

# Bayesian uncertainty-weighted loss for improved generalisability on polyp segmentation task

Rebecca S. Stone(✉), Pedro E. Chavarrias-Solano, Andrew J. Bulpitt, David C. Hogg, Sharib Ali(✉)<sup>[0000-0003-1313-3542]</sup>

School of Computing, University of Leeds, LS2 9JT, Leeds, UK  
sc16rsmly, s.s.ali@leeds.ac.uk

**Abstract.** While several previous studies have devised methods for segmentation of polyps, most of these methods are not rigorously assessed on multi-center datasets. Variability due to appearance of polyps from one center to another, difference in endoscopic instrument grades, and acquisition quality result in methods with good performance on in-distribution test data, and poor performance on out-of-distribution or underrepresented samples. Unfair models have serious implications and pose a critical challenge to clinical applications. We adapt an implicit bias mitigation method which leverages Bayesian epistemic uncertainties during training to encourage the model to focus on underrepresented sample regions. We demonstrate the potential of this approach to improve generalisability without sacrificing state-of-the-art performance on a challenging multi-center polyp segmentation dataset (PolypGen) with different centers and image modalities.

## 1 Introduction

Colorectal cancer (CRC) is the third most common cancer worldwide [27] with early screening and removal of precancerous lesions (colorectal adenomas such as “polyps”) suggesting longer survival rates. While surgical removal of polyps (polypectomy) is a standard procedure during colonoscopy, detecting these and their precise delineation, especially for sessile serrated adenomas/polyps, is extremely challenging. Over a decade, advanced computer-aided methods have been developed and most recently machine learning (ML) methods have been widely developed by several groups. However, the translation of these technologies in clinical settings has still not been fully achieved. One of the main reasons is the generalisability issues with the ML methods [2]. Most techniques are built and adapted over carefully curated datasets which may not match the natural occurrences of the scene during colonoscopy.

Recent literature demonstrates how intelligent models can be systematically unfair and biased against certain subgroups of populations. In medical imaging, the problem is prevalent across various image modalities and target tasks; for example, models trained for lung disease prediction [25], retinal diagnosis [6], cardiac MR segmentation [23], and skin lesion detection [1,17] are all subject to biased performance against one or a combination of underrepresented gender,

age, socio-economic, and ethnic subgroups. Even under the assumption of an ideal sampling environment, a perfectly balanced dataset does not ensure unbiased performance as relative quantities are not solely responsible for bias [31,19]. This, and the scarcity of literature exploring bias mitigation for polyp segmentation in particular, strongly motivate the need for development and evaluation of mitigation methods which work on naturally occurring diverse colonoscopy datasets such as PolypGen [3].

## 2 Related work

Convolutional neural networks have recently worked favourably towards the advancement of building data-driven approaches to polyp segmentation using deep learning. These methods [18,34] are widely adapted from the encoder-decoder U-Net [24] architecture. Moreover, addressing the problem of different polyp sizes using multi-scale feature pruning methods, such as atrous-spatial pyramid pooling in DeepLabV3 [8] or high-resolution feature fusion networks like HR-Net [28] have been used by several groups for improved polyp segmentation. For example, MSRFNet [29] uses feature fusion networks between different resolution stages. Recent work on generalisability assessment found that methods trained on specific centers do not tend to generalise well on unseen center data or different naturally occurring modalities such as sequence colonoscopy data [2]. These performance gaps were reported to be large (drops of nearly 20%).

Out-of-distribution (OOD) generalisation and bias mitigation are challenging, open problems in the computer vision research community. While in the bias problem formulation, models wrongly correlate one or more spurious (non-core) features with the target task, the out-of-distribution problem states that test data is drawn from a separate distribution than the training data. Some degree of overlap between the two distributions in the latter formulation exists, which likely includes the core features. Regardless of the perspective, the two problems have clear similarities, and both result in unfair models which struggle to generalise for certain sub-populations. In the literature, many works focus on OOD detection, through normal or modified softmax outputs [13], sample uncertainty thresholds from Bayesian, ensemble, or other models [20,14,7], and distance measures in feature latent space [12]. Other approaches tackle the more difficult problem of algorithmic mitigation through disentangled representation learning, architectural and learning methods, and methods which optimise for OOD generalisability directly [26].

