# Benchmarking Scalable Epistemic Uncertainty Quantification in Organ Segmentation

Jadie Adams[1,2] and Shireen Y. Elhabian[1,2]

[1] Scientific Computing and Imaging Institute, University of Utah, UT, USA
[2] School of Computing, University of Utah, UT, USA
`jadie.adams@utah.edu, shireen@sci.utah.edu`

**Abstract.** Deep learning based methods for automatic organ segmentation have shown promise in aiding diagnosis and treatment planning. However, quantifying and understanding the uncertainty associated with model predictions is crucial in critical clinical applications. While many techniques have been proposed for epistemic or model-based uncertainty estimation, it is unclear which method is preferred in the medical image analysis setting. This paper presents a comprehensive benchmarking study that evaluates epistemic uncertainty quantification methods in organ segmentation in terms of accuracy, uncertainty calibration, and scalability. We provide a comprehensive discussion of the strengths, weaknesses, and out-of-distribution detection capabilities of each method as well as recommendations for future improvements. These findings contribute to the development of reliable and robust models that yield accurate segmentations while effectively quantifying epistemic uncertainty.

## 1 Introduction

Deep learning systems have made significant strides in automating organ segmentation from 3D medical images. Segmentation networks can be efficiently integrated into image processing pipelines, facilitating research and clinical use (i.e., tumor segmentation in radiotherapy [30] and hippocampus segmentation for neurological disease analysis [6]). However, these systems also introduce new challenges and risks compared to traditional segmentation processes, including issues of bias, errors, and lack of transparency. Deep networks are prone to providing overconfident estimates and thus cannot be blindly trusted in sensitive decision-making scenarios without the safeguard of granular uncertainty quantification (UQ) [20,21]. UQ is the process of estimating and representing the uncertainty associated with predictions made by deep neural networks. UQ provides necessary insight into the reliability and confidence of the model's predicted segmentation. In the context of organ segmentation, areas near organ boundaries can be uncertain due to the low contrast between the target organ and surrounding tissues [33]. Pixel or voxel-level uncertainty estimates can be used to identify potential incorrect regions or guide user interactions for refinement [33,27]. This enables quality control of the segmentation process and the detection of out-of-distribution (OOD) samples.

Two forms of uncertainty are distinguished in deep learning frameworks: aleatoric and epistemic[14]. *Aleatoric uncertainty* refers to the inherent uncertainty in the input data distribution that cannot be reduced [16] (i.e., uncertainty resulting from factors like image acquisition noise, over-exposure, occlusion, or a lack of visual features [32]). Aleatoric uncertainty is typically quantified by adjusting the model to be stochastic (predicting a distribution rather than a point-wise estimate [16]) or by methods such as test time augmentation [33,31]. *Epistemic uncertainty* is model-based and arises from a lack of knowledge or uncertainty about the model's parameters due to limited training data or model complexity. Capturing epistemic uncertainty is considerably more difficult as it cannot be learned as a function of the input but rather requires fitting a distribution over model parameters. Several approaches have been proposed to accomplish this, but many of them significantly increase the computational cost and memory footprint and may impact prediction accuracy [7]. There is no ubiquitous method for epistemic UQ in segmentation networks, as each proposed technique has its own trade-offs and limitations.

This study benchmarks Bayesian and frequentist epistemic UQ methods for organ segmentation from 3D CT scans in terms of scalability, segmentation accuracy, and uncertainty calibration using multiple datasets. While previous benchmarks (e.g., [25,29,24]) have been conducted on small subsets of such methods, there is a need for a comprehensive evaluation. To the best of our knowledge, this work provides the most extensive benchmarking of scalable methods for epistemic UQ in medical segmentation. The key contributions are as follows:

1. We conduct a benchmark evaluation of scalable methods for epistemic UQ in medical image segmentation, including deep ensemble [28], batch ensemble [34], Monte Carlo dropout [10], concrete dropout [11], Rank-1 Bayesian Neural Net (BNN) [7], latent posterior BNN [9], Stochastic Weight Averaging (SWA) [15], SWA Gaussian (SWAG) [22], and Multi-SWAG [35].
2. We evaluate these methods in detecting out-of-distribution (OOD) instances, which is an important aspect of robust uncertainty estimation.
3. We provide a comprehensive discussion of the strengths and weaknesses of the evaluated methods, enabling a better understanding of their performance characteristics and potential improvements.
4. To facilitate further research and reproducibility, we provide an open-source PyTorch implementation of all benchmarked methods.[*]

## 2    Epistemic Uncertainty Quantification Techniques

Modeling epistemic uncertainty in a scalable manner poses significant challenges as it entails placing distributions over model weights. Both Bayesian and frequentist methods have been proposed, with Bayesian approaches aiming to directly estimate the posterior distribution over the model's parameters, while frequentist methods use ensembles of models to approximate the posterior empirically.

