Tran, Tania Bedrax Weiss, and Balaji Lakshmi-  
388 narayanan. Simple and principled uncertainty estimation with deterministic deep learning via  
389 distance awareness. *Advances in Neural Information Processing Systems*, 33:7498–7512, 2020.
- 390 [41] Alexander Lyzhov, Yuliya Molchanova, Arsenii Ashukha, Dmitry Molchanov, and Dmitry  
391 Vetrov. Greedy policy search: A simple baseline for learnable test-time augmentation. In  
392 *Conference on Uncertainty in Artificial Intelligence*, pages 1308–1317. PMLR, 2020.
- 393 [42] Zhaoyang Ma, Yunfeng Yan, Jianmin Lin, and Dongfang Ma. A multi-scale and multi-layer  
394 feature extraction network with dual attention for tropical cyclone intensity estimation. *IEEE*  
395 *Transactions on Geoscience and Remote Sensing*, 2024.
- 396 [43] Wesley J Maddox, Pavel Izmailov, Timur Garipov, Dmitry P Vetrov, and Andrew Gordon  
397 Wilson. A simple baseline for bayesian uncertainty in deep learning. *Advances in neural*  
398 *information processing systems*, 32, 2019.
- 399 [44] M. Maskey, R. Ramachandran, I. Gurung, B. Freitag, M. Ramasubramanian, and J. Miller.  
400 Tropical Cyclone Wind Estimation Competition Dataset. [https://doi.org/10.34911/  
401 rdnt.xs53up](https://doi.org/10.34911/rdnt.xs53up), 2021.
- 402 [45] Riccardo Miotto, Fei Wang, Shuang Wang, Xiaoqian Jiang, and Joel T Dudley. Deep learning for  
403 healthcare: review, opportunities and challenges. *Briefings in bioinformatics*, 19(6):1236–1246,  
404 2018.
- 405 [46] Zachary Nado, Neil Band, Mark Collier, Josip Djolonga, Michael W Dusenberry, Sebastian  
406 Farquhar, Qixuan Feng, Angelos Filos, Marton Havasi, Rodolphe Jenatton, et al. Uncer-  
407 tainty baselines: Benchmarks for uncertainty & robustness in deep learning. *arXiv preprint*  
408 *arXiv:2106.04015*, 2021.
- 409 [47] Yuhao Nie, Xiatong Li, Andea Scott, Yuchi Sun, Vignesh Venugopal, and Adam Brandt. Skipp’d:  
410 A sky images and photovoltaic power generation dataset for short-term solar forecasting. *Solar*  
411 *Energy*, 255:171–179, 2023.
- 412 [48] Yuhao Nie, Yuchi Sun, Yuanlei Chen, Rachel Orsini, and Adam Brandt. Pv power output  
413 prediction from sky images using convolutional neural network: The comparison of sky-  
414 condition-specific sub-models and an end-to-end model. *Journal of Renewable and Sustainable*  
415 *Energy*, 12(4), 2020.
- 416 [49] David A Nix and Andreas S Weigend. Estimating the mean and variance of the target proba-  
417 bility distribution. In *Proceedings of 1994 ieee international conference on neural networks*  
418 *(ICNN’94)*, volume 1, pages 55–60. IEEE, 1994.
- 419 [50] Yaniv Ovadia, Emily Fertig, Jie Ren, Zachary Nado, David Sculley, Sebastian Nowozin, Joshua  
420 Dillon, Balaji Lakshminarayanan, and Jasper Snoek. Can you trust your model’s uncertainty?  
421 evaluating predictive uncertainty under dataset shift. *Advances in neural information processing*  
422 *systems*, 32, 2019.
- 423 [51] Quentin Paletta, Guillaume Arbod, and Joan Lasenby. Benchmarking of deep learning irradiance  
424 forecasting models from sky images—an in-depth analysis. *Solar Energy*, 224:855–867, 2021.
- 425 [52] Quentin Paletta, Anthony Hu, Guillaume Arbod, and Joan Lasenby. Eclipse: Envisioning cloud  
426 induced perturbations in solar energy. *Applied Energy*, 326:119924, 2022.

