

# Laplacian Segmentation Networks: Improved Epistemic Uncertainty from Spatial Aleatoric Uncertainty

Kilian Zepf<sup>1,\*</sup>, Selma Wanna<sup>2,\*</sup>, Marco Miani<sup>1</sup>, Juston Moore<sup>3</sup>,  
Jes Frellsen<sup>1</sup>, Søren Hauberg<sup>1</sup>, Aasa Feragen<sup>1</sup>, Frederik Warburg<sup>1</sup>

<sup>1</sup>Technical University of Denmark, <sup>2</sup>The University of Texas at Austin, <sup>3</sup>Los Alamos National Laboratory

## Abstract

Out of distribution (OOD) medical images are frequently encountered, e.g. because of site- or scanner differences, or image corruption. OOD images come with a risk of incorrect image segmentation, potentially negatively affecting downstream diagnoses or treatment. To ensure robustness to such incorrect segmentations, we propose Laplacian Segmentation Networks (LSN) that jointly model epistemic (model) and aleatoric (data) uncertainty in image segmentation. We capture data uncertainty with a spatially correlated logit distribution. For model uncertainty, we propose the first Laplace approximation of the weight posterior that scales to large neural networks with skip connections that have high-dimensional outputs. Empirically, we demonstrate that modelling spatial pixel correlation allows the Laplacian Segmentation Network to successfully assign high epistemic uncertainty to out-of-distribution objects appearing within images.

## 1. Introduction

Image segmentation is a core component in the biomedical image analysis toolbox, used extensively both to quantify organs for use in scientific studies, as well as to inform clinicians by outlining relevant parts of the image. Its widespread use places high requirements for its safe and interpretable operation. However, neural networks, which form the backbones of most modern segmentation models, are infamous for being overconfident in their predictions outside the training distribution (Hendrycks and Gimpel, 2016). As a result, downstream predictions can be highly unreliable even though they might come with high accuracy on in-distribution data.

Figure 1 shows an example derived from a U-net model. While the network has never seen images corrupted by syn-

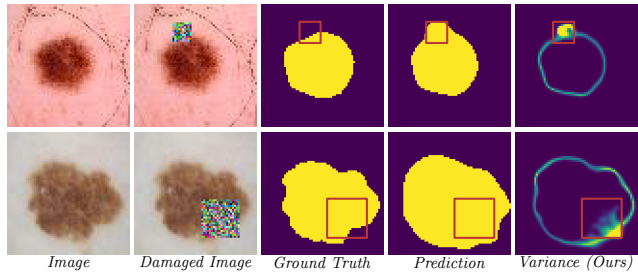


Figure 1: Corrupted images can lead to confident and wrong predictions by a U-net (fourth column). In automated downstream tasks, where predictions are not inspected by experts, this can remain undiscovered, leading to critical mistakes and wrong follow up treatment of patients. Our model assigns high epistemic uncertainty to the corrupted images.

thetic noise during training it will confidently predict pixels replaced by noise as skin lesions. This renders every following analysis based on the prediction unreliable. While the example is synthetic, similar effects occur in images corrupted, for example, by polluted cameras, or data loss.

**In this work**, we introduce Laplace approximations (LA) for epistemic uncertainty quantification in binary image segmentation models. Current Laplace approximations (Daxberger et al., 2021) scale quadratically with the output dimension of the neural network, which prevent their usage in segmentation. We develop a fast Hessian approximation for deep architectures with skip connections, which are an integral part of modern segmentation networks such as the U-net (Ronneberger et al., 2015). This enables us to combine the aleatoric logit distribution of Stochastic Segmentation Networks (SSN) (Monteiro et al., 2020) with post-hoc Laplace approximations for epistemic uncertainty quantification. Specifically, we propose to find the Gaussian approximation on the mean network of the SSN. This can be viewed as combining an evidential deep learning (EDL) model (Sensoy et al., 2018; Ulmer et al., 2023) with a LA approximation of epistemic uncertainty rather than the point estimate prevalently used in EDL, and, to the best of our

\* denotes equal contribution. Correspondence to: kmze@dtu.dk, slwanna@utexas.edu, mmia@dtu.dk, frwa@dtu.dk

knowledge, this is the first work to propose this combination. We demonstrate that the spatial correlation in the aleatoric component is essential for the quantification of epistemic uncertainty in the Laplace approximation.

Joint modelling of epistemic and aleatoric segmentation uncertainty is not new (Kendall and Gal, 2017; Popescu et al., 2021; Joy et al., 2021; Fuchs et al., 2021). The binary cross entropy loss typically used to train segmentation models makes an assumption of independent pixels. Segmentation models using this loss tend to get stuck in local minima when the target is small compared to the background, unless the classes are reweighted in the loss (Lin et al., 2017). We show experimentally that this changes when you spatially correlate pixels: The model can converge without reweighting the smaller class. As a result, the loss landscape also looks different, leading to a change in the LA estimation of epistemic uncertainty. In our experiments, we observe an effect of this in the form of improved OOD behavior.

In a series of experiments on two medical binary segmentation tasks, we show that the proposed method outperforms baselines on a variety of OOD scenarios, including the detection of distribution shifts on datasets and highlighting corrupted parts of an image, by assigning high uncertainty to respective pixels. Providing a more reliable model uncertainty, the proposed Laplacian Segmentation Networks contributes to the applicability and safety of binary segmentation models in practice.

## 2. Background and Related Work

While several different sources and taxonomies of uncertainties have been proposed (Kiureghian and Ditlevsen, 2009; Gawlikowski et al., 2021), the Bayesian framework (Bishop, 2006; Kendall and Gal, 2017) distinguishes between two of them: aleatoric and epistemic. These appear directly in the predictive distribution

$$p(y|x, D) = \int \underbrace{p(y|x, \theta)}_{\text{aleatoric}} \underbrace{p(\theta|D)}_{\text{epistemic}} d\theta, \quad (1)$$

where  $(x, y)$  is an input-output pair,  $D$  is the training data and  $\theta$  the model parameters. The aleatoric (likelihood) density  $p(y|x, \theta)$  models noise or variation in the data from which the model learns. This is typically caused either by ambiguities in the image or in the definition of the object to be annotated, for example tumor boundaries which are hard to define due to gradual tissue infiltration. The posterior density  $p(\theta|D) \propto p(D|\theta)p(\theta)$  reflects the epistemic uncertainty of the estimated model, which quantifies the degree to which the model itself should be trusted. Here  $p(\theta)$  is the prior over the parameters and  $p(D|\theta)$  is the likelihood of the training data.

The intractable posterior distribution over the parameters prohibits the evaluation of the integral in Eq. (1). To over-

come this, it is necessary to utilize suitable approximations for the predictive distribution itself or parts of the integral on the right hand side of Eq. (1).

In standard image segmentation, the posterior distribution  $p(\theta|D)$  is usually approximated with a Dirac distribution  $\delta(\theta - \hat{\theta})$ , where  $\hat{\theta}$  is the maximum a posteriori (MAP) or maximum likelihood (ML) estimate, which simplifies the posterior predictive to  $p(y|x, D) = p(y|x, \hat{\theta})$ . The additional assumption of pixel-wise conditional independence within the segmentation mask gives the log-likelihood function  $\log p(y|x, \theta) = \sum_{s=1}^S \log p(y_s|x_s, \theta)$ , where  $s$  indexes pixels. If we assume a flat prior  $p(\theta) = 1$ , we recover the frequently used negative log-likelihood loss

$$L(\theta) = -\log p(D|\theta) = -\sum_{i=1}^N \sum_{s=1}^S \log p(y_{s,i}|x_{s,i}, \theta), \quad (2)$$

where  $i$  indexes the training data. Note that this equivalent to the cross-entropy loss for a Bernoulli likelihood.

While a joint modelling objective for aleatoric and epistemic uncertainty for regression and classification tasks in computer vision was suggested by Kendall and Gal (2017), most research has focused on either one or the other (Kohl et al., 2018; Monteiro et al., 2020).

Methods that mainly focus on modelling aleatoric uncertainty in image segmentation have been proposed based on mixing deterministic segmentation architectures with generative components such as variational autoencoders and normalizing flows (Baumgartner et al., 2019; Kohl et al., 2019; Selvan et al., 2020). Probabilistic graphical models and combinations with neural networks have also been introduced for this purpose, but their application is restricted due to computational expense during inference time (Batra et al., 2012; Kirillov et al., 2015, 2016; Arnab et al., 2018; Kamnitsas et al., 2017; Kirillov et al., 2015). Ensemble and multi-head models (Lakshminarayanan et al., 2017; Rupprecht et al., 2017; Lee et al., 2016) have been applied for both aleatoric and epistemic uncertainty quantification, training on multiple- or same annotations for aleatoric or epistemic uncertainty quantification, respectively. In particular, this results in a frequentist approach to estimate  $p(\theta|D)$ . A Bayesian way to approximate the posterior is Monte-Carlo Dropout, a variational method based on Bernoulli distributions, which can be easily implemented for neural networks (Gal and Ghahramani, 2016, 2015, 2016). Model uncertainty is retrieved by averaging multiple forward passes, while setting weights randomly to 0 with probability  $p$ , both during inference time and training. Approaches that model both types of uncertainty in one model have been introduced based on the combination of Mean-Variance networks with diagonal covariance matrix and Dropout (Kendall and Gal, 2017) and Gaussian-Process based convolutional layers (Popescu et al., 2021).