Similarly, several categories of bias mitigation methods exist. Some methods rely on two or more models, one encouraged to learn the biased correlations of the majority, and the other penalised for learning the correlations of the first [21,16]. Other approaches modify the objective loss functions to reward learning core rather than spurious features [33,22], or by neutralising representations to remove learned spurious correlations [10]. Others use data augmentation [6], or explore implicit versions of up-weighting or re-sampling underrepresented samples by discovering sparse areas of the feature space [4] or dynamically identifying

samples more likely to be underrepresented [30]. De-biasing methods leveraging Bayesian model uncertainties [15,5,30] provide the added benefits of uncertainty estimations which are useful in clinical application for model interpretability and building user confidence.

To tackle the generalisability problem for polyp segmentation, we consider the diversity of features in a multi-centre polyp dataset [3]. Our contributions can be listed as: 1) adapting an implicit bias mitigation strategy in [30] from a classification to a segmentation task; 2) evaluating the suitability of this approach on three separate test sets which have been shown to be challenging generalisation problems. Our experiments demonstrate that our method is comparable and in many cases even improves the performance compared to the baseline state-of-the-art segmentation method while decreasing performance discrepancies between different test splits.

### 3 Method

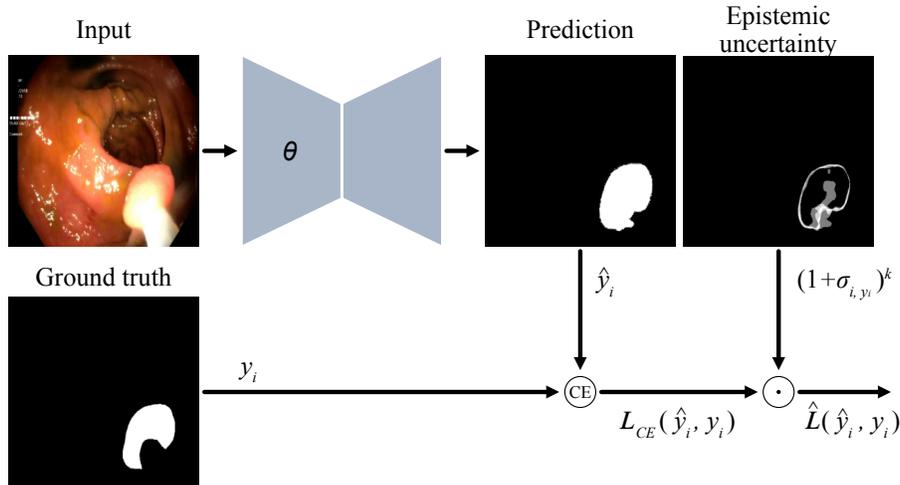
The encoder-decoder architecture for semantic segmentation has been widely explored in medical image analysis. In our approach we have used DeepLabV3 [9] as baseline model that has SOTA performance on the PolypGen dataset [3]. We then apply a probabilistic model assuming a Gaussian prior on all trainable weights (both encoder and decoder) that are updated to the posterior using the training dataset. For the Bayesian network with parameters  $\theta$ , posterior  $p(\theta \mid D)$ , training data with ground truth segmentation masks  $D = (X, Y)$ , and sample  $x_i$ , the predictive posterior distribution for a given ground truth segmentation mask  $y_i$  can be written as:

$$p(y_i \mid D, x_i) = \int p(y_i \mid \theta, x_i)p(\theta \mid D)d\theta \quad (1)$$

While Monte-Carlo dropout [11] at test-time is a popular approach to approximating this intractable integral, we choose stochastic gradient Monte-Carlo sampling MCMC (SG-MCMC [32]) for a better posterior. Stochastic gradient over mini-batches includes a noise term approximating the gradient over the whole training distribution. Furthermore, the cyclical learning rate schedule introduced in [35] known as cyclical SG-MCMC, or cSG-MCMC, allows for faster convergence and better exploration of the multimodal distributions prevalent in deep neural networks. Larger learning step phases provide a warm restart to the subsequent smaller steps in the sampling phases.

The final estimated posterior of the Bayesian network,  $\Theta = \{\theta_1, \dots, \theta_M\}$ , consists of  $M$  moments sampled from the posterior taken during the sampling phases of each learning cycle. With functional model  $\Phi$  representing the neural network, the approximate predictive mean  $\mu_i$  for one sample  $x_i$  is:

$$\mu_i \approx \frac{1}{M} \sum_{m=1}^M \Phi_{\theta_m}(x_i) \quad (2)$$



**Fig. 1:** Pixel-wise weighting of cross entropy (CE) loss contribution based on epistemic uncertainty maps for each training sample; the model is encouraged to focus on regions for which it is more uncertain.