- 427 [53] Theodore Papamarkou, Maria Skoularidou, Konstantina Palla, Laurence Aitchison, Julyan  
428 Arbel, David Dunson, Maurizio Filippone, Vincent Fortuin, Philipp Hennig, Aliaksandr Hubin,  
429 et al. Position paper: Bayesian deep learning in the age of large-scale AI. *arXiv preprint*  
430 *arXiv:2402.00809*, 2024.
- 431 [54] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,  
432 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative  
433 style, high-performance deep learning library. *Advances in neural information processing*  
434 *systems*, 32, 2019.
- 435 [55] Claudio Persello, Jan Dirk Wegner, Ronny Hänsch, Devis Tuia, Pedram Ghamisi, Mila Koeva,  
436 and Gustau Camps-Valls. Deep learning and earth observation to support the sustainable  
437 development goals: Current approaches, open challenges, and future opportunities. *IEEE*  
438 *Geoscience and Remote Sensing Magazine*, 10(2):172–200, 2022.
- 439 [56] Harry A Pierson and Michael S Gashler. Deep learning in robotics: a review of recent research.  
440 *Advanced Robotics*, 31(16):821–835, 2017.
- 441 [57] Ritesh Pradhan, Ramazan S Aygun, Manil Maskey, Rahul Ramachandran, and Daniel J Ce-  
442 cil. Tropical cyclone intensity estimation using a deep convolutional neural network. *IEEE*  
443 *Transactions on Image Processing*, 27(2):692–702, 2017.
- 444 [58] David Radke, Anna Hessler, and Dan Ellsworth. Firecast: Leveraging deep learning to predict  
445 wildfire spread. In *IJCAI*, pages 4575–4581, 2019.
- 446 [59] Stephan Rasp, Peter D Dueben, Sebastian Scher, Jonathan A Weyn, Soukayna Mouatadid, and  
447 Nils Thuerey. Weatherbench: a benchmark data set for data-driven weather forecasting. *Journal*  
448 *of Advances in Modeling Earth Systems*, 12(11):e2020MS002203, 2020.
- 449 [60] Markus Reichstein, Gustau Camps-Valls, Bjorn Stevens, Martin Jung, Joachim Denzler, Nuno  
450 Carvalhais, and fnm Prabhat. Deep learning and process understanding for data-driven earth  
451 system science. *Nature*, 566(7743):195–204, 2019.
- 452 [61] Hippolyt Ritter, Aleksandar Botev, and David Barber. A scalable laplace approximation for  
453 neural networks. In *6th International Conference on Learning Representations, ICLR 2018-*  
454 *Conference Track Proceedings*, volume 6. International Conference on Representation Learning,  
455 2018.
- 456 [62] Yaniv Romano, Evan Patterson, and Emmanuel Candes. Conformalized quantile regression.  
457 *Advances in neural information processing systems*, 32, 2019.
- 458 [63] Franko Schmähling, Jörg Martin, and Clemens Elster. A framework for benchmarking uncer-  
459 tainty in deep regression. *Applied Intelligence*, pages 1–14, 2022.
- 460 [64] Florian Seligmann, Philipp Becker, Michael Volpp, and Gerhard Neumann. Beyond deep  
461 ensembles: A large-scale evaluation of bayesian deep learning under distribution shift. *Advances*  
462 *in Neural Information Processing Systems*, 36, 2024.
- 463 [65] Mrinank Sharma, Sebastian Farquhar, Eric Nalisnick, and Tom Rainforth. Do bayesian neural  
464 networks need to be fully stochastic? In *International Conference on Artificial Intelligence and*  
465 *Statistics*, pages 7694–7722. PMLR, 2023.
- 466 [66] Robert H Simpson. The hurricane disaster—potential scale. *Weatherwise*, 27(4):169–186, 1974.
- 467 [67] Yuchi Sun, Gergely Szűcs, and Adam R Brandt. Solar pv output prediction from video streams  
468 using convolutional neural networks. *Energy & Environmental Science*, 11(7):1811–1818, 2018.
- 469 [68] Yuchi Sun, Vignesh Venugopal, and Adam R Brandt. Short-term solar power forecast with deep  
470 learning: Exploring optimal input and output configuration. *Solar Energy*, 188:730–741, 2019.
- 471 [69] Joost van Amersfoort, Lewis Smith, Andrew Jesson, Oscar Key, and Yarin Gal. On fea-  
472 ture collapse and deep kernel learning for single forward pass uncertainty. *arXiv preprint*  
473 *arXiv:2102.11409*, 2021.