## 2.1. Post-hoc Laplace Approximation

For Bayesian Neural Networks (Bishop, 2006) Laplace’s method (MacKay, 1992) is used to approximate the intractable posterior distribution  $p(\theta|D)$  with a Gaussian centered at a local mode  $\theta_{\text{MAP}}$ . In a first step the mode of the posterior can be determined by iterative numerical optimization such as stochastic or conjugate gradient decent. Note that for the binary cross entropy loss case described in Eq. (2) maximizing the log posterior is equivalent to minimizing the loss function. A second order Taylor expansion around  $\theta_{\text{MAP}}$  gives

$$\log p(\theta|D) \simeq \log(\theta_{\text{MAP}}|D) - \frac{1}{2}(\theta - \theta_{\text{MAP}})^\top \mathbf{H}(\theta - \theta_{\text{MAP}}) \quad (3)$$

where the first derivative vanishes because  $\theta_{\text{MAP}}$  is assumed to be a local mode. The Hessian  $\mathbf{H}$  in Eq. (3) is defined as

$$\mathbf{H} = -\nabla_\theta \nabla_\theta \log p(\theta|D) \Big|_{\theta=\theta_{\text{MAP}}}. \quad (4)$$

Applying the exponential function to Eq. (3) gives an approximation  $q(\theta)$  which is proportional to the density function of a multivariate normal distribution. The normalizing factor can be found by evaluating the determinant of  $\mathbf{H}$  giving the final approximation to the posterior as

$$q(\theta) = \mathcal{N}(\theta|\theta_{\text{MAP}}, \mathbf{H}^{-1}). \quad (5)$$

Evaluating the Hessian for such large functions is computationally infeasible because of the quadratic complexity in the network parameters and the usually large output dimensions in vision tasks. Approximations exist to overcome this (Botev, 2020) and frameworks for applying them are available (Detlefsen et al., 2021; Daxberger et al., 2021).

We propose a fast Hessian approximation for models with skip connections (Sec. 3.2) that scales linearly with parameters and output resolution and makes a post-hoc Laplace approximation of the parameter posterior feasible. Further, we extend the mean-variance network in Kendall and Gal (2017) from heteroscedastic pixel-wise to spatially correlated aleatoric uncertainty using SSN. We quantify epistemic uncertainty by approximating the posterior on the mean network, and show that this yields reliable OOD detection.

## 3. Laplacian Segmentation Network

To model the epistemic uncertainty in Eq. (1) with a Laplace approximation, we show that recent progress in scaling LA to images (Miani et al., 2022) can be extended to segmentation networks using skip connections.

### 3.1. Laplace Approximation of Mean Network

We can reformulate the integral for the predictive distribution in Eq. (1) by integrating over logits  $\eta$  to obtain

$$p(y|x, D) = \iint p(y|\eta)p(\eta|x, \theta)p(\theta|D) d\eta d\theta. \quad (6)$$

This further factorization of Eq. (1) corresponds to the model formulation used in Evidential Deep Learning (Sensoy et al., 2018; Ulmer et al., 2023). Following Monteiro et al. (2020) and Kendall and Gal (2017), we model the conditional distribution over logits  $p(\eta|x, \theta)$  as a normal distribution parametrized by neural networks  $\mu$  and  $\Sigma$  :

$$\eta|x \sim \mathcal{N}(\mu(x, \theta_1), \Sigma(x, \theta_2)), \quad (7)$$

and assume pixel-wise independence for the predicted labels given the logits. Thus, we can model  $p(y|\eta)$  for each pixel  $i$  as a Bernoulli distribution parametrized by the softmax of the respective logit. Since the size of the covariance matrix  $\Sigma$  scales quadratically with the number of pixels in the image, we use the low-rank parameterisation of Monteiro et al. (2020)

$$\Sigma(x) = D(x) + P(x)^T P(x), \quad (8)$$

i.e. the variance network  $\Sigma(x)$  is implemented with two networks  $D(x)$  and  $P(x)$ .

The vectors  $\theta_1 \in \Theta_1 = \mathbb{R}^T$  and  $\theta_2 \in \Theta_2 = \mathbb{R}^T$  parameterize the mean and variance networks (c.f. Eq. 7) and share the first  $t$  entries, i.e. we define the shared weight vector  $\theta_t$  of the network by

$$\theta_t := (\theta_{1_1}, \dots, \theta_{1_t}) = (\theta_{2_1}, \dots, \theta_{2_t}) \in \Theta_t = \mathbb{R}^t. \quad (9)$$

Then  $\theta \in \Theta = \mathbb{R}^{(t+2 \cdot (T-t))}$  contains all model parameters

$$\theta := (\theta_t, \theta_{1_{t+1}}, \dots, \theta_{1_T}, \theta_{2_{t+1}}, \dots, \theta_{2_T}). \quad (10)$$

The post-hoc Laplace approximation first finds a mode  $\theta_{\text{MAP}}$  by minimizing the loss function

$$\begin{aligned} \mathcal{L}(\theta) &= -\log \mathbb{E}_{p(\eta|x, \theta)} [p(y|\eta)] - \log p(\theta) \approx \\ &= -\log \text{sumexp}_{m=1}^M \left( \sum_{s=1}^S \log p(y_s | \eta_s^{(m)}) \right) + \log(M), \end{aligned} \quad (11)$$

where  $M$  logits  $\eta$  are sampled from the distribution in Eq. (7) and where the term  $-\log p(\theta)$  vanishes assuming a flat prior  $p(\theta) = 1$ . Since current algorithms for fast Hessian computations have no implementation for this loss function, we instead make use of the shared weights in the parameter vectors. For  $t \gg T - t$ , the loss landscape is dominantly defined by the shared parameters. Current SSN implementations usually fulfil this criteria, since they estimate the mean and variance of the logit distribution based on the feature maps of a deep deterministic segmentation model, while using only one convolutional layer each for mean and variance estimation. We therefore discard the entries of the variance network on the parameter vector  $\theta_{\text{MAP}}$ , i.e. we set

$$\theta_{\text{MAP}}^* := \theta_{\text{MAP}} \Big|_{(\theta_t, \theta_{1_{t+1}}, \dots, \theta_{1_T})} \in \Theta_{\text{mean}} = \mathbb{R}^T. \quad (12)$$

We can make use of the fact that the SSN loss function reduces to the binary cross entropy loss under zero variance, which allows us to fall back on the fast Hessian computation frameworks available. The posterior is then found by the Laplace approximation as described in Sec. 2.1 resulting in a Gaussian approximation in the parameter space  $\Theta_{\text{mean}}$

$$q(\theta^*) = \mathcal{N}(\theta^* | \theta_{\text{MAP}}^*, \mathbf{H}^{*-1}), \quad (13)$$

with  $\mathbf{H}^*$  defined as

$$\mathbf{H}^* = -\nabla_{\theta} \nabla_{\theta} \log p(\theta^* | D) \Big|_{\theta^* = \theta_{\text{MAP}}^*}. \quad (14)$$

Additionally to the aleatoric logit distribution, we can now investigate the epistemic component in form of the Laplace approximation for a given sample during inference time. Figure 2 gives an schematic overview of the proposed Laplacian Segmentation Network (LSN).

### 3.2. Fast Hessian Approximations for Segmentation Networks with Skip Connections

Computation of second order derivatives for Segmentation Networks is expensive due to the vast amount of parameters and pixels in the output. Standard methods approximate the Hessian with the diagonal of the Generalized Gauss Newton (GGN) matrix (Foresee and Hagan, 1997; Botev, 2020). This approximation, besides enforcing positive definiteness, also allows for an efficient backpropagation-like algorithm. The required compute scales linearly in the number of parameters and quadratic in the number pixels. The quadratic dependency is prohibitive already with images of size  $64 \times 64$ . We therefore make use of the diagonal backpropagation (DB) proposed by Miani et al. (2022), which returns a trace-preserving approximation of the diagonal of the GGN. The complexity of this approximation scales linearly with the number of pixels, allowing the computation of the Hessian also for larger images. The idea is to add a diagonal operator  $\mathfrak{D}$  in-between each backpropagation step. For each layer  $l$

$$\begin{aligned} [\nabla_{\theta} \nabla_{\theta} \log p(\theta | D)]_l &\stackrel{\text{GGN}}{\approx} [J_{\theta} f_{\theta}(x)^{\top} H^{(L)} J_{\theta} f_{\theta}(x)]_l = \quad (15) \\ &= J_{\theta} f^{(l)\top} \left( \prod_{i=l+1}^L J_x f^{(i)\top} H^{(L)} \prod_{i=L}^{l+1} J_x f^{(i)} \right) J_{\theta} f^{(l)} \\ &\stackrel{\text{DB}}{\approx} J_{\theta} f^{(l)\top} \mathfrak{D} \left( J_x f^{(l+1)\top} \mathfrak{D} (\dots) J_x f^{(l+1)} \right) J_{\theta} f^{(l)} \end{aligned}$$

where  $H^{(L)}$  is the Hessian of the binary cross entropy loss with respect to the logits, which can be expressed in closed form as diagonal plus outer product matrix.

Moreover, we extend the `StochMan` library (Detlefsen et al., 2021) with support for skip-connection layers. For a given submodule  $f_{\theta}$ , a skip-connection layer  $SC_f$  concatenates the function with the identity, such that

$SC_f(x) = (f_{\theta}(x), x)$ . The Jacobian is then  $J_x SC_f(x) := (J_x f_{\theta}(x), \mathbb{I}_x)$ . We exploit the block structure and efficiently backpropagate the diagonal only. With a recursive call on the submodule  $f$ , the backpropagation supports nested skip-connections, i.e. when some submodules of  $f$  are skip-connections as well. This unlocks the use of various curvature-based methods for segmentation architectures with skip connection in future research. For a technical description of the used Hessian approximation we refer to A.2.