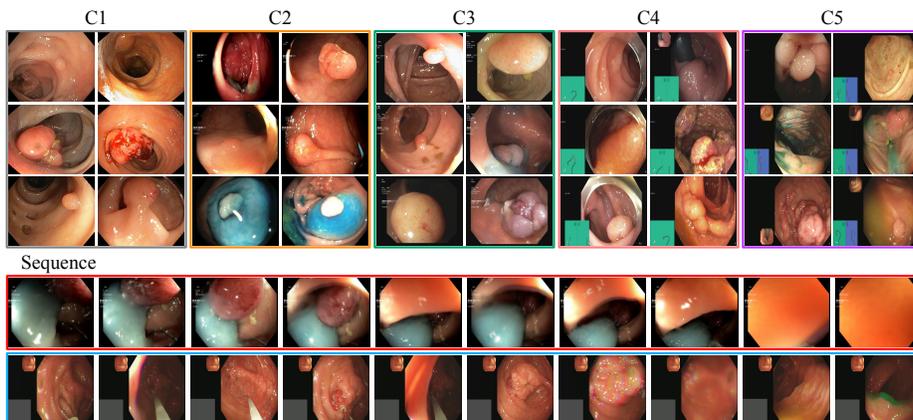
We can derive a segmentation prediction mask  $\hat{y}_i$  from  $\mu_i$  by taking the maximum output between the foreground and background channels. The epistemic uncertainty mask corresponding to this prediction (Equation 3) represents the *model uncertainty* for the predicted segmentation mask, the variance in the predictive distribution for that sample.

$$\sigma_i \approx \frac{1}{M} \sqrt{\sum_{m=1}^M (\Phi_{\theta_m}(x_i) - \mu_i)^2} \quad (3)$$

We add epistemic uncertainty-weighted sample loss [30] that identifies high-uncertainty sample regions during training. It also scales the pixel-wise contribution of these regions to the loss computation via a simple weighting function (Equation 4). This unreduced cross-entropy loss is then averaged over each image and batch (see Fig. 1).

$$\hat{L}(\hat{y}_i, y_i) = L_{CE}(\hat{y}_i, y_i) * (1.0 + \sigma_{i,y_i})^\kappa \quad (4)$$

The shift by a constant (1.0) normalises the values, ensuring that the lowest uncertainty samples are never irrelevant to the loss term.  $\kappa$  is a tunable debiasing parameter;  $\kappa = 1$  being a normal weighting, whereas  $\kappa \rightarrow \infty$  increases the importance of high-uncertainty regions. As too large a  $\kappa$  results in degraded performance due to overfitting, the optimal value is determined by validation metrics.



**Fig. 2:** Samples from the EndoCV2021 dataset; from (*top*) C1-5 single frames and (*bottom*) C1-5-SEQ; (*top*) highlights the data distribution of each center (C1-C5), which consists of curated frames with well-defined polyps; (*bottom*) demonstrates the variability of sequential data due to the presence of artifacts, occlusions, and polyps with different morphology.

## 4 Experiments and results

### 4.1 Dataset and experimental setup

PolypGen [3] is an expert-curated polyp segmentation dataset comprising of both single frames and sequence frames (frames sampled at every 10 frames from video) from over 300 unique patients across six different medical centers. The natural data collection format is video from which single frames and sequence data are hand-selected. The single frames are clearer, better quality, and with polyps in each frame including polyps of various sizes (10k to 40k pixels), and also potentially containing additional artifacts such as light reflections, blue dye, partial view of instruments, and anatomies such as colon linings and mucosa covered with stool, and air bubbles (Fig. 2). The sequence frames are more challenging and contain more negative samples without a polyp and more severe artifacts, which are a natural occurrence in colonoscopy. Our training set includes 1449 single frames from five centers (C1 to C5) and we evaluate on the three tests sets used for generalisability assessment in literature [2,3].