- 474 [70] Joost Van Amersfoort, Lewis Smith, Yee Whye Teh, and Yarin Gal. Uncertainty estimation  
475 using a single deep deterministic neural network. In *International Conference on Machine*  
476 *Learning*, pages 9690–9700. PMLR, 2020.
- 477 [71] Hongxin Wei and Jianguo Huang. TorchCP: A library for conformal prediction based on  
478 PyTorch, 2024.
- 479 [72] Ross Wightman. PyTorch Image Models. [https://github.com/rwightman/  
480 pytorch-image-models](https://github.com/rwightman/pytorch-image-models), 2019.
- 481 [73] Andrew Gordon Wilson, Zhiting Hu, Ruslan Salakhutdinov, and Eric P Xing. Deep kernel  
482 learning. In *Artificial intelligence and statistics*, pages 370–378. PMLR, 2016.
- 483 [74] Anthony Wimmers, Christopher Velden, and Joshua H Cossuth. Using deep learning to estimate  
484 tropical cyclone intensity from satellite passive microwave imagery. *Monthly Weather Review*,  
485 147(6):2261–2282, 2019.
- 486 [75] Chang-Jiang Zhang, Xiao-Jie Wang, Lei-Ming Ma, and Xiao-Qin Lu. Tropical cyclone intensity  
487 classification and estimation using infrared satellite images with deep learning. *IEEE Journal*  
488 *of Selected Topics in Applied Earth Observations and Remote Sensing*, 14:2070–2086, 2021.
- 489 [76] Jinsong Zhang, Rodrigo Verschae, Shohei Nobuhara, and Jean-François Lalonde. Deep photo-  
490 voltaic nowcasting. *Solar Energy*, 176:267–276, 2018.

491 **NeurIPS Paper Checklist**

492 **1. Claims**

493 Question: Do the main claims made in the abstract and introduction accurately reflect the  
494 paper's contributions and scope?

495 Answer: [Yes]

496 Justification: We introduce the toolbox in Section 2 and analyze its usefulness on three  
497 example datasets in Section 4.

498 Guidelines:

- 499 • The answer NA means that the abstract and introduction do not include the claims  
500 made in the paper.
- 501 • The abstract and/or introduction should clearly state the claims made, including the  
502 contributions made in the paper and important assumptions and limitations. A No or  
503 NA answer to this question will not be perceived well by the reviewers.
- 504 • The claims made should match theoretical and experimental results, and reflect how  
505 much the results can be expected to generalize to other settings.
- 506 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
507 are not attained by the paper.

508 **2. Limitations**

509 Question: Does the paper discuss the limitations of the work performed by the authors?

510 Answer: [Yes]

511 Justification: We discuss limitations, which is the embedding into the PyTorch framework,  
512 within a paragraph of Section 2.

513 Guidelines:

- 514 • The answer NA means that the paper has no limitations, while the answer No means  
515 that the paper has limitations, but those are not discussed in the paper.
- 516 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 517 • The paper should point out any strong assumptions and how robust the results are to  
518 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
519 model well-specification, asymptotic approximations only holding locally). The authors  
520 should reflect on how these assumptions might be violated in practice and what the  
521 implications would be.
- 522 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
523 only tested on a few datasets or with a few runs. In general, empirical results often  
524 depend on implicit assumptions, which should be articulated.
- 525 • The authors should reflect on the factors that influence the performance of the approach.  
526 For example, a facial recognition algorithm may perform poorly when image resolution  
527 is low or images are taken in low lighting. Or a speech-to-text system might not be  
528 used reliably to provide closed captions for online lectures because it fails to handle  
529 technical jargon.
- 530 • The authors should discuss the computational efficiency of the proposed algorithms  
531 and how they scale with dataset size.
- 532 • If applicable, the authors should discuss possible limitations of their approach to  
533 address problems of privacy and fairness.
- 534 • While the authors might fear that complete honesty about limitations might be used by  
535 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
536 limitations that aren't acknowledged in the paper. The authors should use their best  
537 judgment and recognize that individual actions in favor of transparency play an impor-  
538 tant role in developing norms that preserve the integrity of the community. Reviewers  
539 will be specifically instructed to not penalize honesty concerning limitations.