## 4. Results

In the following, we investigate the derived LSN and test it in different practical scenarios. We start by illustrating the effect of introducing spatial correlation between pixels in scenarios with strong class imbalance. When segmenting small objects from background, reweighting classes is often necessary to optimize the network; this, however, has an undesired effect on the estimated epistemic uncertainties. We show how LSNs incorporating spatial correlation do not have this problem, and illustrate how this results in epistemic uncertainty with superior OOD behavior.

Next, we compare how the two types of uncertainties behave when modelled together and how and whether they differ. Finally we investigate out-of-distribution capabilities of the considered models on three different tasks. First we show with a toy example to which extent the epistemic components of the models are able to detect corrupted white noise boxes that are added to the images. Second, we assess the ability to detect distribution shifts by comparing the ISIC dataset against two similar datasets containing images of skin without lesions.

### 4.1. Baselines & Data

We compare our suggested model (LSN) to the following baselines in our experiments: (1) Ensemble of U-nets (Ensemble), where each member is trained on the same dataset (2) U-net with Monte-Carlo Dropout (U-net + Dropout) (3) U-net with post-hoc Laplace approximation (U-net + LA) (4) SSN with Dropout on the mean network (SSN + Dropout). Additionally, we consider the case of a diagonal covariance matrix, ignoring spatial correlation, for LSN and SSN + Dropout. We indicate this with a (*diag*) tag.

Training and experimental scripts are implemented in PyTorch and will be made available upon publication. All models share the same U-net backbone with five encoding blocks for comparability. Each block contains two convolution layers and uses the hyperbolic tangent activation function. For the estimation of  $\mu$ ,  $D$  and  $P$  the last feature map of the U-net is passed through three separate  $1 \times 1$  convolution layers, respectively. The ISIC19 dataset (Combalia et al., 2019; Codella et al., 2018; Tschandl et al., 2018) is scaled to a resolution of  $64 \times 64$  and is split into 11028 images for training and 1379 images for each, validation and

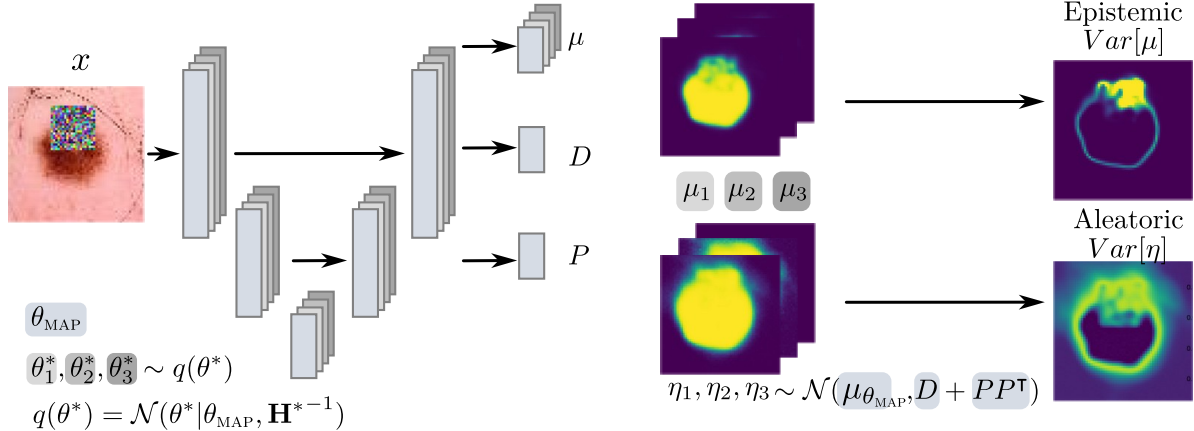


Figure 2: **Model overview.** Epistemic uncertainty maps are retrieved by sampling mean networks from the Laplace approximation  $q(\theta^*)$  and calculating the variance on their outputs for  $x$ . The aleatoric logit distribution is predicted on the parameter configuration  $\theta_{\text{MAP}}$  akin to Monteiro et al. (2020).

Table 1: Modelling the spatial correlation in the aleatoric logit distribution makes scaling the binary cross entropy loss unnecessary. IoU on the ISIC validation set after training for all models and different scaling factors of the loss function. Other learning parameters are kept equal.

Model	Scaling Factor			
	1	2	4	8
Ensemble	$< 10^{-10}$	$< 10^{-10}$	0.74	0.72
U-net + Dropout	$< 10^{-10}$	$< 10^{-10}$	0.64	0.61
SSN (diag) + Dropout	$< 10^{-10}$	$< 10^{-10}$	0.69	0.65
SSN + Dropout	0.73	0.72	0.68	0.60
U-net + LA	$< 10^{-10}$	$< 10^{-10}$	0.67	0.70
LSN (diag)	$< 10^{-10}$	$< 10^{-10}$	0.67	0.65
LSN	0.71	0.65	0.58	0.59

testing. We refer to appendix A.1 for further details on the training procedure.

## 4.2. Spatial Correlation affects Class Imbalance affects Epistemic Uncertainty

Class imbalance is a common challenge when segmenting small objects, often treated with methods of cost sensitive learning (Elkan, 2001; Kukar et al., 1998) deriving loss functions that penalize wrongly classified pixels of different classes differently (Lin et al., 2017). For the binary cross entropy loss, this is typically done by scaling the class-wise loss contributions with a factor that represents the imbalance between target and background.

To compensate for class imbalance in the ISIC dataset (about 1:8), all models were trained with binary cross en-

trophy, reweighting the underrepresented class with different factors. Table 1 shows the validation set intersection-over-union (IoU) after the last epoch for different scaling factors of the loss function. We find that only the models including spatial variation reach a useful optimum when reweighting is not applied, illustrating that including spatial correlation into the modelling has a desired effect on training with imbalanced classes. We hypothesize that this is because the high correlation between pixels in the background effectively reduces their contribution to the loss.

But this spatial correlation does not only affect our ability to optimize models; this also affects our epistemic uncertainty estimates: The alternative rescaling of parts of the loss has an immediate effect on epistemic uncertainty, since the Gaussian Approximation found post-hoc depends on the curvature of the loss landscape. In Fig. 3 we illustrate this effect with a toy example, showing the landscape a function (left) and of a rescaled version of the same function (right), along with their post-hoc fitted Laplace approximations. As shown in the illustration, the scaled function exhibits higher curvature, and hence a lower epistemic uncertainty, than the original. In other words, we find that including spatial correlation between pixels in the aleatoric logit distribution should affect the resulting epistemic uncertainty given by a post-hoc Laplace approximation. In the next sections, we shall see this empirically by comparing LSN to the version LSN (diag), where the aleatoric uncertainty is modelled without spatial correlation.

## 4.3. How do aleatoric and epistemic uncertainty differ when modelled jointly?

The proposed model can quantify aleatoric and epistemic uncertainty via the logit distribution of the SSN and the post-hoc Laplace approximation. To evaluate how similar

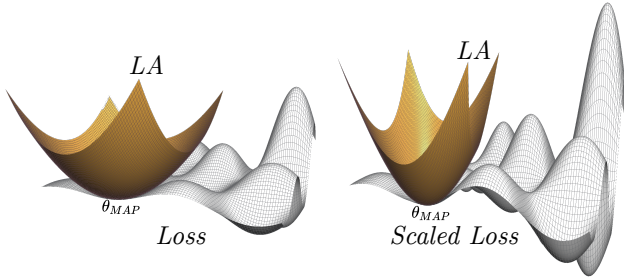


Figure 3: Illustration of loss landscape (grey) with fitted post-hoc Laplace approximation (orange) around a local minimum  $\theta_{MAP}$ . On the right the loss landscape is scaled by a factor  $> 1$ , as done when correcting for class imbalance. Note that the found Gaussian approximation’s variance depends on the curvature at  $\theta_{MAP}$ .

the modelled distributions are we evaluate them in terms of segmentation performance. For each image in the ISIC test set, we draw 50 samples from the logit distribution, parameterized by  $\theta_{MAP}$ . We then draw 50 mean networks from the parameter distribution given by the Laplace approximation. Retrieving samples from both components of the model is illustrated in Fig. 2. Figure 4 shows Precision-Recall curves for the proposed model, as well as the baseline models that use the SSN and therefore model both uncertainties. Note, that the SSN (diag) + Dropout corresponds to the heteroscedastic Bayesian Neural Network with Dropout derived by Kendall and Gal (2017).

First, we find that the aleatoric component, targeting the data variation, performs better in terms of segmentation performance than the epistemic component for our model. Further, the epistemic component of the LSN with diagonal covariance matrix performs worse (lower AUC), underpinning that the scaling of the loss has an effect on the post-hoc Laplace approximation. For the Dropout-based methods, on the other hand, we see virtually no difference between the epistemic and aleatoric components. This is in line with the findings of Kendall and Gal (2017), that Dropout as the epistemic component imitates the aleatoric logit distribution, unable to capture the different types of uncertainty.

