# Analyzing Overfitting under Class Imbalance in Neural Networks for Image Segmentation

Zeju Li, Konstantinos Kamnitsas and Ben Glocker

*Abstract*—**Class imbalance poses a challenge for developing unbiased, accurate predictive models. In particular, in image segmentation neural networks may overfit to the foreground samples from small structures, which are often heavily under-represented in the training set, leading to poor generalization. In this study, we provide new insights on the problem of overfitting under class imbalance by inspecting the network behavior. We find empirically that when training with limited data and strong class imbalance, at test time the distribution of logit activations may shift across the decision boundary, while samples of the well-represented class seem unaffected. This bias leads to a systematic under-segmentation of small structures. This phenomenon is consistently observed for different databases, tasks and network architectures. To tackle this problem, we introduce new asymmetric variants of popular loss functions and regularization techniques including a large margin loss, focal loss, adversarial training, mixup and data augmentation, which are explicitly designed to counter logit shift of the under-represented classes. Extensive experiments are conducted on several challenging segmentation tasks. Our results demonstrate that the proposed modifications to the objective function can lead to significantly improved segmentation accuracy compared to baselines and alternative approaches.**

*Index Terms*—**overfitting, class imbalance, image segmentation.**

## I. INTRODUCTION

**T**HE success of convolutional neural networks (CNNs) is strongly linked with the availability of large scale, representative datasets. However, in many real-world applications such as medical image segmentation, the availability of large, annotated datasets is still limited. But even when there is a sufficient number of images available, the fundamental problem of class imbalance remains where region-of-interests (ROIs) (i.e. foreground classes) are heavily under-represented in the training data [3], [30]. Similar to [9], the class imbalance ratio of one image can be defined as the ratio between the number of pixels of the background class (which is commonly the most frequent class) and the number of pixels of different object classes. The class imbalance ratio of a whole dataset would then be reported as the average class imbalance ratio of all images in the set. Class imbalance ratios of 100:1 or higher are not uncommon in applications such as lesion segmentation, as shown in Table I.

When the model is trained with imbalanced datasets, it can overfit to the training samples from the under-represented classes and may not generalize well during test time. However, the effects of overfitting under class imbalance on the model behavior is not well understood. In this study, we investigate how the distribution of activations of the classification layer (*logits*) changes when the model is trained using different amounts of training data with strong class imbalance. As the model is trained with fewer training data and overfit the under-represented classes more, we find that the model projects unseen samples of the under-represented classes closer to and even across the decision boundary, while samples of the over-represented classes remain unaffected. This biased distribution shift leads to under-segmentation of under-represented class. Current solutions to address class imbalance or to mitigate overfitting do not explicitly consider this asymmetric logit shift and are unable to lead to significant improvements, as we show through an extensive set of experiments.

This study sheds new light on the problem of overfitting in the presence of class imbalance by making the following key contributions: 1) Via inspection of the network behavior on four segmentation tasks and datasets, and two popular model architectures, we conclude that overfitting under class imbalance consistently leads to decreased performance on under-represented classes specifically in terms of low sensitivity; 2) We identify the shift in the logit distribution of unseen test samples of under-represented classes as a result of overfitting under class imbalance; 3) Base on our observations, we propose simple yet effective asymmetric variants of five loss functions and regularization techniques which are explicitly designed to change the network behavior yielding improved segmentation accuracy for the under-represented classes.

This article is an extension of our earlier work presented at MICCAI 2019 [27]. We extend our previous work on multiple aspects: 1) We provide a more detailed analysis including experiments on two additional datasets; 2) We further include a 3D U-Net to confirm that our observations hold across different network architectures; 3) We explore the proposed training objectives with the Dice loss in addition to cross-entropy; 4) We enrich the experiments by adding comparisons when training with F-score, extend experiments to multi-class segmentation, and also evaluate another regularization method which is noted as asymmetric augmentation. Our findings here confirm our initial observations about the biased behavior of neural networks. The behavior of logit distribution shift is consistently observed across different types of data, tasks and architectures. Our work highlights the importance of the issue of overfitting under class imbalance. The quantitative evaluation further supports our proposal of taking class imbalance into account when designing the learning objective.

Z. Li, K. Kamnitsas and B. Glocker are with the BioMedIA Group, Department of Computing, Imperial College London, SW72AZ, United Kingdom. E-mail: zeju.li18@imperial.ac.uk.

## II. RELATED WORK

### A. Class imbalance

Class imbalance, which has been the focus of previous works [4], [19], is a common issue in image classification and image segmentation. Compared with the literature on class imbalance, a key contribution of this study is the focus on the model behaviour when it overfits to the under-represented classes with a detailed analysis and potential solutions. In the following, we discuss related work categorized by different methodological approaches.