The first test dataset has 88 single frames from an unseen center C6 (C6-SIN), and the second has 432 frames from sequence data also from unseen center C6 (C6-SEQ). Here, the first test data (C6-SIN) comprises of hand selected images from the colonoscopy videos while the second test data (C6-SEQ) includes short sequences (every 10<sup>th</sup> frame of video) mimicking the natural occurrence of the procedure. The third test dataset includes 124 frames but from seen centers C1 - C5; however, these are more challenging as they contain both positive and negative samples with different levels of corruption that are not as present in

Dataset	Method	JAC	Dice	F2	PPV	Recall	Accuracy
C6-SIN	SOTA	0.738±0.3	0.806±0.3	0.795±0.3	<b>0.912±0.2</b>	0.793±0.3	0.979±0.1
	BayDeepLabV3+	0.721±0.3	0.790±0.3	<b>0.809±0.3</b>	0.836±0.2	<b>0.843±0.3</b>	<b>0.977±0.1</b>
	Ours	<b>0.740±0.3</b>	<b>0.810±0.3</b>	<u>0.804±0.3</u>	<u>0.903±0.1</u>	<u>0.806±0.3</u>	<b>0.977±0.1</b>
C1-5-SEQ	SOTA	<b>0.747±0.3</b>	<b>0.819±0.3</b>	<b>0.828±0.3</b>	<u>0.877±0.2</u>	<u>0.852±0.3</u>	0.960±0.0
	BayDeepLabV3+	0.708±0.3	0.778±0.3	0.805±0.3	0.784±0.3	<b>0.885±0.2</b>	<b>0.963±0.0</b>
	Ours	<u>0.741±0.3</u>	<u>0.810±0.3</u>	<u>0.815±0.3</u>	<b>0.888±0.2</b>	0.836±0.3	<u>0.961±0.0</u>
C6-SEQ	SOTA	0.608±0.4	0.676±0.4	0.653±0.4	<u>0.845±0.3</u>	0.719±0.3	0.964±0.1
	BayDeepLabV3+	<u>0.622±0.4</u>	<u>0.682±0.4</u>	<u>0.669±0.4</u>	0.802±0.3	<b>0.764±0.3</b>	<u>0.965±0.1</u>
	Ours	<b>0.637±0.4</b>	<b>0.697±0.4</b>	<b>0.682±0.4</b>	<b>0.858±0.3</b>	<u>0.741±0.3</u>	<b>0.967±0.1</b>

**Table 1:** Evaluation of the state-of-the-art deterministic DeepLabV3+, Bay-DeepLabV3+, and our proposed BayDeepLabV3+Unc, showing mean and standard deviations across the respective test dataset samples. **First** and **second** best results for each metric per dataset formatted.

the curated single frame training set. As no C6 samples nor sequence data are present in the training data, these test sets present a challenging generalisability problem.<sup>1</sup>