What does this tell us? We see that the aleatoric and epistemic components of the LSN clearly capture different aspects of the model’s uncertainty. To better understand where this difference comes from, Figure 5 shows epistemic and aleatoric variance maps of sigmoids for a given ISIC test set sample. While the epistemic component assigns low variance to this in-distribution sample, the aleatoric component expresses the variation of plausible predictions for the image. We find this behavior attractive: In the Bayesian framework, the aleatoric component should play the role of the predictor, as it imitates the distribution of annotations

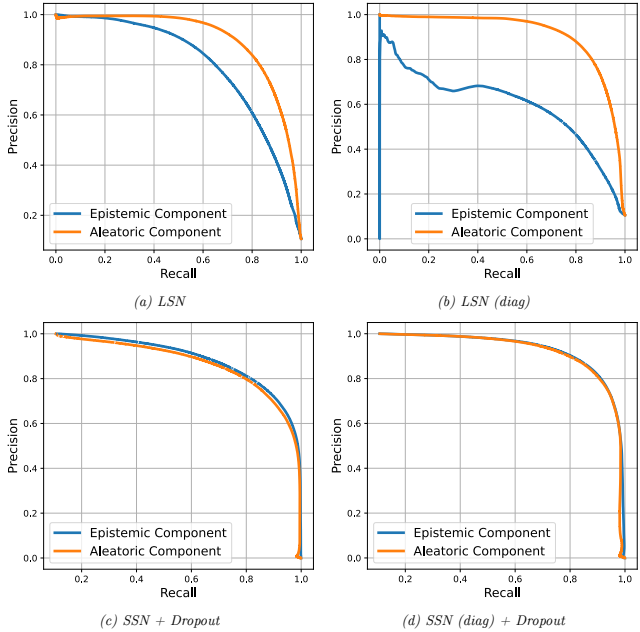


Figure 4: Precision-Recall plots for the ISIC test dataset show that aleatoric and epistemic variance in the LSN models (first row) are not the same, opposed to the SSN + Dropout models. The better performance of the aleatoric component (higher AUC) indicates that the model learns to separate data variation from model uncertainty. This is not the case for models with Dropout as an epistemic component, confirming the findings by Kendall and Gal (2017).

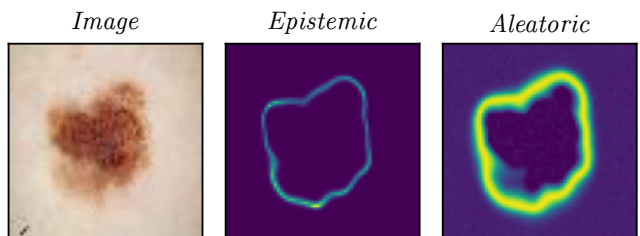


Figure 5: Epistemic and aleatoric uncertainties differ for the LSN model. While the aleatoric component assigns high variance to the boundary, expressing how plausible segmentations could differ, the epistemic variance for this in-distribution sample is low.

found in the data. The epistemic component, conversely, should capture the model uncertainty, enabling the model to identify OOD samples. In the following experiments we will investigate to which extent the proposed model’s epistemic component accomplishes this requirement.

#### 4.4. Detection of Within-Image Corruption

To assess how well the different models are able to assign high epistemic variance to corrupted parts of an image, we distort in-distribution images by either adding squared white noise boxes, or cropping out black boxes, to the samples of the ISIC test set. This allows us to directly compare the epistemic components’ variance predictions for ID and OOD pixels and investigate whether they are able to assign high variance to the affected area within the image. Figure 6 shows a sample from the ISIC dataset with its corrupted counterpart, along with the respective models’ epistemic variance maps. The ensemble assigns high variance to the whole lesion since some members end up in suboptimal local minima and predict background only. Any meaningful epistemic uncertainty captured by the remaining models is clouded by this. While all models assign variance to the boundary of the lesion, the LSN assigns the variance more accurately to the white noise box. To quantify whether the models assign higher variance to the OOD samples compared to their ID counterparts, we define the Pixel Ratio for a given image  $x$  as

$$Pixel\ Ratio(x) = \frac{\sum_{s=1}^S \text{Var}(\sigma(\eta_s(x_{OOD})))}{\sum_{s=1}^S \text{Var}(\sigma(\eta_s(x)))}, \quad (16)$$

where  $S$  is the number of pixels and  $x_{OOD}$  the corrupted version of the image. While the Pixel Ratio detects increased variance across the entire image, it does not quantify whether the increase in prediction variance is localized at the corrupted box. To this end, we define the Box Ratio as

$$Box\ Ratio(x_{OOD}) = \frac{\sum_{k=1}^K \text{Var}(\sigma(\eta_k(x_{OOD})))}{\sum_{s=1}^S \text{Var}(\sigma(\eta_s(x_{OOD})))} \quad (17)$$

where  $K$  is the index set for those pixels lying within the white noise box. The Box Ratio gives an indication to what extent the considered method assigns high variance to those pixels that have been distorted. In Table 2, we provide both measures for the case of an added white noise box or a black crop box, where size and position of the box vary randomly for each image in the ISIC test set. Specifically, we sample positions  $[0, 40]$ ,  $[0, 40]$  and side lengths  $[10, 20]$  of the boxes uniformly. All models assign higher variance to the OOD samples, indicated by Pixel Ratio values larger than 1. However, the LSN is able to better locate the corrupted boxes on average. The total sum of variance from ID to OOD does not increase much for the LSN, indicated by low Pixel Ratios. At the same time the high Box Ratio indicates that the LSN moves the variance from the boundary of the lesion to the corrupted box. This is also visible in the qualitative result in Fig. 6. This is illustrated by the examples in Figure 6, where we for instance see clearly that while the U-net + Dropout gets a top score on the Pixel Ra-

Table 2: Average ratio of pixel sum of variance maps for ID and OOD samples (Pixel Ratio) and ratio between correctly assigned variance to the corrupted box and pixel sum of variance maps (Box Ratio) for all models on ISIC. Results are shown for OOD images generated by adding white noise box (Noise) or a black crop box (Black) at a randomly assigned position and size.

Model	Pixel Ratio $\uparrow$		Box Ratio $\uparrow$	
	Black	Noise	Black	Noise
Ensemble	1.00	0.95	0.14	0.13
U-net + Dropout	<b>1.82</b>	<b>1.80</b>	0.13	0.14
SSN (diag) + Dropout	1.09	1.20	0.15	0.17
SSN + Dropout	1.51	1.60	0.14	0.20
U-net + LA	1.14	1.47	0.17	0.24
LSN (diag)	1.33	1.23	0.23	0.24
LSN	1.29	1.64	<b>0.26</b>	<b>0.29</b>

tio, its low Box Ratio score uncovers the fact that the model did not really localize the corrupted box at all.

#### 4.5. Distribution Shifts on ISIC dataset

Detecting gradual distribution shifts is critical in medical imaging, since downstream tasks might be sensitive to small systematic input changes. For example, a slightly flatter camera angle directly influences predicted area of a skin lesion, distorting decisions in automated systems. We evaluate OOD performance of all models, trained on the ISIC 2018 dataset and evaluated on an in-distribution test set and three OOD datasets: Derm-Skin (DERM), Clin-Skin (CLINIC) (Pacheco et al., 2020) and the PAD-UFES-20 dataset (Pacheco et al., 2020). The Derm-Skin dataset contains 1,565 images of healthy skin, cropped out of the ISIC dataset. The Clin-Skin dataset contains 723 images showing healthy skin gathered from social networks. The PAD-UFES-20 dataset contains 1570 photos of skin lesions collected from smartphone cameras. Figure 7 shows OOD detection capabilities for all models, while investigating the effects of different modeling on the epistemic component’s predictive entropy. The distributions reported in Figure 7 relate to the average per pixel predictive entropy per image.

The models with the best separation between ID and OOD entropy are the Ensemble, SSN (diag) + Dropout, LSN and LSN (diag) methods. However, even among the best performing methods, there remains significant overlap in predictive entropy between the ID and OOD datasets. This lack in separation motivates future work in exploring alternative uncertainty quantification metrics.