*1) Re-weighting:* A common approach to tackle class imbalance is class-level re-weighting, which assigns higher weights or higher sampling probability to the under-represented classes based on sample frequency [41], [46] or advanced rules [9]. In this study, explore re-weighting as a baseline approach in all experiments where we train the models with patches which are separately sampled from different classes with the same probability. Beyond that, sample-level re-weighting strategies are also proposed to build a balanced model. For example, hard sample mining was proposed to avoid the dominant effect of majority classes [11]. Similarly, focal loss and its variants were proposed to weight difficult samples over easy samples [1], [14], [29], [43] to steer the learning towards small objects. However, the under-represented samples are not necessarily difficult to predict during training. In fact, as we show empirically, the training samples of the under-represented classes are learned well due to overfitting. In this case, we find that a focal loss may even decrease the performance when processing imbalanced datasets because it reduces the focus on the under-represented samples. Therefore, in this study, we improve upon focal loss by removing the attenuation of under-represented classes. Margin based loss functions were proposed to learn discriminative embeddings and widely adopted to metric-learning and face recognition [10], [31]. Margin losses can also been seen as a kind of re-weighting approach which changes the magnitude of the gradient of the network output by multiplying a scalar, as we show in the supplementary Section VII. In this study, we propose to only assign margins for the under-represented classes. The design of uneven margins for imbalanced datasets was first proposed in [26] for perception. Recently, Large Margin Local Embedding (LMLE) was proposed to put more constraints for the under-represented classes by only applying multiple margins to the minority classes, with a computationally expensive metric-learning based framework [17]. More recently, two concurrent studies were also proposed to set larger margins for the under-represented classes from the perspective of uncertainty [23] or generalization bound [5], [49]. In this study, we empirically show that one should not assign margins for the over-represented classes based on the observations of asymmetric logit distribution under class imbalance.

*2) Data synthesis:* Our work is related to data synthesis methods [6], [13] which generate synthetic samples of the minority classes based on intra-class relationship between samples to increase the variance of under-represented class. In addition, we create synthetic samples in the latent feature space rather than image space and provide two new ways to synthesize samples of the under-represented classes for modern machine learning models. We also propose to adopt stronger data augmentation for the under-represented classes by changing the augmentation probabilities to alleviate overfitting.

*3) Other methods:* The above mentioned methods are all based on changing the training data distribution to tackle class imbalance. In contrast, some other approaches try to counter class imbalance by modifying the training strategy. Specifically, [15] firstly trained their model with data which is sampled from each class with the same probability. Thereafter, they only retrain the output layer with uniformly sampled data while freezing all other network parameters. In this way, they could separately learn a diverse representation and a classifier for realistic data distribution. Similar strategies, which aim to change the decision boundary at test-time, are also proposed recently for long-tailed recognition [21], [38], [48]. These approaches are complementary to our proposed solutions and could be combined. Other learning paradigms such as meta-learning [42] and transfer learning [32] were also recently proposed for long-tail learning, but these are outside the scope of this paper.

*4) Segmentation:* The problem of class imbalance in image segmentation is different from that in image recognition [32] because the dominating class in image segmentation is the background class with diverse characteristics, and its segmentation accuracy is highly robust. In contrast, the accuracy for the majority classes in long-tailed image recognition can degrade with common techniques such as re-weighting [21]. In addition, the evaluation of segmentation performance mostly relies on the foreground classes, and therefore, the focus is on improving accuracy in those classes. For example, recent studies proposed to provide a better trade-off between sensitivity and precision for segmentation [14], [33]. However, these strategies yield little improvements when processing highly imbalanced datasets, as shown in our experiments. This is because a deep neural network may achieve near-perfect training accuracy even for the under-represented samples without benefitting from the modified loss function. Class imbalance in image segmentation has been also approached via a boundary loss [22]. However, it is only applicable to segmentation. The authors show promising results with sufficient training data, but the model may still be prone to overfit the under-represented class with limited dataset. Other work adopted multi-stage approaches with candidate proposals and background suppression [36], [39]. However, the candidate prediction process may still suffer from class imbalance. In addition, any missed candidates in one stage cannot be recovered in a later stage. In contrast, the solutions proposed here are general loss functions that can be incorporated into any model or learning approach and is applicable beyond image segmentation.

### B. Regularization techniques

To improve generalization of deep neural networks, a number of regularization techniques are available. This includes

dropout [37], weight decay [24], data augmentation [7], [8], data mixing [45], [47], and adversarial training [12], [44]. However, most of these techniques were proposed for general image classification tasks where class imbalance is not explicitly addressed. It is also unclear how these techniques affect the network behavior in this setting.

## III. OVERFITTING UNDER CLASS IMBALANCE AND ITS EFFECT ON SEGMENTATION PERFORMANCE

To explore the effects of overfitting on the network behavior, we train CNNs using different amounts of data, on segmentation tasks that exhibit strong class imbalance. We conduct experiments on challenging segmentation tasks using data from the Multimodal Brain Tumor Image Segmentation (BRATS) challenge [2], the Anatomical Tracings of Lesions After Stroke (ATLAS) dataset [28], small organ segmentation (data from [25]) and Kidney Tumor Segmentation (KiTS) [16]. The statistics of those four datasets are summarized in Table I. To ensure our findings generalize across models, in our investigation we employ two convolutional network architectures that have been proven potent in a variety of segmentation tasks: We employ a DeepMedic architecture [20] for the experiments on brain lesions and multi-organ segmentation tasks on which it has previously shown high performance [20], [35], and a well configured 3D U-Net [18] for the experiments on kidney tumor segmentation on KiTS19 data, which is the base model of the winning entry of KiTS19 challenge [16]. The detailed network configurations are summarized in Section V.