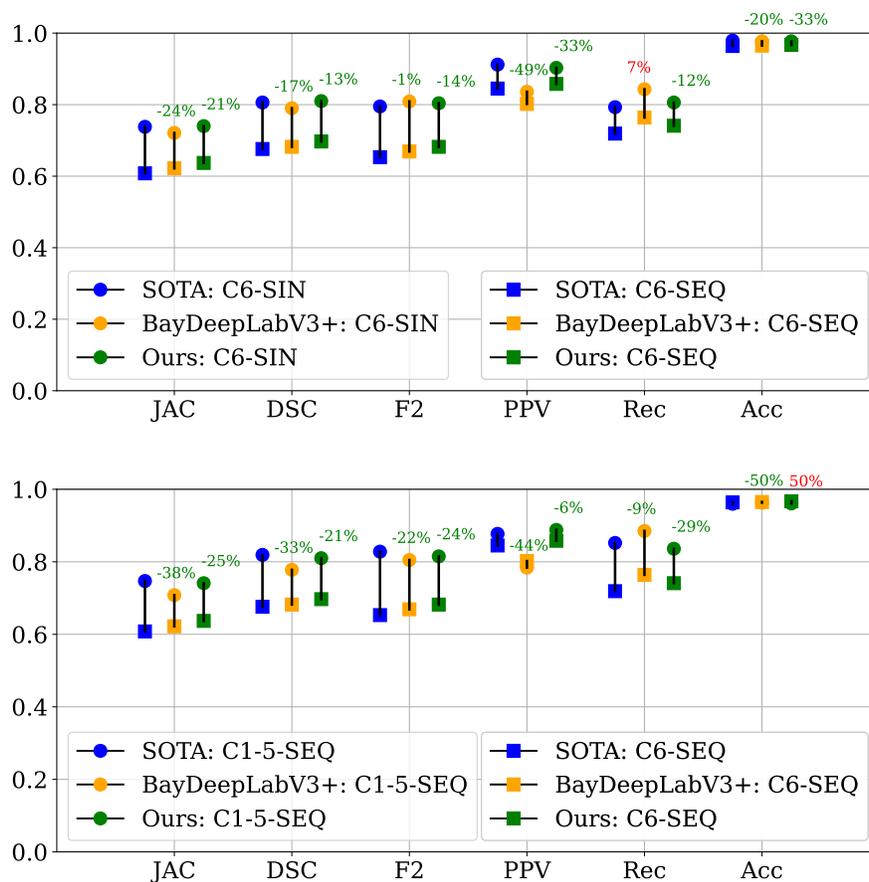
Training was carried out on several IBM Power 9 dual-CPU nodes with 4 NVIDIA V100 GPUs. Validation metrics were used to determine optimal models for all experiments with hyper-parameters chosen via grid search. Perhaps due to some frames containing very large polyps with high uncertainties, we found that the gradients of Bayesian models with uncertainty-weighted loss (BayDeepLabV3+Unc) occasionally exploded during the second learning cycle, and clipping the absolute gradients at 1.0 for all weights prevented this issue. All Bayesian DeepLabV3+ (BayDeepLabV3+) models had 2 cycles, a cycle length of 550 epochs, noise control parameter  $\alpha = 0.9$ , and an initial learning rate of 0.1. For BayDeepLabV3+Unc, we found optimal results with de-biasing tuning parameter  $\kappa = 3$ . Posterior estimates for BayDeepLabV3+ and Bay-DeepLabV3+Unc included 6 and 4 samples per cycle, respectively.

## 4.2 Results

We use the state-of-the-art deterministic model<sup>2</sup> and checkpoints to evaluate on the three test sets, and compare against the baseline Bayesian model Bay-DeepLabV3+ and BayDeepLabV3+Unc with uncertainty-weighted loss.

<sup>1</sup> C1-5-SEQ and C6-SEQ data are referred to as DATA3 and DATA4, respectively, in [2]

<sup>2</sup> <https://github.com/sharib-vision/PolypGen-Benchmark>



**Fig. 3:** Performance gaps of the three models (state-of-the-art deterministic DeepLabV3+, BayDeepLabV3+, and BayDeepLabV3+Unc) between the three different test sets; (*top*) comparing performance on single vs. sequence frames from out-of-distribution test set C6 (C6-SIN vs. C6-SEQ), and (*bottom*) sequence frames from C1 - C5 vs. unseen C6 (C1-5-SEQ vs. C6-SEQ). The subtext above bars indicates the percent decrease in performance gap compared to SOTA; a larger percent decrease and shorter vertical bar length indicate better generalisability.

We report results for Jaccard index (JAC), Dice coefficient (Dice),  $F_\beta$ -measure with  $\beta = 2$  (F2), positive predictive value (PPV), recall (Rec), and mean pixel-wise accuracy (Acc). PPV in particular has high clinical value as it indicates a more accurate delineation for the detected polyps. Recall and mean accuracy are less indicative since the majority of frames are background in the segmentation task and these metrics do not account for false positives. A larger number of

false positive predictions can cause inconvenience to endoscopists during colonoscopic procedure and hence can hinder clinical adoption of methods. Figure 3 illustrates that our approach maintains SOTA performance across most metrics and various test settings, even outperforming in some cases; simultaneously, the performance gaps between different test sets representing different challenging features (1) image modalities (single vs. sequence frames) and (2) source centers (C1 - C5 vs. C6) are significantly decreased. Simply turning the SOTA model Bayesian improves the model’s ability to generalise, yet comes with a sacrifice in performance across metrics and datasets. Our proposed uncertainty-weighted loss achieves better generalisability without sacrificing performance (also see Table 1). We note performance superiority to SOTA especially on C6-SEQ, approximately 3% improvement on Dice. We can also observe slight improvement on PPV for test sets with sequence (both held-out data and unseen centre data). Finally, we note that in clinical applications, the uncertainty maps for samples during inference could be useful for drawing clinicians’ attention towards potentially challenging cases, increasing the likelihood of a fairer outcome.