---

[*] Source code is publicly available: `https://github.com/jadie1/MedSegUQ`

In Bayesian deep learning, obtaining an analytical solution for the posterior is often intractable, necessitating the use of approximate posterior inference techniques such as variational inference [3]. The most common Bayesian technique for UQ is **Monte Carlo (MC) dropout** sampling, as it provides a fast, scalable solution for approximate variational inference [10]. In MC dropout, uncertainty is captured by the spread of predictions resulting from sampled dropout masks in inference. However, obtaining well-calibrated epistemic UQ with dropout requires a time-consuming grid search to tune layer-wise dropout probabilities. **Concrete dropout** [11] was proposed to address this limitation by automatically optimizing layer-wise dropout probabilities along with the network weights. Certain Bayesian approaches for approximate inference are excluded from this benchmark due to their limited scalability and tendency to underfit [7,10,19], such as sampling-based Markov chain Monte Carlo [5], Bayes by Backprop [4], and variational inference based methods [3]. These techniques rely on structured or factorized distributions with tied parameters, have a high computational cost, and are generally slow to converge [10]. Additionally, previous work has shown that such methods perform similarly to the much more lightweight MC dropout approach in medical image segmentation UQ [26].

**Deep ensembles** are an effective and popular frequentist method for UQ [13]. Ensembling involves training multiple independent networks (or ensemble members) with different initialization then aggregating predictions for improved robustness [8]. The spread or variability among the predictions of the ensemble members effectively captures the epistemic uncertainty [19]. It has been shown that deep ensemble models can provide a better approximation than standard Bayesian methods [35]. The main drawback of deep ensembles lies in their computational and memory costs, which increase linearly with the number of ensemble members. To address the trade-off between accuracy and scalability, **batch ensemble** [34] has been proposed. Batch ensemble [34] compromises between a single network and an ensemble by employing shared weight matrices and lightweight rank-1 ensemble members. The concept of batch ensembling has also been applied to improve the scalability of Bayesian Neural Networks (BNNs). **Rank-1 BNN** [7] reduces computational complexity by utilizing a rank-1 parameterization in variational inference. **Latent Posterior BNN (LP-BNN)** [9] was proposed to improve scalability further by learning the posterior distribution of lower-dimensional latent variables derived from rank-1 vectors. This is accomplished by training layer-wise variational autoencoders (VAEs) [18] on the rank-1 vectors. The latent space of these VAEs can then be sampled, providing a distribution of rank-1 weights.

Additional methods of epistemic UQ have been developed based on the stochastic weight averaging (SWA) technique [15]. SWA was proposed to enhance generalization in deep learning by estimating the mean of the stationary distribution of SGD iterates. In SWA, final model weights are defined by averaging the weights traversed during SGD after initial convergence. **SWA-Gaussian (SWAG)** [22] fits a Gaussian distribution to the traversed weights, creating an approximate posterior distribution over the weights that can be sam-

pled for Bayesian model averaging. It has been shown that combining traditional Bayesian methods with ensembling improves the fidelity of approximate inference via multimodal marginalization, resulting in a more robust, accurate model [35]. Based on this observation, **Multi-SWAG** [35] was proposed as an ensemble of SWAG models. These approaches offer alternative ways to capture epistemic uncertainty by leveraging the ensemble characteristics and combining them with Bayesian principles. They aim to improve the fidelity of approximate inference and provide scalable solutions for epistemic UQ in deep learning tasks.

## 3   Experimental Design

We utilize the residual U-Net architecture originally proposed for cardiac left-ventricle segmentation [17] as a base architecture to compare the epistemic UQ techniques. This model is comprised of residual units of 3D convolutional layers with batch normalization and PReLU activation. As a baseline for UQ calibration analysis, we consider the predicted segmentation probabilities. Specifically, for the **base** model, we quantify voxel-wise UQ as: $UQ = 1 - C$, where confidence, $C$, is the maximum of the foreground and background softmax probabilities. This is not strictly a measure of epistemic UQ, as there is no notion of posterior approximation and marginalization. However, this formulation of UQ has been shown to correlate with prediction error and is useful for OOD detection [12] and thus provides an evaluation baseline.