TABLE I

THE STATISTICS AND CLASS IMBALANCE RATIOS OF THE FOUR DATASETS USED IN THIS STUDY. CLASS IMBALANCE RATIO IS DEFINED AS THE AVERAGE RATIO BETWEEN THE NUMBER OF THE BACKGROUND (BG) PIXELS AND THE FOREGROUND (FG) PIXELS OVER ALL IMAGES.

| Dataset | Total FG pixels | Total BG pixels | Class imbalance ratio (avg. $\pm$ std.) |
|---|---|---|---|
| BRATS | $1.2 \times 10^7$ | $253.2 \times 10^7$ | $712.8 \pm 1463.1$ |
| ATLAS | $4.6 \times 10^6$ | $1908.6 \times 10^6$ | $1768.7 \pm 6710.8$ |
| KiTS-Kidney | $13.5 \times 10^6$ | $1662.0 \times 10^6$ | $122.8 \pm 69.3$ |
| KiTS-Tumor | $2.9 \times 10^6$ | $1662.0 \times 10^6$ | $6736.6 \pm 1522.9$ |
| Abdomen-Gallbladder | $1.0 \times 10^5$ | $2630.5 \times 10^5$ | $4887.5 \pm 4854.7$ |
| Abdomen-Aorta | $3.8 \times 10^5$ | $2630.5 \times 10^5$ | $829.2 \pm 686.1$ |
| Abdomen-Vena cava | $3.4 \times 10^5$ | $2630.5 \times 10^5$ | $832.2 \pm 278.2$ |
| Abdomen-Vein | $1.4 \times 10^5$ | $2630.5 \times 10^5$ | $2196.5 \pm 792.2$ |

The observations on the test and training set are summarized in Fig. 1. With less training data, we notice a clear decrease of segmentation accuracy on test data while the accuracy on training data increases due to easier overfitting, as expressed by DSC (defined as $DSC = 2\frac{\text{sensitivity} \cdot \text{precision}}{\text{sensitivity} + \text{precision}}$). We observe that overfitting leads to a reduction of sensitivity while precision remains largely stable. In all settings and tasks, the specificity of the foreground always remains near-perfect ($>0.999$) with different amounts of training data, indicating the predictions of background samples are stable (not shown in Fig. 1 to avoid cluttering). We also provide some corresponding visualization examples in Fig 2.

Our observations on four different datasets show that this behavior is consistent when foreground classes are under-represented and is not specific to particular tasks and model

architectures. Our findings reveal that models that overfit to imbalanced training data have a bias to under-segment the under-represented class on unseen test data.

### A. Logit distribution shift

To obtain a better understanding of the network behavior after training on imbalanced data, we monitor the logit distribution when processing training and unseen test samples. The observations we make for the tasks of brain tumor core, kidney tumor and brain stroke lesion segmentation are summarized in Fig. 3, 4 and 5, respectively. We notice that the logit distribution of foreground samples shifts significantly towards and even across the decision boundary, while the logit distribution of background samples remains stable. The shift of the foreground logits results in a higher number of false negatives, which causes a drastic decrease of sensitivity (calculated as $\frac{\text{TP}}{\text{TP}+\text{FN}}$). This biased logit shift under class imbalance may also occur in other tasks such as image classification. However, it is particularly prevalent in image segmentation with small structures-of-interest.

We find that this shift of logits correlates with how much a model overfits to the under-represented class. Training with less data leads to more overfitting, and the logit distribution shift becomes larger. Moreover, we find that the logit shift also correlates with the size of structures represented by the foreground class. The rarest class shifts the most, as shown in the right part of Fig. 4.

In image segmentation, a CNN is optimized to push the logits of different classes away from each other and far from the decision boundary. It is relatively easy for a deep CNN to build an embedding for the training samples from the under-represented class because it just needs to build a set of case-specific filters to facilitate memorization. For example, as a CNN will only observe very few training samples of the foreground class, a CNN can dedicate specific model parameters to memorize all foreground samples, even if the individual patterns are rather complex. Specifically, we find a CNN seems to be more confident about foreground samples during training, mapping them farther away from the decision boundary when overfitting, as shown in Fig. 3, Fig. 4 and Fig. 5. However, these tailored filters will not generalize to unseen test data. Therefore, the activations for test samples of the under-represented class are smaller in magnitude (sub-optimal pattern matching of filters and unseen samples), leading to the observed distribution shift. In contrast, a CNN has to build generic filters for a well represented class to represent many different characteristics of the same class, leading to good generalization. Such filters will map unseen samples to similar locations in logit space and no shift between the embeddings of training and test samples is observed. As a result of class imbalance *and* overfitting, a CNN may underperform on the under-represented class while still generalizing well for the well-represented class.