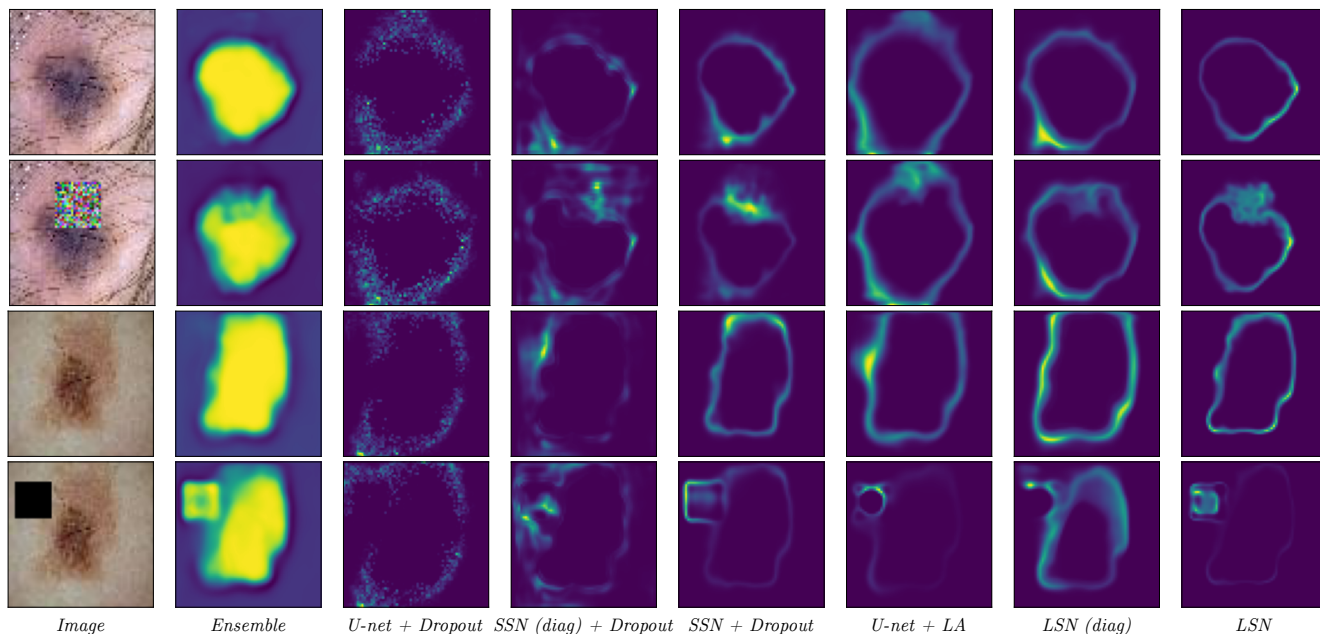


Figure 6: Predicted epistemic variance maps of two samples from the ISIC test set and counterparts with added white noise box and black crop for the different models. LSN assigns high variance to the white noise box as well as the black crop.

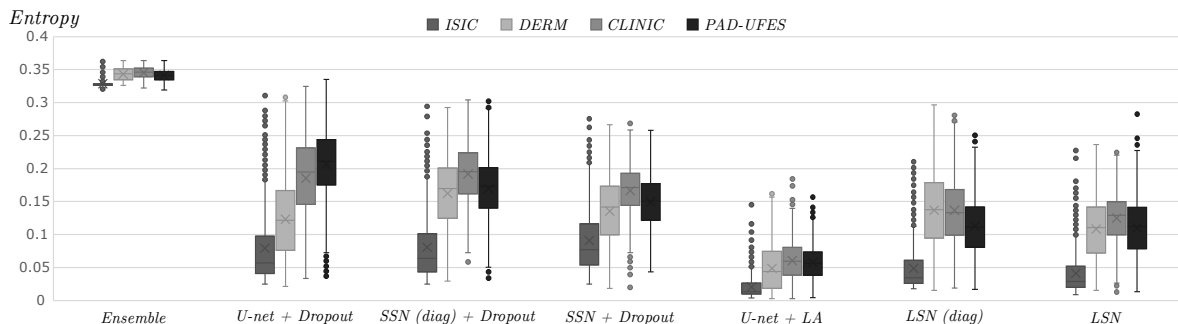


Figure 7: Box plot of entropy for all models on the in-distribution ISIC test set and the OOD datasets DERM, CLINIC and PAD-UFES-20. All models assign higher entropy to the OOD datasets.

## 5. Discussion and Conclusion

In this paper, we have demonstrated how Laplace approximations can scale to image segmentation tasks, through a trace-preserving diagonal Hessian approximation. Importantly, this scales linearly with the number of image pixels, unlike past work which exhibited a quadratic complexity. We have demonstrated across different datasets and quality measures that this successfully captures the epistemic uncertainty of estimated model parameters, in contrast to existing methods.

While we have illustrated a promising potential for the LSN models, our current implementation still has some limitations. First, we currently tackle only binary segmentation problems, which are common in medical imaging. Moreover, our Laplace approximation currently estimates uncertainty in the mean function of a mean-variance network.

However, a fully Bayesian approach would apply a Laplace approximation to the entire model, giving even more precise estimates of uncertainty.

We make the interesting experimental discovery that accurately capturing aleatoric uncertainty is essential to recover useful epistemic uncertainty estimates. In particular, we find that an explicit model of spatial correlation is essential. One hypothesis is that the epistemic uncertainty is forced to explain whichever uncertainty is left unexplained by the aleatoric component. A coarse (aleatoric) likelihood model that ignores spatial correlation, thus, force the epistemic component to capture spatial data correlations even if this is not epistemic. Such a hypothesis would suggest that richer correlation structures than those considered here (e.g. based on attention mechanisms (Vaswani et al., 2017)) could yield even finer grained epistemic uncertainties.



## **Acknowledgements**

This work was supported by a research grant (42062) from VILLUM FONDEN. This project received funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme (grant agreement 757360). The work was partly funded by the Novo Nordisk Foundation through the Center for Basic Machine Learning Research in Life Science (NNF20OC0062606). The authors acknowledge the Pioneer Centre for AI, DNRF grant number P1. This research used resources provided by the Darwin testbed at Los Alamos National Laboratory (LANL) which is funded by the Computational Systems and Software Environments subprogram of LANL’s Advanced Simulation and Computing program (NNSA/DOE). This work was supported by the Laboratory Directed Research and Development program of LANL under project number 20210043DR. LANL is operated by Triad National Security, LLC, for the National Nuclear Security Administration of the U.S. Department of Energy (Contract No. 89233218CNA000001).

## A. Appendix

### A.1. Implementation and Training Details

All architectures were trained using the Adam optimizer (Kingma and Ba, 2014) and its default PyTorch configuration. For models with Dropout the learning rate was 0.005, for all other models the learning rate was 0.0001. We used a batch size of 32 and trained all models for 60 epochs. These hyperparameters were found with a grid search on the validation set. Computations were performed on an internal GPU cluster. To retrieve epistemic uncertainties we sample 50 samples from the parameter distributions given by Dropout and the post-hoc Laplace approximation.

### A.2. Fast Hessian Approximation

The aim of this section is to clarify the following sequence of approximations for the Hessian  $H_\theta$

$$H_\theta \stackrel{(1)}{\approx} \text{GGN}_\theta \stackrel{(2)}{\approx} \text{GGND}_\theta \stackrel{(3)}{\approx} \text{DB}_\theta$$

where  $\text{GGN}_\theta$  and  $\text{GGND}_\theta$  refer to the Generalized Gauß Newton matrix and its Diagonal (Botev, 2020) and  $\text{DB}_\theta$  is the matrix found by the diagonal backpropagation algorithm (Miani et al., 2022). Each of the above approximations trade exactness of the found solution for speed of its computation. Specifically, the approximation (1) yields a matrix  $\text{GGN}_\theta$  that is always positive semi-definite. Its diagonal given by (2) is commonly used for the Laplace approximation in neural networks, since its Inverse is computationally feasible, often referred to as diagonal Laplace. In the following we develop the terminology needed to explain our contribution in extending the diagonal backpropagation algorithm  $\text{DB}_\theta$  to nested skip-connection layers.

### A.3. Notation

Consider a neural network (NN)  $f_\theta : \mathcal{X} \rightarrow \mathcal{Y}$  with  $L$  layers. The parameter  $\theta = (\theta_1, \dots, \theta_L) \in \Theta$  is the concatenation of the parameters for each layer  $i \in \{1, \dots, L\}$ . The NN is a composition of  $L$  functions  $f^{(1)}, f^{(2)}, \dots, f^{(L)}$ , where  $f^{(i)}$  is parametrized by  $\theta_i$ .

$$f_\theta := f_{\theta_L}^{(L)} \circ f_{\theta_{L-1}}^{(L-1)} \circ \dots \circ f_{\theta_2}^{(2)} \circ f_{\theta_1}^{(1)}.$$

Since we need explicit access to the intermediate values, we call the input  $x_0 \in \mathcal{X}$  and iteratively define  $x_i := f_{\theta_i}^{(i)}(x_{i-1})$  for  $i = 1, \dots, L$ , such that the NN output is  $x_L \in \mathcal{Y}$ . This notation can be visually presented as

$$\begin{array}{ccccccc} \mathcal{X} & \xrightarrow{\quad f_\theta \quad} & & & & & \mathcal{Y} \\ x_0 & \xrightarrow{f_{\theta_1}^{(1)}} x_1 & \longrightarrow & \dots & \longrightarrow & x_{i-1} \xrightarrow{f_{\theta_i}^{(i)}} x_i & \longrightarrow & \dots & \longrightarrow & x_{L-1} \xrightarrow{f_{\theta_L}^{(L)}} x_L \end{array}$$

We highlight that there is no restriction of what a layer  $f^{(i)}$  can be, in particular, it can be a composition of functions itself, translating to an object of the *sequential* class in PyTorch. This recursive structure of the function, which is replicated by its implementation in PyTorch, will be key in extending the diagonal backpropagation method to nested skip-connections.

We define the *diagonal* operator  $\mathcal{D}$  for a quadratic matrix  $M$  of size  $m \in \mathbb{N}^+$

$$\mathcal{D} : \mathbb{R}^{m \times m} \rightarrow \mathbb{R}^{m \times m}$$

by

$$[\mathcal{D}(M)]_{ij} := \begin{cases} M_{ij} & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases} \quad \forall i, j = 1, \dots, m$$

Note that this is a trace-preserving operator, a property that will be inherited by our final approximation. Moreover, note that applying this operator induces a systematic bias, specifically it decreases the Von Neumann entropy of the matrix. This is due to a smoothing effect on the eigenspectrum.