We benchmark the following aforementioned scalable methods for epistemic UQ: Ensemble[19], Batch Ensemble[34], MC Dropout[10], Concrete Dropout[11], Rank1 BNN[7], LP-BNN[9], SWAG[22], and Mulit-SWAG[35]. In implementing dropout and rank1-based methods, dropout and batch ensemble are applied to every convolutional layer respectively. Models are trained on $(96 \times 96 \times 96)$ patches of the input images scaled to an intensity of $[0, 1]$. A validation set is used to assess the convergence of each model. Models are trained until the performance on the validation set has not improved in 100 epochs. The model weights resulting from the epoch with the best validation performance are used in the evaluation of held-out testing data. Implementation and tuned hyperparameter values for each model are provided in the GitHub repository.

### 3.1   Datasets

We utilize two open-source datasets from the Medical Segmentation Decathlon [32] in evaluation: the spleen and pancreas. These datasets comprise of 3D CT images and corresponding manual ground truth segmentations. The **spleen** dataset was selected to provide a typical medical image analysis scenario where data is scarce and varied. There are only 41 instances in the spleen dataset, and the size of the spleen versus the background varies widely. The **pancreas** dataset was selected to evaluate OOD detection accuracy. This dataset contains cancerous cases with segmentations of both the pancreas organ and tumor masses. We analyze the accuracy of joint segmentation of the pancreas and tumors, holding out the cases with the largest tumors an OOD test set.

### 3.2 Metrics

In 3D image segmentation, accuracy is typically assessed via the **Dice Similarity Coefficient (DSC)** between manual annotations and the model's prediction. The DSC metric captures the percentage of voxel-wise agreement between a predicted segmentation and its corresponding ground truth. In these experiments the target organs are small compared to the total image size. Because of this, we only include the foreground in DSC calculations so that the DSC value is not overwhelmed by the background signal.

To measure overall uncertainty calibration, we consider the correlation between estimated epistemic uncertainty and prediction error (100 - DSC). We report the **Pearson correlation coefficient (r)** between error and sum of the uncertainty map, where a higher r-value indicates better correlation.

Finally, we assess segmentation accuracy and uncertainty quality jointly via **error-retention curves** [19,23]. Error-retention curves depict a given error metric over a dataset as ground-truth labels replace a model's predictions in order of decreasing estimated uncertainty. The **area under the curve (AUC)** is reduced as the overall error is decreased, as well as the correlation between error and uncertainty is increased. We report the **area under the error retention curve (R-AUC)** using 100 - DSC as the error metric. Smaller R-AUC indicates both better segmentation accuracy and uncertainty calibration.

## 4 Results

### 4.1 Spleen

Because the spleen dataset is comprised of only 41 examples, we employ K-folds to define the training, validation, and test sets. The data is split into 70% train, 10% validation, and 20% held-out test using five different folds. For each fold, a separate model of every form is trained. In this manner, results are reported across the entire dataset, where the predicted segmentation of each image is acquired from the model for which the image was held out. The results are reported in Table 1, and qualitative visualizations are provided in Fig. 1. For Ensemble and Batch Ensemble models four ensemble members are used, thus four predictions are averaged and used for UQ estimation.

For methods that fit a distribution over weights, any number of weight samples can be used to provide an average prediction and UQ estimation. Additional samples improve accuracy and calibration but also increase inference time. In evaluating such methods, we elect to use 4 and 30 samples for a comprehensive comparison.