While the negative effect of class imbalance and overfitting on model performance is well known, to our knowledge there has been little work investigating the specifics how the network behavior is affected. Only by understanding better
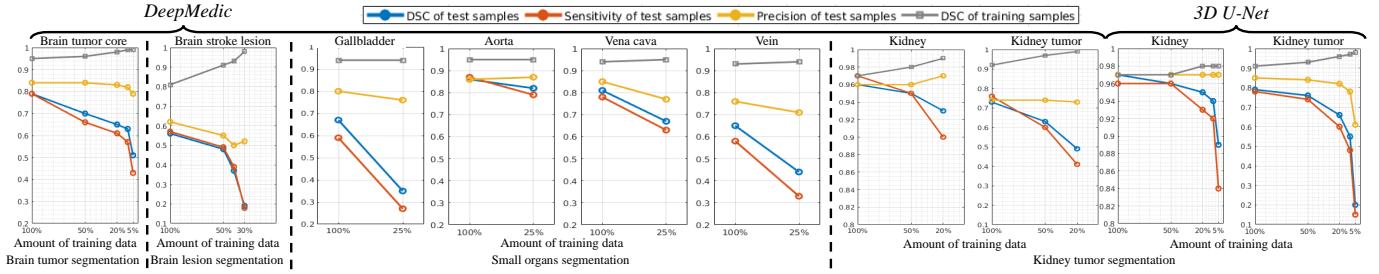
Fig. 1. Performance on brain tumor core, brain stroke lesion, small organs, kidney and kidney tumor segmentation with varying amounts of training data. The foreground (FG) and background (BG) samples are highly imbalanced, as noted below each subfigure. With less training data, performance drops due to the decrease of sensitivity, while the precision is largely retained.
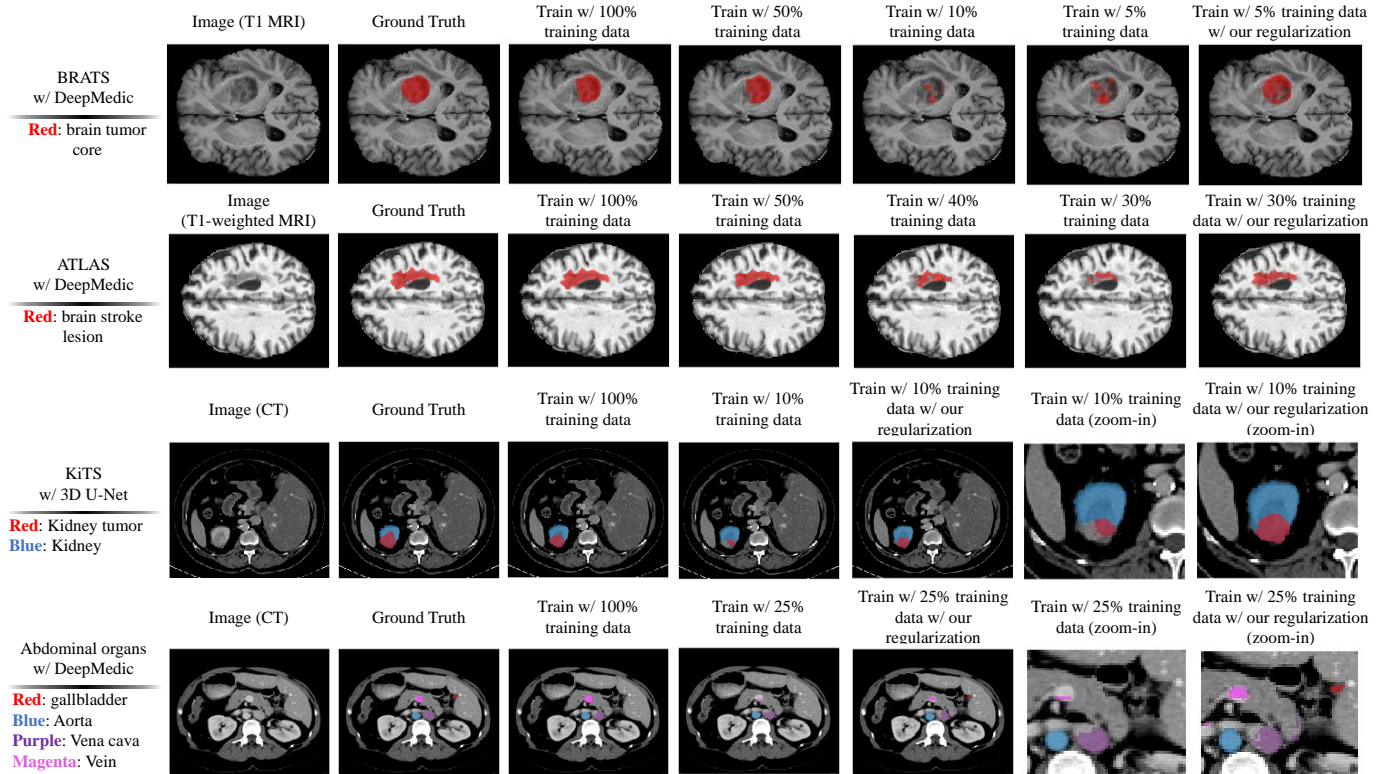