#### A.4. Jacobian of the Neural Network

We are interested in the Jacobian  $J_\theta f_\theta(x_0)$  of the NN with respect to the parameter  $\theta$ . Each column of the Jacobian is the derivative of the output vector w.r.t. a single parameter. We can then group the parameters (i.e. columns) layer by layer

$$\begin{aligned} J_\theta f_\theta(x_0) &= \left( J_{\theta_1} f_\theta(x_0) \mid \dots \mid J_{\theta_i} f_\theta(x_0) \mid \dots \mid J_{\theta_L} f_\theta(x_0) \right) \\ &= \left( J_{\theta_1} \left( f_{\theta_1}^{(1)} \circ \dots \circ f_{\theta_L}^{(L)} \right) (x_0) \mid \dots \mid J_{\theta_i} \left( f_{\theta_i}^{(i)} \circ \dots \circ f_{\theta_L}^{(L)} \right) (x_{i-1}) \mid \dots \mid J_{\theta_L} f_{\theta_L}^{(L)}(x_{L-1}) \right), \end{aligned}$$

where the second equality comes from the fact that each layer only depends on its respective parameters, i.e.

$$J_{\theta_j} f_{\theta_i}^{(i)}(x_{i-1}) = 0 \quad \text{if } i \neq j.$$

Exploiting this block-structure, we can focus on a single layer Jacobian  $J_{\theta_i} f_\theta(x_0)$ , and concatenate them afterwards. With the chain rule we get

$$J_{\theta_i} f_\theta(x_0) = J_{\theta_i} \left( f_{\theta_i}^{(i)} \circ \dots \circ f_{\theta_L}^{(L)} \right) (x_{i-1}) = \left( \prod_{j=L}^{i+1} J_{x_{j-1}} f_{\theta_j}^{(j)}(x_{j-1}) \right) J_{\theta_i} f_{\theta_i}^{(i)}(x_{i-1}). \quad (18)$$

The intuition for the chain rule is that the Jacobian  $J_{\theta_i} f_\theta(x_0)$  for layer  $i$  is the composition of the Jacobians w.r.t. the *input*  $J_{x_{j-1}} f_{\theta_j}^{(j)}(x_{j-1})$  of subsequent layers  $j = L, L-1, \dots, i+2, i+1$ , times the Jacobian w.r.t. the *parameters*  $J_{\theta_i} f_{\theta_i}^{(i)}(x_{i-1})$  of the specific layer  $i$ . Thus, we can reuse computation for one layer to improve the computation of other layers, specifically the product of Jacobians w.r.t. the input.

#### A.5. Generalized Gauss Newton method

For the Laplace approximation of a given loss function  $\mathcal{L} : \mathcal{Y} \rightarrow \mathbb{R}$  we need to compute the Hessian of the composition of NN and loss with respect to the parameters  $\nabla_\theta^2 \mathcal{L}(f_\theta(x_0)) \in \mathbb{R}^{|\theta| \times |\theta|}$ . In the following we denote the length of vector  $v$  by  $|v|$ . According to the chain rule it holds, that

$$\underbrace{\nabla_\theta^2 \mathcal{L}(f_\theta(x_0))}_{=: H_\theta} = \underbrace{J_\theta f_\theta(x_0)^\top \cdot \nabla_{x_L}^2 \mathcal{L}(x_L) \cdot J_\theta f_\theta(x_0)}_{=: \text{GGN}_\theta} + \sum_{o=1}^{|x_L|} [\nabla_{x_L} \mathcal{L}(x_L)]_o \cdot \nabla_\theta^2 [f_\theta(x_0)]_o, \quad (19)$$

where  $[v]_o$  refers to the  $o$ -th component of vector  $v$ . In this sense, the Generalized Gauss-Newton matrix  $\text{GGN}_\theta$  is commonly used as a positive-definite approximation of the full Hessian  $H_\theta$ . The positive definiteness follows from the positive definiteness of the Hessian of the loss function with respect to the NN output  $H^\mathcal{L} := \nabla_{x_L}^2 \mathcal{L}(x_L) \in \mathbb{R}^{|x_L| \times |x_L|}$ , which holds for common losses like Mean Squared Error or Cross-Entropy.

Consider a layer-by-layer block structure of the Generalized Gauss-Newton

$$\text{GGN}_\theta = \begin{pmatrix} J_{\theta_1} f_\theta(x_0)^\top H^\mathcal{L} J_{\theta_1} f_\theta(x_0) & J_{\theta_1} f_\theta(x_0)^\top H^\mathcal{L} J_{\theta_2} f_\theta(x_0) & \dots & J_{\theta_1} f_\theta(x_0)^\top H^\mathcal{L} J_{\theta_L} f_\theta(x_0) \\ J_{\theta_2} f_\theta(x_0)^\top H^\mathcal{L} J_{\theta_1} f_\theta(x_0) & J_{\theta_2} f_\theta(x_0)^\top H^\mathcal{L} J_{\theta_2} f_\theta(x_0) & & \\ \vdots & & \ddots & \\ J_{\theta_L} f_\theta(x_0)^\top H^\mathcal{L} J_{\theta_1} f_\theta(x_0) & & & J_{\theta_L} f_\theta(x_0)^\top H^\mathcal{L} J_{\theta_L} f_\theta(x_0) \end{pmatrix}.$$

Its *block-diagonal* approximation

$$\text{GGNB}_\theta = \begin{pmatrix} J_{\theta_1} f_\theta(x_0)^\top H^\mathcal{L} J_{\theta_1} f_\theta(x_0) & 0 & \dots & 0 \\ 0 & J_{\theta_2} f_\theta(x_0)^\top H^\mathcal{L} J_{\theta_2} f_\theta(x_0) & & \\ \vdots & & \ddots & \\ 0 & & & J_{\theta_L} f_\theta(x_0)^\top H^\mathcal{L} J_{\theta_L} f_\theta(x_0) \end{pmatrix}$$

is usually considered as a cheaper-to-compute approximation of the full  $\text{GGN}_\theta$ . According to chain rule in Eq. (18), we can explicitly write the diagonal block  $\text{GGNB}_\theta^{(i)} = J_{\theta_i} f_\theta(x_0)^\top H^\mathcal{L} J_{\theta_i} f_\theta(x_0)$  of the  $i$ -th layer as

$$\text{GGNB}_\theta^{(i)} = J_{\theta_i} f_{\theta_i}^{(i)}(x_{i-1})^\top \left( \prod_{j=i+1}^L J_{x_{j-1}} f_{\theta_j}^{(j)}(x_{j-1})^\top \right) \cdot H^\mathcal{L} \cdot \left( \prod_{j=L}^{i+1} J_{x_{j-1}} f_{\theta_j}^{(j)}(x_{j-1}) \right) J_{\theta_i} f_{\theta_i}^{(i)}(x_{i-1}). \quad (20)$$

This can be rearranged as a sequence of  $(J^\top M J)$ -like operators iteratively applied to  $H^\mathcal{L}$ , as

$$\begin{aligned} \text{GGNB}_\theta^{(i)} = J_{\theta_i} f_{\theta_i}^{(i)}(x_{i-1})^\top & \left( J_{x_i} f_{\theta_{i+1}}^{(i+1)}(x_i)^\top \left( J_{x_{i+1}} f_{\theta_{i+2}}^{(i+2)}(x_{i+1})^\top \left( \dots \left( J_{x_{L-1}} f_{\theta_L}^{(L)}(x_{L-1})^\top \cdot H^\mathcal{L} \cdot \right. \right. \right. \\ & \left. \left. \left. \cdot J_{x_{L-1}} f_{\theta_L}^{(L)}(x_{L-1}) \right) \dots \right) J_{x_i} f_{\theta_{i+1}}^{(i+1)}(x_i) \right) J_{x_{i+1}} f_{\theta_{i+2}}^{(i+2)}(x_{i+1}) \right) J_{\theta_i} f_{\theta_i}^{(i)}(x_{i-1}). \end{aligned} \quad (21)$$

From this expression, we can build an efficient backpropagation-like algorithm to compute  $\text{GGNB}_\theta$ .

---

**Algorithm 1** Computation of  $\text{GGNB}_\theta$  (exact backpropagation)

---

```

M = Hℒ
for j = L, L - 1, ..., 1 do
  GGNBθ(j) = Jθj fθj(j)(xj-1)⊤ · M · Jθj fθj(j)(xj-1)
  M = Jxj-1 fθj(j)(xj-1)⊤ · M · Jxj-1 fθj(j)(xj-1)
end for
GGNBθ = (GGNBθ(1), ..., GGNBθ(L))
return GGNBθ

```

---

Note the two memory bottlenecks of this algorithm: after each step  $j$  the involved matrices have sizes

$$\text{GGNB}_\theta^{(j)} \in \mathbb{R}^{|\theta_j| \times |\theta_j|} \quad M \in \mathbb{R}^{|x_{j-1}| \times |x_{j-1}|}.$$

The former one is commonly avoided by just computing the diagonal of each block  $\text{GGNB}_\theta^{(j)}$ , and thus just computing the diagonal of the Generalized Gauss-Newton matrix, which we can refer to as  $\text{GGND}_\theta := \mathcal{D}(\text{GGN}_\theta) = \mathcal{D}(\text{GGNB}_\theta)$ . Using this approach together with a Laplace approximation of the loss landscape is called *diagonal Laplace* and scales linearly in the number of parameter  $|\theta|$ .