## 5 Conclusion

We have motivated the critical problem of model fairness in polyp segmentation on a multi-center dataset, and modified a Bayesian bias mitigation method to our task. The results on three challenging test sets show strong potential for improving generalisability while maintaining competitive performance across all metrics. Furthermore, the proposed mitigation method is implicit, not requiring comprehensive knowledge of biases or out-of-distribution features in the training data. This is of particular importance in the medical community given the sensitivity and privacy issues limiting collection of annotations and metadata. Our findings are highly relevant to the understudied problem of generalisation across high variability colonoscopy images, and we anticipate future work will include comparisons with other methods to improve generalisability and an extension to the approach. We also anticipate having access to additional test data for more in-depth analysis of the results.

**Acknowledgements** R. S. Stone is supported by an Ezra Rabin scholarship.

## References

1. Abbasi-Sureshjani, S., Raumanns, R., Michels, B.E., Schouten, G., Cheplygina, V.: Risk of training diagnostic algorithms on data with demographic bias. In: Interpretable and Annotation-Efficient Learning for Medical Image Computing: Third International Workshop, iMIMIC 2020, Second International Workshop, MIL3ID 2020, and 5th International Workshop, LABELS 2020. pp. 183–192. Springer (2020)
2. Ali, S., Ghatwary, N., Jha, D., Isik-Polat, E., Polat, G., Yang, o.: Assessing generalisability of deep learning-based polyp detection and segmentation methods through a computer vision challenge. arXiv preprint arXiv:2202.12031 (2022)

3. Ali, S., Jha, D., Ghatwary, N., Realdon, S., Cannizzaro, R., Salem, O.E., Lamarque, D., Daul, C., Riegler, M.A., Anonsen, K.V., et al.: A multi-centre polyp detection and segmentation dataset for generalisability assessment. *Scientific Data* **10**(1), 75 (2023)
4. Amini, A., Soleimany, A.P., Schwarting, W., Bhatia, S.N., Rus, D.: Uncovering and mitigating algorithmic bias through learned latent structure. In: *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*. pp. 289–295 (2019)
5. Branchaud-Charron, F., Atighehchian, P., Rodríguez, P., Abuhamad, G., Lacoste, A.: Can active learning preemptively mitigate fairness issues? *arXiv preprint arXiv:2104.06879* (2021)
6. Burlina, P., Joshi, N., Paul, W., Pacheco, K.D., Bressler, N.M.: Addressing artificial intelligence bias in retinal diagnostics. *Translational Vision Science & Technology* **10**(2), 13–13 (2021)
7. Cao, S., Zhang, Z.: Deep hybrid models for out-of-distribution detection. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 4733–4743 (2022)
8. Chen, L.C., Papandreou, G., Kokkinos, I., Murphy, K., Yuille, A.L.: Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **40**(4), 834–848 (2018)
9. Chen, L.C., Zhu, Y., Papandreou, G., Schroff, F., Adam, H.: Encoder-decoder with atrous separable convolution for semantic image segmentation. In: *Proceedings of the European conference on computer vision (ECCV)*. pp. 801–818 (2018)
10. Du, M., Mukherjee, S., Wang, G., Tang, R., Awadallah, A., Hu, X.: Fairness via representation neutralization. *Advances in Neural Information Processing Systems* **34**, 12091–12103 (2021)
11. 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)
12. Gonzalez, C., Gotkowski, K., Bucher, A., Fischbach, R., Kaltenborn, I., Mukhopadhyay, A.: Detecting when pre-trained nnu-net models fail silently for covid-19 lung lesion segmentation. In: *Medical Image Computing and Computer Assisted Intervention–MICCAI 2021*. pp. 304–314. Springer (2021)
13. Hendrycks, D., Mazeika, M., Kadavath, S., Song, D.: Using self-supervised learning can improve model robustness and uncertainty. *Advances in neural information processing systems* **32** (2019)
14. Jungo, A., Balsiger, F., Reyes, M.: Analyzing the quality and challenges of uncertainty estimations for brain tumor segmentation. *Frontiers in neuroscience* **14**, 282 (2020)
15. Khan, S., Hayat, M., Zamir, S.W., Shen, J., Shao, L.: Striking the right balance with uncertainty. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 103–112 (2019)
16. Kim, N., Hwang, S., Ahn, S., Park, J., Kwak, S.: Learning debiased classifier with biased committee. *arXiv preprint arXiv:2206.10843* (2022)
17. Li, X., Cui, Z., Wu, Y., Gu, L., Harada, T.: Estimating and improving fairness with adversarial learning. *arXiv preprint arXiv:2103.04243* (2021)
18. Mahmud, T., Paul, B., Fattah, S.A.: Polypsegnet: A modified encoder-decoder architecture for automated polyp segmentation from colonoscopy images. *Computers in Biology and Medicine* **128**, 104119 (2021)