Fig. 2. Visualization of different datasets and segmentation results with different portions of training data. With less training data, the models are prone to under-segment the under-represented classes. The proposed regulation methods can alleviate the overfitting of under-represented classes and provide segmentation results with higher sensitivity and overall accuracy. Best viewed in color.

the implication, we can devise mitigation strategies. Previous loss functions and regularization techniques that aim to prevent overfitting did not take the behavior that we observe into account, and thus show limited success for improving segmentation accuracy in the setting of limited data with strong class imbalance. Here, we propose solutions via new asymmetric variants of existing objective functions leading to better feature embeddings for the under-represented samples, leading to significant improvements in segmentation accuracy for small structures-of-interest.

## IV. TACKLING OVERFITTING UNDER CLASS IMBALANCE WITH ASYMMETRIC OBJECTIVE FUNCTIONS

Based on our observations above about the biased behavior of CNNs, we design modifications to existing loss functions

and training strategies to prevent the logit distribution shift. Specifically, we add a bias for the under-represented class. Although the original techniques were proposed for different purposes, our modifications share a common goal: keep the logit activations of the under-represented class away from the decision boundary. Even if the logit of a foreground sample shifts towards the decision boundary as long as it does not cross it, its prediction remains correct (cf. Fig. 6).

### A. Asymmetric large margin loss

We consider a CNN for the task of semantic segmentation. For a training dataset $\{(\boldsymbol{x}_i, \boldsymbol{y}_i)\}_{i=1}^N$ with $N$ samples, we denote a training sample with $\boldsymbol{x}_i$ and its corresponding one-hot vector $\boldsymbol{y}_i$. If $c$ is the total number of classes of the task, $\boldsymbol{y}_i$ has $c$ elements, with its $j$'th element $y_{ij} \in \{0, 1\}$ corresponding
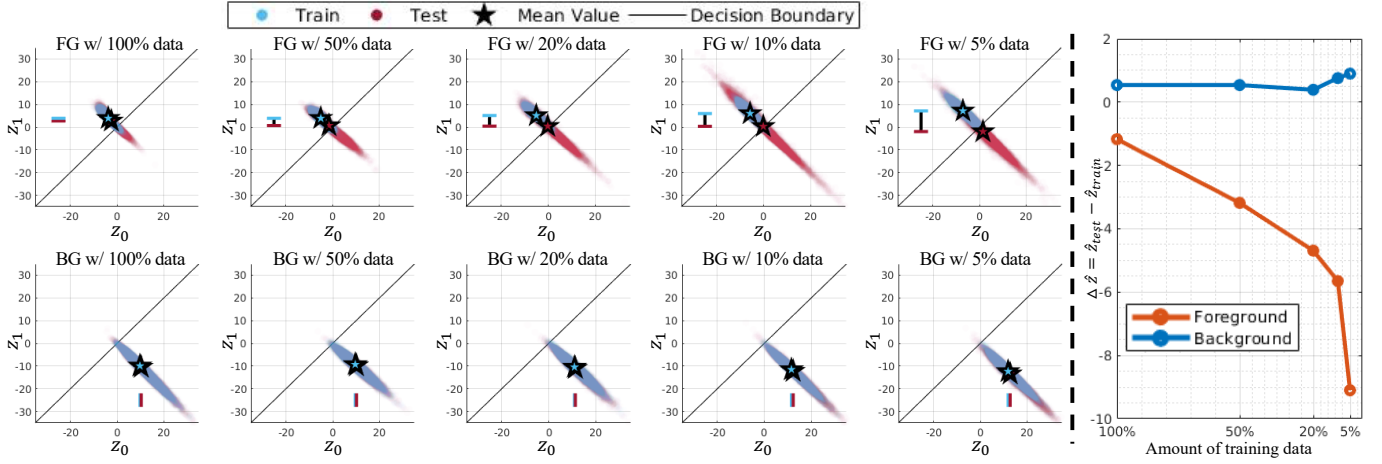