Table 1: Spleen results: Mean and standard deviation values across held-out data, best values in bold.

| | DSC ↑ | r ↑ | R-AUC ↓ |
|---|---|---|---|
| Base | 89.15±9.55 | 0.56 | 2.54±4.52 |

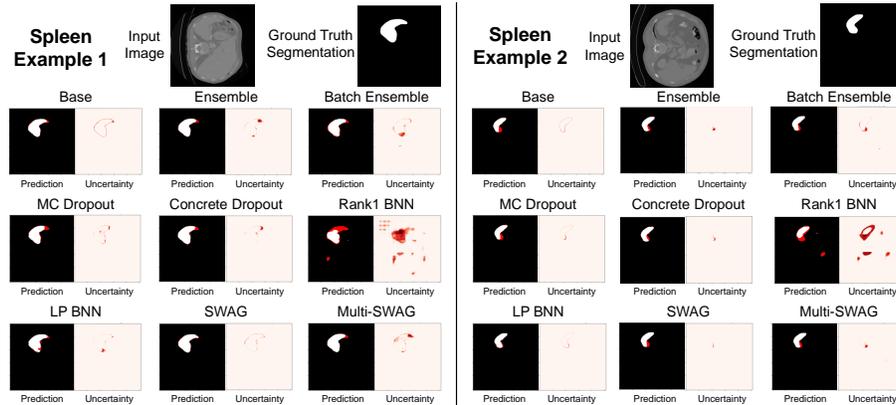| Method | 4 Samples | | | 30 Samples | | |
|---|---|---|---|---|---|---|
| | DSC ↑ | r ↑ | R-AUC ↓ | DSC ↑ | r ↑ | R-AUC ↓ |
| Ensemble[19] | **92.77**±5.34 | 0.37 | 0.71±1.09 | N/A | N/A | N/A |
| Batch Ensemble[34] | 86.87±14.21 | 0.40 | 0.58±1.09 | N/A | N/A | N/A |
| MC Dropout[10] | 86.16±19.07 | -0.07 | 1.54±4.06 | 86.40±18.33 | -0.08 | 1.31±3.79 |
| Concrete Dropout[11] | 90.21±6.42 | 0.45 | **0.42±0.47** | 90.23±6.49 | 0.59 | **0.28±0.35** |
| Rank1 BNN[7] | 55.01±19.17 | 0.10 | 3.69±3.47 | 66.36±16.84 | 0.11 | 0.79±1.45 |
| LP-BNN[9] | 87.77±10.66 | 0.57 | 0.98±1.27 | 87.74±10.83 | 0.47 | 0.93±1.48 |
| SWAG[22] | 87.80±15.53 | 0.18 | 0.99±2.47 | 87.80±15.69 | 0.20 | 0.74±1.71 |
| Mulit-SWAG[35] | 92.49±10.18 | **0.69** | 0.69±0.84 | **93.11±10.07** | **0.64** | 0.57±0.77 |

Fig. 1: Slices of two spleen examples are provided with the segmentation resulting from each model with error overlaid in red. Additionally, predicted uncertainty maps are shown where darker red indicates higher uncertainty.

## 4.2   Pancreas

The pancreas dataset is used to analyze the robustness and uncertainty calibration of the various methods in the case of OOD examples. To this end, we calculate the ratio of tumor to pancreas voxels in the ground truth segmentations. We hold out the 50 instances with the largest tumor ratio as a test set. This provides an OOD test set, as the models are trained only on examples with smaller tumors. The remaining 231 image/segmentation pairs are randomly split into a single training, validation, and in-distribution (ID) test set using a 70%, 10%, 20% split. Fig. 2 displays the distributions of tumor-to-pancreas ratios in the ID and OOD test sets. The results are reported on both the ID and OOD test set in Table 2. Additionally, the correlation between the tumor-to-pancreas ratio and the estimated uncertainty across both test sets is reported in Table 3. We expect well-calibrated UQ to correlate with the tumor-to-pancreas ratio as the models are not exposed examples with large tumors in training. However, none of the epistemic UQ quantification methods proved a strong correlation with the ratio, suggesting these models are not effective in accurate OOD detection.

Table 2: Pancreas results: Mean and standard deviation values across both held-out test sets are reported with the best values marked in bold.