---

**Algorithm 2** Computation of  $\text{GGND}_\theta$  (exact backpropagation)

---

```

M = Hℒ
for j = L, L - 1, ..., 1 do
  GGNDθ(j) = ℰ ( Jθj fθj(j)(xj-1)⊤ · M · Jθj fθj(j)(xj-1) )
  M = Jxj-1 fθj(j)(xj-1)⊤ · M · Jxj-1 fθj(j)(xj-1)
end for
GGNDθ = (GGNDθ(1), ..., GGNDθ(L))
return GGNDθ

```

---

The memory bottlenecks of this modified version of Eq. (1) are

$$\text{GGND}_\theta^{(j)} \in \mathbb{R}^{|\theta_j|} \quad M \in \mathbb{R}^{|x_{j-1}| \times |x_{j-1}|}.$$

However, the quadratic dependence on the size of  $x_j$  is still potentially problematic. In an encoder-decoder architecture, as in our case, the maximum of this size is realized in the input  $x_0$  and output  $x_L$ . For image-based networks, this size corresponds to the number of pixels. With *reasonable* resources this backpropagation algorithm is feasible for images with size  $28 \times 28$ , but it became infeasible for larger images. In order to deal with high-resolution images we need to avoid the quadratic dependency on the number of pixels in the output.

### A.6. Fast diagonal backprop

The diagonal backpropagation computes an approximation of the Generalized Gauss-Newton. Inspired by Eq. (21), it is defined, for each layer  $i = 1, \dots, L$ , as

$$\text{DB}_\theta^{(i)} = J_{\theta_i} f_{\theta_i}^{(i)}(x_{i-1})^\top \mathcal{D} \left( J_{x_i} f_{\theta_{i+1}}^{(i+1)}(x_i)^\top \mathcal{D} \left( J_{x_{i+1}} f_{\theta_{i+2}}^{(i+2)}(x_{i+1})^\top \mathcal{D} \left( \dots \mathcal{D} (J_{x_{L-1}} f_{\theta_L}^{(L)}(x_{L-1})^\top \cdot \mathcal{D}(H^\mathcal{L}) \cdot \right. \right. \right. \\ \left. \left. \left. \cdot J_{x_{L-1}} f_{\theta_L}^{(L)}(x_{L-1}) \right) \dots \right) J_{x_i} f_{\theta_{i+1}}^{(i+1)}(x_i) \right) J_{x_{i+1}} f_{\theta_{i+2}}^{(i+2)}(x_{i+1}) \right) J_{\theta_i} f_{\theta_i}^{(i)}(x_{i-1}). \quad (22)$$

As before, we can build an efficient backpropagation-like algorithm from this expression to compute  $\text{DB}_\theta$ .

---

#### Algorithm 3 Computation of $\text{DB}_\theta$ (approximated backpropagation)

---

```

M = Hℒ
for j = L, L - 1, ..., 1 do
  DBθ(j) = ℒ ( Jθj fθj(j)(xj-1)⊤ · M · Jθj fθj(j)(xj-1) )
  M = ℒ ( Jxj-1 fθj(j)(xj-1)⊤ · M · Jxj-1 fθj(j)(xj-1) )
end for
DBθ = (DBθ(1), ..., DBθ(L))
return DBθ

```

---

The memory bottlenecks of this algorithm are

$$\text{DB}_\theta^{(j)} \in \mathbb{R}^{|\theta_j|} \quad M \in \mathbb{R}^{|x_{j-1}|}$$

thus allowing for linear scaling in both parameter and number of pixels in the output.

A key aspect, that enables this algorithm to be so efficient, is the simultaneous diagonalization and computation of the Jacobian product. Instead of first computing the full matrix and then discarding all non-diagonal elements, we directly compute only the diagonal elements of the products. However, this means that this operation needs to be implemented separately for every type of layer. Moreover, all diagonal matrices are stored implicitly, i.e. diagonal matrices are stored as vectors of the diagonal entries.

### A.7. Skip-connections

In PyTorch, objects of the class *module* can be nested, allowing to create classes in a tree-like structure. All modules, but the root of such a tree are usually referred to as *submodules*. It's then possible to form sequential neural networks by composing different *submodules*, each representing a (potentially parametric) function as part of the architecture. For a given submodule  $g_\theta : \mathbb{R}^I \rightarrow \mathbb{R}^O$ , a skip-connection layer  $\text{SC}(g)_\theta$  concatenates the submodule with the identity. It is defined as

$$\text{SC}(g)_\theta : \mathbb{R}^I \longrightarrow \mathbb{R}^{O+I} \quad (23) \\ x \longmapsto (g_\theta(x), x)$$

Note that this is a submodule itself, and its parameter  $\theta$  is the same as the parameter of the used submodule  $g_\theta$ .

The Jacobian with respect to the parameter is the same as the Jacobian of the submodule

$$J_{\theta} \text{SC}(g)_{\theta}(x) = J_{\theta} g_{\theta}(x) \in \mathbb{R}^{O \times |\theta|} \quad (24)$$

while the Jacobian with respect to the input is

$$J_x \text{SC}(g)_{\theta}(x) = \begin{pmatrix} J_x g_{\theta}(x) \\ \mathbb{I} \end{pmatrix} \in \mathbb{R}^{(O+I) \times I}, \quad (25)$$

where  $\mathbb{I} \in \mathbb{R}^{I \times I}$  is the identity matrix of the same size as the input.

In order to perform the fast diagonal backpropagation through this kind of layer we need to compute the  $\mathcal{D}(J^{\top} \cdot M \cdot J)$  operator efficiently. The one with respect to the parameter is straightforward: being the Jacobian identical to the one of the submodule  $g$ , we simply need to call the operator implemented on the submodule. The one with respect to the input, on the other hand, is slightly different. We can exploit the same block structure on the matrix  $M \in \mathbb{R}^{(O+I) \times (O+I)}$  as

$$M = \begin{pmatrix} M_{11} & M_{12} \\ M_{21} & M_{22} \end{pmatrix}, \quad \text{where } M_{11} \in \mathbb{R}^{O \times O}, M_{22} \in \mathbb{R}^{I \times I}. \quad (26)$$

We can then explicitly write

$$\begin{aligned} J_x \text{SC}(g)_{\theta}(x)^{\top} \cdot M \cdot J_x \text{SC}(g)_{\theta}(x) &= \begin{pmatrix} J_x g_{\theta}(x)^{\top} & | & \mathbb{I} \end{pmatrix} \begin{pmatrix} M_{11} & M_{12} \\ M_{21} & M_{22} \end{pmatrix} \begin{pmatrix} J_x g_{\theta}(x) \\ \mathbb{I} \end{pmatrix} \\ &= J_x g_{\theta}(x)^{\top} M_{11} J_x g_{\theta}(x) + M_{12} J_x g_{\theta}(x) + J_x g_{\theta}(x)^{\top} M_{21} + M_{22} \end{aligned}$$

and thus, assuming  $M$  is already in diagonal form, which implies  $M_{12} = 0, M_{21} = 0$ , we have

$$\mathcal{D}(J_x \text{SC}(g)_{\theta}(x)^{\top} \cdot M \cdot J_x \text{SC}(g)_{\theta}(x)) = \mathcal{D}(J_x g_{\theta}(x)^{\top} M_{11} J_x g_{\theta}(x)) + \mathcal{D}(M_{22}). \quad (27)$$

This means that we can efficiently perform this backpropagation by calling the same operator on the submodule  $g$ , which returns  $\mathcal{D}(J_x g_{\theta}(x)^{\top} M_{11} J_x g_{\theta}(x))$  and then adding the diagonal  $\mathcal{D}(M_{22})$ .

### A.7.1 Recursiveness

We highlight that this efficient implementation builds on a recursive structure, which unlocks several possibilities. The idea of the library is that every module has implemented the two backpropagation operators:  $\mathcal{D}(J_x^{\top} \cdot M \cdot J_x)$  with respect to the input and  $\mathcal{D}(J_{\theta}^{\top} \cdot M \cdot J_{\theta})$  with respect to the parameter. Most importantly, these operators have to share the same syntax across all modules. For basic modules such as LINEAR, CONV2D, RELU, TANH, these operators are hardcoded with tensor operations. The Skip-connection module however, has these operators coded with a recursive call on the submodule, as in Eq. (27). In the Sequential module the two operators are implemented with a backpropagation-like structure, as in algorithm (3), with a recursive call on all the submodules in the sequence (in reverse order in the for loop).

In our work we use this to implement the skip-connections of the U-net (Ronneberger et al., 2015). The submodule of a skip-connection is a sequential, which contains convolutions and other layers, as well as another skip-connection with the same structure.