Fig. 3. (Left part) Activations of the classification layer (logit $z_0$ for background, logit $z_1$ for brain tumor core) when processing (top) tumor and (bottom) background samples of BRATS with DeepMedic, using different amounts of training data. The CNN maps training and testing samples of the background class to similar logit values. However, mean activation for testing data shifts significantly for the tumor class towards and sometimes across the decision boundary. (Right part) The shift of mean value of logits observed when processing training and testing data ($\Delta \hat{z} = |\hat{z}_{test}| - |\hat{z}_{train}|$).
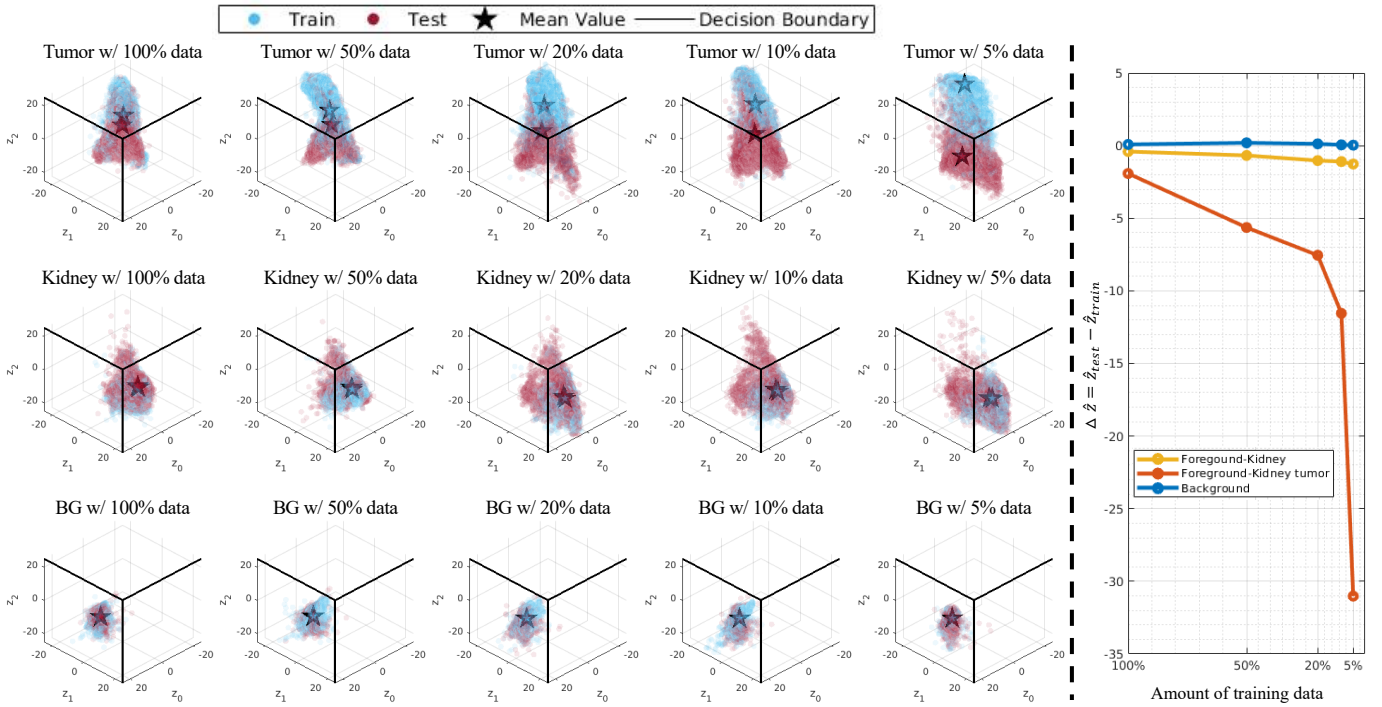


Fig. 4. (Left part) Activations of the classification layer (logit $z_0$ for background, logit $z_1$ for kidney, logit $z_2$ for kidney tumor) when processing (top) tumor, (middle) kidney and (bottom) background samples of KiTS with 3D U-Net, using different amounts of training data. The CNN also fails to map the training and testing samples of the tumor class in a similar position. (Right part) The shift of mean value of logits.

to the $j$'th class. $y_{ij}$ equals to 1 if $j$ is the real class of $\boldsymbol{x_i}$, or 0 otherwise. With this notation, the cross-entropy (CE) loss can be written as[1]:

$$L_{CE}(\boldsymbol{x_i}, \boldsymbol{y_i}) = -\sum_{j=1}^{c} y_{ij} \log(p_{ij}), \quad (1)$$

where $p_{ij}$ is the predicted probability by the network that the real class of $\boldsymbol{x_i}$ is $j$. Probability $p_{ij}$ is commonly obtained via

a softmax function over the $c$ activations $\{(z_{ij})_{j=1}^{c} \in \mathbb{R}^c\}$ that the network outputs for $\boldsymbol{x_i}$ at its last layer. These activations are called the *logits*. With this, $p_{ij}$ is given by:

$$p_{ij} = \frac{e^{z_{ij}}}{\sum_{j=1}^{c} e^{z_{ij}}}. \quad (2)$$

Besides CE, the smooth version of the DSC metric is an alternative choice for the loss function which is widely used for medical image segmentation [33]. DSC loss can be calculated in the form of $1 - \text{DSC} = \frac{\text{FP+FN}}{2\,\text{TP+FP+FN}}$), which is:

[1]We formulate CE as sum over classes to make class specific modifications.
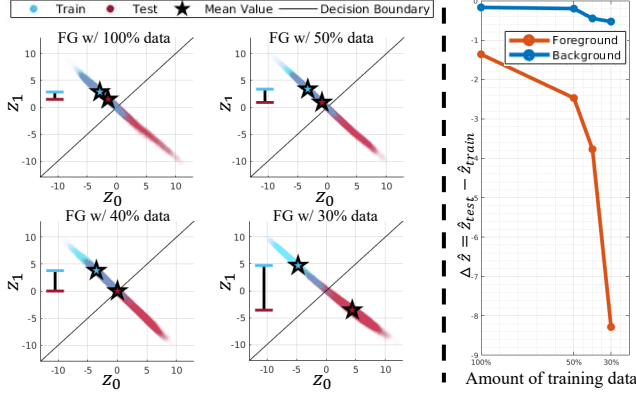
Fig. 5. (Left part) Activations of the classification layer when processing lesion samples of ATLAS with DeepMedic, using different amounts of training data. (Right part) The shift of mean value of logits.

$$L_{DSC}(\boldsymbol{x_i}, \boldsymbol{y_i}) = \sum_{j=1}^{c} \Big( \frac{(1 - y_{ij})p_{ij} + y_{ij}(1 - p_{ij})}{2y_{ij}p_{ij} + (1 - y_{ij})p_{ij} + y_{ij}(1 - p_{ij})} \Big). \quad (3)$$

The large margin loss was proposed for increasing the Euclidean distances between logits for different classes to learn discriminative features [40]. Symmetrically, it is implemented by adding a margin on the logits of every class:

$$L_{CE_M}(\boldsymbol{x_i}, \boldsymbol{y_i}) = -\sum_{j=1}^{c} y_{ij} \log(q_{ij}), \quad (4)$$

in which we require:

$$q_{ij} = \frac{\mathrm{e}^{z_{ij} - y_{ij}m}}{\sum_{j=1}^{c} \mathrm{e}^{z_{ij} - y_{ij}m}}, \quad (5)$$

where $m$ is a hyper-parameter for the margin. Although the large margin loss encourages the model to map different classes away from each other, the decision boundary remains in the center. According to our observations, class imbalance causes shifts of unseen foreground samples towards the background class. To mitigate this, a regularizer may aim to move the decision boundary closer to the background class. Our asymmetric modification only sets the margin for the rare classes. We define $\boldsymbol{r}$ as a one-hot vector with $c$ elements, with its $j$'th element $r_j \in \{0, 1\}$ corresponding to the $j$'th class and $r_j$ equals to 1 if $j$ is taken as the rare class. With the indication of $\boldsymbol{r}$, we derive the asymmetric large margin loss as:

$$\hat{L}_{CE_M}(\boldsymbol{x_i}, \boldsymbol{y_i}) = -\sum_{j=1}^{c} y_{ij} \log(\hat{q}_{ij}), \quad (6)$$

where we require:

$$\hat{q}_{ij} = \frac{\mathrm{e}^{z_{ij} - y_{ij}r_j m}}{\sum_{j=1}^{c} \mathrm{e}^{z_{ij} - y_{ij}r_j m}}, \quad (7)$$

In this study, we define $r_j$ as 1 for the foreground samples and 0 for the background samples. In other applications, $r_j$ can also be defined as a continuous variable indicating the rarity

of the classes with $r_j \in [0, 1]$, for methods in Section IV-A, IV-B and IV-C. Similarly, the symmetric and asymmetric large margin loss for DSC loss can be derived by substituting equation 5:

$$L_{DSC_M}(\boldsymbol{x_i}, \boldsymbol{y_i}) = \sum_{j=1}^{c} \Big( \frac{(1 - y_{ij})q_{ij} + y_{ij}(1 - q_{ij})}{2y_{ij}q_{ij} + (1 - y_{ij})q_{ij} + y_{ij}(1 - q_{ij})} \Big), \quad (8)$$

and

$$\hat{L}_{DSC_M}(\boldsymbol{x_i}, \boldsymbol{y_i}) = \sum_{j=1}^{c} \Big( \frac{(1 - y_{ij})\hat{q}_{ij} + y_{ij}(1 - \hat{q}_{ij})}{2y_{ij}\hat{q}_{ij} + (1 - y_{ij})\hat{q}_{ij} + y_{ij}(1 - \hat{q}_{ij})} \Big). \quad (9)$$