540 **3. Theory Assumptions and Proofs**

541 Question: For each theoretical result, does the paper provide the full set of assumptions and  
542 a complete (and correct) proof?



543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593  
594  
595  
596

Answer: [NA]

Justification: The paper provides a toolbox and does not include theoretical results.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

**4. Experimental Result Reproducibility**

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: All experiments are reproducible with the presented toolbox and the provided code. Further, we describe the experimental setups in Section 3 and in the supplement and reference related works.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
  - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

**5. Open access to data and code**

597 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
598 tions to faithfully reproduce the main experimental results, as described in supplemental  
599 material?

600 Answer: [Yes]

601 Justification: The whole toolbox is under Apache-2.0 license. The full code for the presented  
602 example experiments, utilizing the toolbox, is provided together with instructions and  
603 explanations: <https://github.com/lightning-uq-box/experiments>.

604 Guidelines:

- 605 • The answer NA means that paper does not include experiments requiring code.
- 606 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/  
607 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 608 • While we encourage the release of code and data, we understand that this might not be  
609 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not  
610 including code, unless this is central to the contribution (e.g., for a new open-source  
611 benchmark).
- 612 • The instructions should contain the exact command and environment needed to run to  
613 reproduce the results. See the NeurIPS code and data submission guidelines ([https://  
614 nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 615 • The authors should provide instructions on data access and preparation, including how  
616 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 617 • The authors should provide scripts to reproduce all experimental results for the new  
618 proposed method and baselines. If only a subset of experiments are reproducible, they  
619 should state which ones are omitted from the script and why.
- 620 • At submission time, to preserve anonymity, the authors should release anonymized  
621 versions (if applicable).
- 622 • Providing as much information as possible in supplemental material (appended to the  
623 paper) is recommended, but including URLs to data and code is permitted.

## 624 6. Experimental Setting/Details

625 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-  
626 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the  
627 results?

628 Answer: [Yes]

629 Justification: The paper mentions experimental setups that is needed to understand the  
630 presented results and further references to works on which the experimental setup builds  
631 upon.

632 Guidelines:

- 633 • The answer NA means that the paper does not include experiments.
- 634 • The experimental setting should be presented in the core of the paper to a level of detail  
635 that is necessary to appreciate the results and make sense of them.
- 636 • The full details can be provided either with the code, in appendix, or as supplemental  
637 material.

## 638 7. Experiment Statistical Significance

639 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
640 information about the statistical significance of the experiments?

641 Answer: [NA]

642 Justification: The experiments are utilized to represent the usability and potential advantages  
643 of the toolbox.

644 Guidelines:

- 645 • The answer NA means that the paper does not include experiments.
- 646 • The authors should answer "Yes" if the results are accompanied by error bars, confi-  
647 dence intervals, or statistical significance tests, at least for the experiments that support  
648 the main claims of the paper.

- 649 • The factors of variability that the error bars are capturing should be clearly stated (for  
650 example, train/test split, initialization, random drawing of some parameter, or overall  
651 run with given experimental conditions).
- 652 • The method for calculating the error bars should be explained (closed form formula,  
653 call to a library function, bootstrap, etc.)
- 654 • The assumptions made should be given (e.g., Normally distributed errors).
- 655 • It should be clear whether the error bar is the standard deviation or the standard error  
656 of the mean.
- 657 • It is OK to report 1-sigma error bars, but one should state it. The authors should  
658 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
659 of Normality of errors is not verified.
- 660 • For asymmetric distributions, the authors should be careful not to show in tables or  
661 figures symmetric error bars that would yield results that are out of range (e.g. negative  
662 error rates).
- 663 • If error bars are reported in tables or plots, The authors should explain in the text how  
664 they were calculated and reference the corresponding figures or tables in the text.