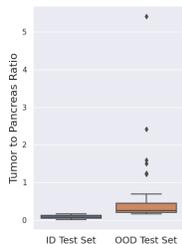| | In-Distribution Test Set | | | | | | | Out-of-Distribution Test Set | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | DSC | R | R-AUC | | | | | DSC | R | R-AUC | | |
| | | Base | 71.02±14.49 | 0.12 | 8.66±7.36 | | | | Base | 67.32±18.35 | 0.04 | 8.95±7.03 | | |
| Model | 4 Samples | | | 30 Samples | | | 4 Samples | | | 30 Samples | | | | |
| | DSC | r | R-AUC | DSC | r | R-AUC | DSC | r | R-AUC | DSC | r | R-AUC | | |
| Ensemble[19] | 72.44±13.54 | 0.04 | 4.70±3.57 | N/A | N/A | N/A | **67.54±19.58** | 0.15 | 5.65±5.64 | N/A | N/A | N/A | | |
| Batch Ensemble[34] | 64.04±16.49 | 0.24 | 6.30±4.38 | N/A | N/A | N/A | 57.08±21.51 | 0.12 | 9.77±10.02 | N/A | N/A | N/A | | |
| MC Dropout[10] | 70.16±13.96 | 0.36 | 9.97±7.43 | 70.22±13.93 | 0.38 | 8.90±6.54 | 67.24±18.13 | 0.30 | 10.63±9.33 | 67.35±18.20 | 0.33 | 9.45±9.18 | | |
| Concrete Dropout[11] | **73.07±11.91** | 0.32 | **4.67±3.13** | **73.06±11.98** | 0.33 | **4.29±2.64** | 66.99±18.40 | 0.08 | 8.39±9.17 | **67.83±18.62** | 0.27 | 7.54±8.71 | | |
| Rank1 BNN[7] | 12.77±11.81 | -0.23 | 11.37±8.99 | 17.76±13.60 | -0.2 | 9.07±8.91 | 9.20±10.17 | -0.33 | 12.91±10.17 | 12.39±11.85 | -0.41 | 9.21±4.72 | | |
| LP-BNN[9] | 65.28±16.18 | 0.36 | 4.92±3.05 | 65.31±16.10 | 0.21 | 4.66±2.94 | 60.39±19.34 | 0.14 | 6.34±5.03 | 60.16±19.78 | 0.19 | 7.26±6.95 | | |
| SWAG[22] | 66.68±17.27 | 0.36 | 9.32±8.24 | 66.69±17.29 | **0.49** | 9.34±8.14 | 63.20±20.90 | **0.31** | 10.06±8.16 | 63.07±21.04 | **0.45** | 9.89±7.86 | | |
| Mulit-SWAG[35] | 69.67±15.06 | **0.39** | 5.57±4.32 | 69.31±15.19 | 0.41 | 5.46±4.24 | 64.94±21.13 | 0.12 | **6.19±5.86** | 64.65±21.28 | 0.14 | **6.01±5.71** | | |

Fig. 2: Pancreas test set ratio box plots.

Table 3: Pearson correlation coefficients across both pancreas test sets are reported, where error = 100 - DSC.

| Model | r values | | |
|---|---|---|---|
| | Ratio/Error | UQ/Error ↑ | UQ/Ratio ↑ |
| Base | 0.42 | -0.08 | -0.06 |
| Ensemble[19] | 0.44 | 0.22 | -0.01 |
| Batch Ensemble[34] | 0.34 | 0.13 | 0.10 |
| MC Dropout[10] | 0.46 | 0.35 | 0.10 |
| Concrete Dropout[11] | 0.48 | 0.28 | 0.08 |
| Rank1 BNN[7] | 0.02 | -0.32 | 0.06 |
| LP-BNN[9] | 0.38 | 0.2 | 0.06 |
| SWAG[22] | 0.34 | **0.46** | **0.11** |
| Mulit-SWAG[35] | 0.39 | 0.22 | -0.04 |

## 4.3 Scalability Comparison

Table 4 reports the average time and memory requirements associated with training and testing each model on the pancreas dataset. Note that for models that can be sampled,the reported inference time is for a single sample. Inference time scales linearly with the number of samples. Additionally, note "params size" refers to the memory required to store the parameters and "pass size" refers to the memory required to perform a forward/backward pass through the network.

Table 4: Scalability comparison: Time reported in seconds and memory size in MB, best values in bold.

| Model | Train epochs | Train time | Inference time | Total params | Params size | Pass size |
|---|---|---|---|---|---|---|
| Base | 222 | 17261 | 0.2587 | **4808917** | **19.24** | **1211.20** |
| Ensemble[19] | 888 | 69045 | 1.0348 | 19235668 | 76.96 | 4844.80 |
| Batch Ensemble[34] | 341 | 45409. | 0.8698 | 4824513 | **19.24** | 4842.81 |
| MC Dropout[10] | **197** | **16841** | **0.2548** | **4808917** | **19.24** | **1211.20** |
| Concrete Dropout[11] | 259 | 36847 | 0.3694 | 4808934 | **19.24** | **1211.20** |
| Rank1 BNN[7] | 1142 | 157100 | 0.7712 | 4835229 | **19.24** | 4844.81 |
| LP-BNN[9] | 881 | 121343 | 0.8742 | 4957940 | 19.77 | 4844.92 |
| SWAG[22] | 422 | 31615 | 0.2921 | 9617834 | 38.48 | **1211.20** |
| Mulit-SWAG[35] | 1688 | 126460 | 1.1684 | 38471336 | 153.92 | 4844.80 |