## References

- Anurag Arnab, Shuai Zheng, Sadeep Jayasumana, Bernardino Romera-Paredes, Måns Larsson, Alexander Kirillov, Bogdan Savchynskyy, Carsten Rother, Fredrik Kahl, and Philip H.S. Torr. Conditional random fields meet deep neural networks for semantic segmentation: Combining probabilistic graphical models with deep learning for structured prediction. *IEEE Signal Processing Magazine*, 35(1):37–52, 2018. 2
- Dhruv Batra, Payman Yadollahpour, Abner Guzman-Rivera, and Gregory Shakhnarovich. Diverse M-Best solutions in Markov random fields. *Computer Vision (ECCV)*, pages 1–16, 2012. 2
- Christian F. Baumgartner, Kerem C. Tezcan, Krishna Chaitanya, Andreas M. Hötker, Urs J. Muehlematter, Khoschy Schawkat, Anton S. Becker, Olivio Donati, and Ender Konukoglu. PHiSeg: Capturing uncertainty in medical image segmentation. *Medical Image Computing and Computer Assisted Intervention (MICCAI)*, pages 119–127, 2019. 2
- Christopher M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer-Verlag, Berlin, Heidelberg, 2006. 2, 3
- Aleksandar Botev. *The Gauss-Newton matrix for deep learning models and its applications*. PhD thesis, UCL (University College London), 2020. 3, 4, 10
- Noel CF Codella, David Gutman, M Emre Celebi, Brian Helba, Michael A Marchetti, Stephen W Dusza, Aadi Kallou, Konstantinos Liopyris, Nabin Mishra, Harald Kittler, et al. Skin lesion analysis toward melanoma detection: A challenge at the 2017 international symposium on biomedical imaging (ISBI), hosted by the international skin imaging collaboration (ISIC). In *IEEE 15th International Symposium on Biomedical Imaging*, pages 168–172. IEEE, 2018. 4
- Marc Combalia, Noel CF Codella, Veronica Rotemberg, Brian Helba, Veronica Vilaplana, Ofer Reiter, Cristina Carrera, Alicia Barreiro, Allan C Halpern, Susana Puig, et al. BCN20000: Dermoscopic lesions in the wild. *arXiv preprint arXiv:1908.02288*, 2019. 4
- Erik Daxberger, Agustinus Kristiadi, Alexander Immer, Runa Eschenhagen, Matthias Bauer, and Philipp Hennig. Laplace redux—effortless Bayesian deep learning. In *NeurIPS*, 2021. 1, 3
- Nicki S. Detlefsen, Alison Pouplin, Cilie W. Feldager, Cong Geng, Dimitris Kalatzis, Helene Hauschultz, Miguel González-Duque, Frederik Warburg, Marco Miani, and Søren Hauberg. Stochman. *GitHub. Note: <https://github.com/MachineLearningLifeScience/stochman/>*, 2021. 3, 4
- Charles Elkan. The foundations of cost-sensitive learning. In *International joint conference on artificial intelligence*, volume 17, pages 973–978. Lawrence Erlbaum Associates Ltd, 2001. 5
- F Dan Foresee and Martin T Hagan. Gauss-Newton approximation to Bayesian learning. In *Proceedings of International Conference on Neural Networks (ICNN’97)*, volume 3, pages 1930–1935. IEEE, 1997. 4
- Moritz Fuchs, Camila Gonzalez, and Anirban Mukhopadhyay. Practical uncertainty quantification for brain tumor segmentation. In *Medical Imaging with Deep Learning*, 2021. 2
- Yarin Gal and Zoubin Ghahramani. Bayesian convolutional neural networks with Bernoulli approximate variational inference. *arXiv preprint arXiv:1506.02158*, 2015. 2
- Yarin Gal and Zoubin Ghahramani. Dropout as a Bayesian approximation: Appendix, 2016. 2
- Yarin Gal and Zoubin Ghahramani. Dropout as a Bayesian approximation: Representing model uncertainty in deep learning. *International Conference on Machine Learning*, pages 1050–1059, 2016. 2
- Jakob Gawlikowski, Cedric Rovile Njiteucheu Tassi, Mohsin Ali, Jongseok Lee, Matthias Humt, Jianxiang Feng, Anna Kruspe, Rudolph Triebel, Peter Jung, Ribana Roscher, et al. A survey of uncertainty in deep neural networks. *arXiv preprint arXiv:2107.03342*, 2021. 2
- Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. *arXiv preprint arXiv:1610.02136*, 2016. 1
- Tinu Theckel Joy, Suman Sedai, and Rahil Garnavi. Analyzing epistemic and aleatoric uncertainty for drusen segmentation in optical coherence tomography images. *arXiv preprint arXiv:2101.08888*, 2021. 2
- Konstantinos Kamnitsas, Christian Ledig, Virginia F.J. Newcombe, Joanna P. Simpson, Andrew D. Kane, David K. Menon, Daniel Rueckert, and Ben Glocker. Efficient multi-scale 3D CNN with fully connected CRF for accurate brain lesion segmentation. *Medical Image Analysis*, 36:61–78, 2017. 2
- Alex Kendall and Yarin Gal. What uncertainties do we need in Bayesian deep learning for computer vision? *Advances in Neural Information Processing Systems*, 30, 2017. 2, 3, 6
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 10
- Alexander Kirillov, Bogdan Savchynskyy, Dmitriy Schlessinger, Dmitry Vetrov, and Carsten Rother. Inferring M-Best diverse labelings in a single one. *International Conference on Computer Vision (ICCV)*, pages 1814–1822, 2015. 2
- Alexander Kirillov, Alexander Shekhovtsov, Carsten Rother, and Bogdan Savchynskyy. Joint M-best-diverse labelings as a parametric submodular minimization. *Advances in Neural Information Processing Systems*, 29, 2016. 2
- Alexander Kirillov, Dmytro Shlezinger, Dmitry P Vetrov, Carsten Rother, and Bogdan Savchynskyy. M-Best-Diverse labelings for submodular energies and beyond. *Advances in Neural Information Processing Systems*, 28, 2015. 2
- Armen Der Kiureghian and Ove Ditlevsen. Aleatory or epistemic? Does it matter? *Structural Safety*, 31(2):105–112, 2009. 2
- Simon Kohl, Bernardino Romera-Paredes, Clemens Meyer, Jeffrey De Fauw, Joseph R Ledsam, Klaus Maier-Hein, SM Eslami, Danilo Jimenez Rezende, and Olaf Ronneberger. A probabilistic U-Net for segmentation of ambiguous images. *Advances in Neural Information Processing Systems*, 31, 2018. 2
- Simon A. A. Kohl, Bernardino Romera-Paredes, Klaus H. Maier-Hein, Danilo Jimenez Rezende, S. M. Ali Eslami, Pushmeet Kohli, Andrew Zisserman, and Olaf Ronneberger. A hierarchical probabilistic U-Net for modeling multi-scale ambiguities. *arXiv preprint arXiv:1905.13077*, 2019. 2
- Matjaz Kukar, Igor Kononenko, et al. Cost-sensitive learning with neural networks. In *ECAI*, volume 15, pages 88–94. Citeseer, 1998. 5
- Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation us-

- ing deep ensembles. *Advances in Neural Information Processing Systems*, 30, 2017. 2
- Stefan Lee, Senthil Purushwalkam Shiva Prakash, Michael Cogswell, Viresh Ranjan, David Crandall, and Dhruv Batra. Stochastic multiple choice learning for training diverse deep ensembles. *Advances in Neural Information Processing Systems*, 29, 2016. 2
- Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017. 2, 5
- David JC MacKay. Bayesian interpolation. *Neural computation*, 4(3):415–447, 1992. 3
- Marco Miani, Frederik Warburg, Pablo Moreno-Muñoz, Nicke Skafté Detlefsen, and Søren Hauberg. Laplacian autoencoders for learning stochastic representations. In *Advances in Neural Information Processing Systems*, 2022. 3, 4, 10
- Miguel Monteiro, Loïc Le Folgoc, Daniel Coelho de Castro, Nick Pawłowski, Bernardo Marques, Konstantinos Kamnitsas, Mark van der Wilk, and Ben Glocker. Stochastic segmentation networks: Modelling spatially correlated aleatoric uncertainty. *Advances in Neural Information Processing Systems*, 33:12756–12767, 2020. 1, 2, 3, 5
- Andre GC Pacheco, Gustavo R Lima, Amanda S Salomao, Breno Krohling, Igor P Biral, Gabriel G de Angelo, Fábio CR Alves Jr, José GM Esgario, Alana C Simora, Pedro BC Castro, et al. Pad-ufes-20: A skin lesion dataset composed of patient data and clinical images collected from smartphones. *Data in brief*, 32:106221, 2020. 7
- Andre G. C. Pacheco, Chandramouli S. Sastry, Thomas Trappenberg, Sageev Oore, and Renato A. Krohling. On out-of-distribution detection algorithms with deep neural skin cancer classifiers. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 3152–3161, 2020. 7
- Sebastian G Popescu, David J Sharp, James H Cole, Konstantinos Kamnitsas, and Ben Glocker. Distributional Gaussian process layers for outlier detection in image segmentation. In *International Conference on Information Processing in Medical Imaging*, pages 415–427. Springer, 2021. 2
- Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. *CoRR*, abs/1505.04597, 2015. 1, 14
- Christian Rupprecht, Iro Laina, Robert DiPietro, Maximilian Baust, Federico Tombari, Nassir Navab, and Gregory D Hager. Learning in an uncertain world: Representing ambiguity through multiple hypotheses. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3591–3600, 2017. 2
- Raghavendra Selvan, Frederik Faye, Jon Middleton, and Akshay Pai. Uncertainty quantification in medical image segmentation with normalizing flows. In *Machine Learning in Medical Imaging*, Lecture Notes in Computer Science, Switzerland, 2020. Springer. 2
- Murat Sensoy, Lance Kaplan, and Melih Kandemir. Evidential deep learning to quantify classification uncertainty. *Advances in neural information processing systems*, 31, 2018. 1, 3
- Philipp Tschandl, Cliff Rosendahl, and Harald Kittler. The HAM10000 dataset, a large collection of multi-source dermatoscopic images of common pigmented skin lesions. *Scientific Data*, 5(1):1–9, 2018. 4
- Dennis Ulmer, Christian Hardmeier, and Jes Frellsen. Prior and posterior networks: A survey on evidential deep learning methods for uncertainty estimation. *arXiv preprint arXiv:2110.03051*, 2023. 1, 3
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 8