19. Mehrabi, N., Morstatter, F., Saxena, N., Lerman, K., Galstyan, A.: A survey on bias and fairness in machine learning. *ACM Computing Surveys (CSUR)* **54**(6), 1–35 (2021)
20. 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)
21. Nam, J., Cha, H., Ahn, S., Lee, J., Shin, J.: Learning from failure: De-biasing classifier from biased classifier. *Advances in Neural Information Processing Systems* **33**, 20673–20684 (2020)
22. Pezeshki, M., Kaba, O., Bengio, Y., Courville, A.C., Precup, D., Lajoie, G.: Gradient starvation: A learning proclivity in neural networks. *Advances in Neural Information Processing Systems* **34**, 1256–1272 (2021)
23. Puyol-Antón, E., Ruijsink, B., Piechnik, S.K., Neubauer, S., Petersen, S.E., Razavi, R., King, A.P.: Fairness in cardiac mr image analysis: an investigation of bias due to data imbalance in deep learning based segmentation. In: *Medical Image Computing and Computer Assisted Intervention—MICCAI 2021*. pp. 413–423. Springer (2021)
24. Ronneberger, O., Fischer, P., Brox, T.: U-Net: Convolutional Networks for Biomedical Image Segmentation, vol. 9351, p. 234–241. Springer International Publishing (2015)
25. Seyyed-Kalantari, L., Liu, G., McDermott, M., Chen, I.Y., Ghassemi, M.: Chexclusion: Fairness gaps in deep chest x-ray classifiers. In: *BIOCOMPUTING 2021: Proceedings of the Pacific Symposium*. pp. 232–243. World Scientific (2020)
26. Shen, Z., Liu, J., He, Y., Zhang, X., Xu, R., Yu, H., Cui, P.: Towards out-of-distribution generalization: A survey. *arXiv preprint arXiv:2108.13624* (2021)
27. Silva, J., Histace, A., Romain, O., Dray, X., Granado, B.: Toward embedded detection of polyps in wce images for early diagnosis of colorectal cancer. *International journal of computer assisted radiology and surgery* **9**, 283–293 (2014)
28. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition (2015)
29. Srivastava, A., Jha, D., Chanda, S., Pal, U., Johansen, H.D., Johansen, D., Riegler, M.A., Ali, S., Halvorsen, P.: MSRF-Net: a multi-scale residual fusion network for biomedical image segmentation. *IEEE Journal of Biomedical and Health Informatics* **26**(5), 2252–2263 (2022)
30. Stone, R.S., Ravikumar, N., Bulpitt, A.J., Hogg, D.C.: Epistemic uncertainty-weighted loss for visual bias mitigation. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 2898–2905 (2022)
31. Wang, M., Deng, W.: Mitigating bias in face recognition using skewness-aware reinforcement learning. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. pp. 9322–9331 (2020)
32. Welling, M., Teh, Y.W.: Bayesian learning via stochastic gradient langevin dynamics. In: *Proceedings of the 28th international conference on machine learning (ICML-11)*. pp. 681–688. Citeseer (2011)
33. Xu, X., Huang, Y., Shen, P., Li, S., Li, J., Huang, F., Li, Y., Cui, Z.: Consistent instance false positive improves fairness in face recognition. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 578–586 (2021)
34. Yeung, M., Sala, E., Schönlieb, C.B., Rundo, L.: Focus u-net: A novel dual attention-gated cnn for polyp segmentation during colonoscopy. *Computers in Biology and Medicine* **137**, 104815 (2021)
35. Zhang, R., Li, C., Zhang, J., Chen, C., Wilson, A.G.: Cyclical stochastic gradient mcmc for bayesian deep learning. *arXiv preprint arXiv:1902.03932* (2019)