## 5 Discussion and Conclusion

We conducted a benchmark of scalable epistemic UQ techniques on two challenging organ segmentation tasks. The spleen experiment represented a low training budget scenario, while the pancreas experiment involved significant variation in the shape and size of the organ and tumor masses. These challenging scenarios, where the base model provides low prediction accuracy, served as stress tests for UQ evaluation. We discuss the performance of each model as follows:

**Base model**: While the base models provided competitive accuracy and r values (suggesting instance-level UQ/error correlation) the R-AUC values were low. This indicates the voxel-level UQ did not correlate well with error and thus could not be used to accurately identify erroneous regions. The base model UQ correlates more with the organ boundary than error, as can be seen in Fig. 1.

**Ensemble[19]**: As expected, the ensemble model provided an accuracy improvement over the base model. However, it came at the expense of scalability, which is an important consideration in practical applications. More scalable methods outperformed ensembling in terms of R-AUC in these experiments.

**Batch Ensemble[34]**: While batch ensemble reduces the memory cost associated with ensembling, it did not provide the same accuracy improvement. This is likely because joint training of the rank1 members proved difficult on these challenging tasks. However, it still provided improved UQ over the base model.

**MC Dropout[10]**: This approach is appealing as it does not increase memory costs or impact the training objective. However, it did not perform as well as concrete dropout, highlighting the importance of tuning layer-wise dropout probabilities. For MC dropout models, we used a dropout rate of 0.1 for all layers. The concrete dropout optimization found a dropout probability of around 0.08 for shallow layers, increasing to around 0.16 for the deepest layer.

**Concrete Dropout[11]**: This technique arguably performed best overall and is desirable as it is scalable and only requires the addition of concrete dropout layers and a loss regularization term.

**Rank1 BNN[7]**: This model did not perform well on either task, especially the pancreas segmentation. While the rank1 parameterization greatly improves the scalability of the BNNs, it does not appear to solve the issue of poor convergence that BNNs are prone to suffer from.

**LP-BNN[9]**: Approximating the posterior in a learned latent space of the rank1 vectors improved convergence, as LP-BNN outperformed Rank1 BNN. However, LP-BNN did not perform as well as other methods with regard to any metrics, likely because training layer-wise VAEs complicates the learning task.

**SWAG[22]**: The SWAG models did not outperform the base models in terms of accuracy as expected. This can be attributed to the fact that the base models used in the evaluation were those resulting from the epoch with the best validation performance, whereas the SWAG weight posterior was fit across the converged SGD trajectory. This technique is desirable because it does not require adapting the architecture in any way and can be considered a post hoc process.

**Mulit-SWAG[35]** Ensembling SWAG models improved the accuracy and UQ calibration but again at the expense of scalability.

This benchmark provides some insights into how UQ methods can be improved. The Multi-SWAG performance reinforces the notion that ensembling Bayesian methods can improve the fidelity of approximate inference by enabling multimodal marginalization. This could be made more scalable by combining SWAG with a batch ensemble model rather than applying naive ensembling. Existing work has also demonstrated that combining ensembling with dropout improves performance on related medical imaging tasks [1]. These experiments additionally demonstrate that LP-BNN is a desirable alternative to Rank1 BNN. However, improvements could be made to the LP-BNN process of learning a low-dimension representation of the rank1 vectors, as layer-wise VAEs increase the training burden and hyperparameters to tune. The pancreas OOD analysis reveals that none of the epistemic UQ methods were effective at detecting instances

with larger tumor sizes than those seen in training. As Table 2 demonstrate, all methods provided better calibrated uncertainty estimates on in domain test data than OOD. Such failure has been noted before, as model misestimation can result in overconfidence in OOD predictions [36,2]. This illustrates the need to consider alternative test statistics and objectives in developing epistemic uncertainty estimation techniques.