## 665 8. Experiments Compute Resources

666 Question: For each experiment, does the paper provide sufficient information on the com-  
667 puter resources (type of compute workers, memory, time of execution) needed to reproduce  
668 the experiments?

669 Answer: [Yes]

670 Justification: We stated the resources (Nvidia A100 GPU 40GB) and the computation time  
671 for all experiments of less than 10 hours when automated run with the UQ toolbox.

672 Guidelines:

- 673 • The answer NA means that the paper does not include experiments.
- 674 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,  
675 or cloud provider, including relevant memory and storage.
- 676 • The paper should provide the amount of compute required for each of the individual  
677 experimental runs as well as estimate the total compute.
- 678 • The paper should disclose whether the full research project required more compute  
679 than the experiments reported in the paper (e.g., preliminary or failed experiments that  
680 didn't make it into the paper).

## 681 9. Code Of Ethics

682 Question: Does the research conducted in the paper conform, in every respect, with the  
683 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

684 Answer: [Yes]

685 Justification: We do not see any potential harm caused by the research process and no  
686 negative societal and potentially harmful consequences.

687 Guidelines:

- 688 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 689 • If the authors answer No, they should explain the special circumstances that require a  
690 deviation from the Code of Ethics.
- 691 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-  
692 eration due to laws or regulations in their jurisdiction).

## 693 10. Broader Impacts

694 Question: Does the paper discuss both potential positive societal impacts and negative  
695 societal impacts of the work performed?

696 Answer: [Yes]

697 Justification: We do not see negative societal impacts and point out the positive impact of  
698 open-source uncertainty quantification frameworks in supporting the work on topics with  
699 positive social impact.

700  
701  
702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752

**Guidelines:**

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

**11. Safeguards**

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: We do not see the potential for direct misuse.

**Guidelines:**

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example, by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

**12. Licenses for existing assets**

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We reference the used data sets and experimental setups from other works. We further reference the frameworks on which our toolbox builds up, Lightning and PyTorch.

**Guidelines:**

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.

- 753
- 754
- 755
- 756
- 757
- 758
- 759
- 760
- 761
- 762
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
  - If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
  - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
  - If this information is not available online, the authors are encouraged to reach out to the asset's creators.

763 **13. New Assets**

764 Question: Are new assets introduced in the paper well documented and is the documentation  
765 provided alongside the assets?

766 Answer: [Yes]

767 Justification: The asset is given by the Lightning UQ Box that is linked and fully open-  
768 source. For the experiments there is further code for reproduction of the experiments  
769 provided.

770 Guidelines:

- 771
- 772
- 773
- 774
- 775
- 776
- 777
- 778
- The answer NA means that the paper does not release new assets.
  - Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
  - The paper should discuss whether and how consent was obtained from people whose asset is used.
  - At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

779 **14. Crowdsourcing and Research with Human Subjects**

780 Question: For crowdsourcing experiments and research with human subjects, does the paper  
781 include the full text of instructions given to participants and screenshots, if applicable, as  
782 well as details about compensation (if any)?

783 Answer: [NA]

784 Justification: The work does not contain crowdsourcing experiments and research with  
785 human subjects.

786 Guidelines:

- 787
- 788
- 789
- 790
- 791
- 792
- 793
- 794
- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
  - Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
  - According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

795 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human  
796 Subjects**

797 Question: Does the paper describe potential risks incurred by study participants, whether  
798 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
799 approvals (or an equivalent approval/review based on the requirements of your country or  
800 institution) were obtained?

801 Answer: [NA]

802 Justification: There were no study participants for this work.

803 Guidelines:

804  
805  
806  
807  
808  
809  
810  
811  
812  
813

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.