In conclusion, our benchmarking study of scalable epistemic uncertainty quantification techniques for challenging organ segmentation tasks highlights the importance of accurate uncertainty estimation in medical image analysis. The insights gained from this study can guide researchers and practitioners in selecting appropriate methods to enhance the reliability and robustness of deep learning models for organ segmentation, ultimately contributing to improved diagnosis and treatment planning in clinical practice.

# References

1. Adams, J., Elhabian, S.: Fully bayesian vib-deepssm. arXiv preprint arXiv:2305.05797 (2023)
2. Besnier, V., Bursuc, A., Picard, D., Briot, A.: Triggering failures: Out-of-distribution detection by learning from local adversarial attacks in semantic segmentation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 15701–15710 (2021)
3. Blei, D.M., Kucukelbir, A., McAuliffe, J.D.: Variational inference: A review for statisticians. Journal of the American Statistical Association **112**(518), 859–877 (2017). https://doi.org/10.1080/01621459.2017.1285773, `https://doi.org/10.1080/01621459.2017.1285773`
4. Blundell, C., Cornebise, J., Kavukcuoglu, K., Wierstra, D.: Weight uncertainty in neural network. In: International conference on machine learning. pp. 1613–1622. PMLR (2015)
5. Chen, T., Fox, E., Guestrin, C.: Stochastic gradient hamiltonian monte carlo. In: International conference on machine learning. pp. 1683–1691. PMLR (2014)
6. Chupin, M., Gérardin, E., Cuingnet, R., Boutet, C., Lemieux, L., Lehéricy, S., Benali, H., Garnero, L., Colliot, O.: Fully automatic hippocampus segmentation and classification in alzheimer's disease and mild cognitive impairment applied on data from adni. Hippocampus **19**(6), 579–587 (2009)
7. Dusenberry, M., Jerfel, G., Wen, Y., Ma, Y., Snoek, J., Heller, K., Lakshminarayanan, B., Tran, D.: Efficient and scalable bayesian neural nets with rank-1 factors. In: International conference on machine learning. pp. 2782–2792. PMLR (2020)

8.  Fort, S., Hu, H., Lakshminarayanan, B.: Deep ensembles: A loss landscape perspective. arxiv 2019. arXiv preprint arXiv:1912.02757 (2019)
9.  Franchi, G., Bursuc, A., Aldea, E., Dubuisson, S., Bloch, I.: Encoding the latent posterior of bayesian neural networks for uncertainty quantification. arXiv preprint arXiv:2012.02818 (2020)
10. Gal, Y., Ghahramani, Z.: Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In: international conference on machine learning. pp. 1050–1059. PMLR (2016)
11. Gal, Y., Hron, J., Kendall, A.: Concrete dropout. Advances in neural information processing systems **30** (2017)
12. Hendrycks, D., Gimpel, K.: A baseline for detecting misclassified and out-of-distribution examples in neural networks. arXiv preprint arXiv:1610.02136 (2016)
13. Hu, R., Huang, Q., Chang, S., Wang, H., He, J.: The mbpep: a deep ensemble pruning algorithm providing high quality uncertainty prediction. Applied Intelligence **49**, 2942–2955 (2019)
14. Hüllermeier, E., Waegeman, W.: Aleatoric and epistemic uncertainty in machine learning: An introduction to concepts and methods. Machine Learning **110**(3), 457–506 (2021)
15. Izmailov, P., Podoprikhin, D., Garipov, T., Vetrov, D., Wilson, A.G.: Averaging weights leads to wider optima and better generalization. arXiv preprint arXiv:1803.05407 (2018)
16. Kendall, A., Gal, Y.: What uncertainties do we need in bayesian deep learning for computer vision? Advances in neural information processing systems **30** (2017)
17. Kerfoot, E., Clough, J., Oksuz, I., Lee, J., King, A.P., Schnabel, J.A.: Left-ventricle quantification using residual u-net. In: Statistical Atlases and Computational Models of the Heart, STACOM 2018, Held in Conjunction with MICCAI 2018. pp. 371–380. Springer (2019)
18. Kingma, D.P., Welling, M.: Auto-encoding variational bayes. Stat **1050**, 1 (2014)
19. Lakshminarayanan, B., Pritzel, A., Blundell, C.: Simple and scalable predictive uncertainty estimation using deep ensembles. Advances in neural information processing systems **30** (2017)
20. Li, B., Qi, P., Liu, B., Di, S., Liu, J., Pei, J., Yi, J., Zhou, B.: Trustworthy ai: From principles to practices. ACM Computing Surveys **55**(9), 1–46 (2023)
21. Liang, W., Tadesse, G.A., Ho, D., Fei-Fei, L., Zaharia, M., Zhang, C., Zou, J.: Advances, challenges and opportunities in creating data for trustworthy ai. Nature Machine Intelligence **4**(8), 669–677 (2022)
22. Maddox, W.J., Izmailov, P., Garipov, T., Vetrov, D.P., Wilson, A.G.: A simple baseline for bayesian uncertainty in deep learning. Advances in neural information processing systems **32** (2019)
23. Malinin, A., Athanasopoulos, A., Barakovic, M., Cuadra, M.B., Gales, M.J., Granziera, C., Graziani, M., Kartashev, N., Kyriakopoulos, K., Lu, P.J., et al.: Shifts 2.0: Extending the dataset of real distributional shifts. arXiv preprint arXiv:2206.15407 (2022)
24. Mehrtash, A., Wells, W.M., Tempany, C.M., Abolmaesumi, P., Kapur, T.: Confidence calibration and predictive uncertainty estimation for deep medical image segmentation. IEEE transactions on medical imaging **39**(12), 3868–3878 (2020)
25. Ng, M., Guo, F., Biswas, L., Petersen, S.E., Piechnik, S.K., Neubauer, S., Wright, G.: Estimating uncertainty in neural networks for cardiac mri segmentation: a benchmark study. IEEE Transactions on Biomedical Engineering (2022)

26. Ng, M., Guo, F., Biswas, L., Wright, G.A.: Estimating uncertainty in neural networks for segmentation quality control. In: 32nd International Conference on Neural Information Processing Systems (NIPS 2018), Montréal, Canada, no. NIPS. pp. 3–6 (2018)

27. Prassni, J.S., Ropinski, T., Hinrichs, K.: Uncertainty-aware guided volume segmentation. IEEE transactions on visualization and computer graphics **16**(6), 1358–1365 (2010)

28. Rahaman, R., et al.: Uncertainty quantification and deep ensembles. Advances in Neural Information Processing Systems **34**, 20063–20075 (2021)

29. Sahlsten, J., Jaskari, J., Wahid, K.A., Ahmed, S., Glerean, E., He, R., Kann, B., Makitie, A.A., Fuller, C.D., Naser, M.A., et al.: Application of simultaneous uncertainty quantification for image segmentation with probabilistic deep learning: Performance benchmarking of oropharyngeal cancer target delineation as a use-case. medRxiv pp. 2023–02 (2023)

30. Savjani, R.R., Lauria, M., Bose, S., Deng, J., Yuan, Y., Andrearczyk, V.: Automated tumor segmentation in radiotherapy. In: Seminars in Radiation Oncology. vol. 32, pp. 319–329. Elsevier (2022)

31. Shanmugam, D., Blalock, D., Balakrishnan, G., Guttag, J.: Better aggregation in test-time augmentation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 1214–1223 (2021)

32. Simpson, A.L., Antonelli, M., Bakas, S., Bilello, M., Farahani, K., Van Ginneken, B., Kopp-Schneider, A., Landman, B.A., Litjens, G., Menze, B., et al.: A large annotated medical image dataset for the development and evaluation of segmentation algorithms. arXiv preprint arXiv:1902.09063 (2019)

33. Wang, G., Li, W., Aertsen, M., Deprest, J., Ourselin, S., Vercauteren, T.: Aleatoric uncertainty estimation with test-time augmentation for medical image segmentation with convolutional neural networks. Neurocomputing **338**, 34–45 (2019)

34. Wen, Y., Tran, D., Ba, J.: Batchensemble: an alternative approach to efficient ensemble and lifelong learning. arXiv preprint arXiv:2002.06715 (2020)

35. Wilson, A.G., Izmailov, P.: Bayesian deep learning and a probabilistic perspective of generalization. Advances in neural information processing systems **33**, 4697–4708 (2020)

36. Zhang, L., Goldstein, M., Ranganath, R.: Understanding failures in out-of-distribution detection with deep generative models. In: International Conference on Machine Learning. pp. 12427–12436. PMLR (2021)