### B. Asymmetric focal loss

The focal loss was proposed for small object detection by reducing the weight for well-classified samples and focusing on samples which are near the decision boundary [29]. It adds attenuation inside the loss function based on the logit activations:

$$L_{CE_{focal}}(\boldsymbol{x_i}, \boldsymbol{y_i}) = -\sum_{j=1}^{c} (1 - p_{ij})^\gamma y_{ij} \log(p_{ij}), \quad (10)$$

where $\gamma$ is the hyper-parameter to control the focus. The symmetric focal loss prevents logits from being too large and makes every class stay near the decision boundary. However, this makes it likely for the unseen foreground samples to shift across the decision boundary. We remove the loss attenuation for the foreground class to keep it away from the decision boundary:

$$\hat{L}_{CE_{focal}}(\boldsymbol{x_i}, \boldsymbol{y_i}) = \sum_{j=1}^{c} \Big( - r_j y_{ij} \log(p_{ij}) \\ - (1 - r_j)(1 - p_{ij})^\gamma y_{ij} \log(p_{ij}) \Big). \quad (11)$$

Inspired by the focal loss [29], related work integrates a similar attenuation term into the DSC loss [1], [43]. In practice, we find that the logarithmic DSC loss [1] significantly changes the magnitude of DSC loss making it difficult to be combined with other losses. The attenuation in the focal Tversky loss [43] is very large and may overly suppress the easier class. Here, we propose another form of DSC loss with an adaptive weight preserving a similar loss magnitude. Specifically, we add the attenuation term to the false negatives part of the function and prevent the network being too confident about its prediction:

$$L_{DSC_{focal}}(\boldsymbol{x_i}, \boldsymbol{y_i}) = \sum_{j=1}^{c} \Big( \frac{(1 - y_{ij})p_{ij} + (1 - p_{ij})^\gamma y_{ij}(1 - p_{ij})}{2y_{ij}p_{ij} + (1 - y_{ij})p_{ij} + y_{ij}(1 - p_{ij})} \Big), \quad (12)$$

Compared with the original version of the CE loss, this formulation for the DSC loss has a similar effect of reducing the penalty for the well-classified samples while keeping the magnitude of the loss similar to the original one, as shown
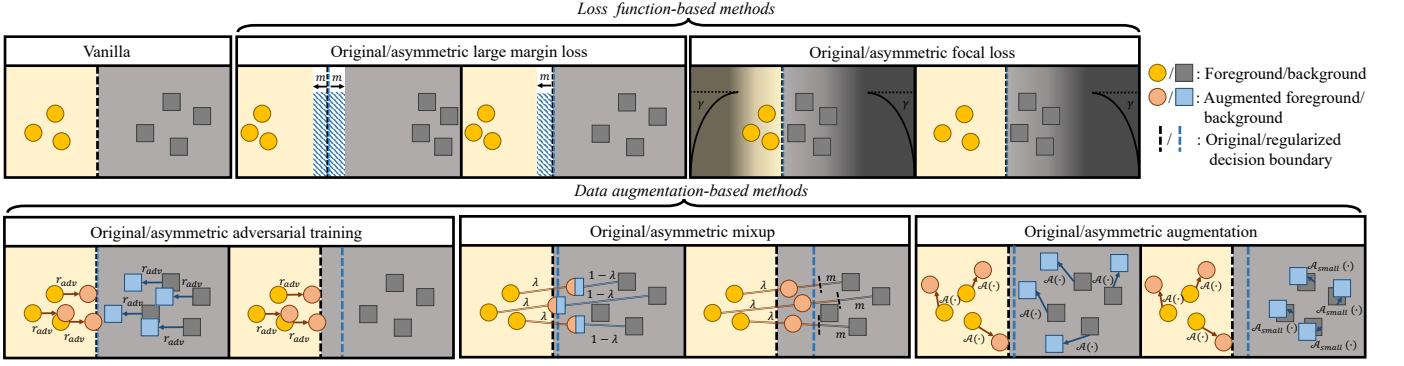
Fig. 6. The illustration of the proposed asymmetric modifications for the existing loss functions and regularization techniques. We make the logit activations of foreground class far away from the decision boundary by setting a bias for the foreground class in different ways.

in Supplementary Fig. 9. We refer to this as the focal DSC loss in the following. Similarly, the asymmetric version of the focal DSC loss is derived by removing the attenuation term for the foreground class:

$$\hat{L}_{DSC_{focal}}(\boldsymbol{x_i}, \boldsymbol{y_i}) = \sum_{j=1}^{c} \Big( \frac{(1-y_{ij})p_{ij} + r_j y_{ij}(1-p_{ij})}{2y_{ij}p_{ij} + (1-y_{ij})p_{ij} + y_{ij}(1-p_{ij})}$$
$$+ \frac{(1-y_{ij})p_{ij} + (1-r_j)(1-p_{ij})^{\gamma} y_{ij}(1-p_{ij})}{2y_{ij}p_{ij} + (1-y_{ij})p_{ij} + y_{ij}(1-p_{ij})} \Big). \tag{13}$$

### C. Asymmetric adversarial training

Adversarial training was proposed to learn more robust classifiers by training with difficult samples [12]. The network is trained by considering adversarial samples as additional training data [34], [44]:

$$L_{adv}(\boldsymbol{x_i}, \boldsymbol{y_i}) = L(\boldsymbol{x_i}, \boldsymbol{y_i}) + L(\boldsymbol{x_i} + l \cdot \frac{\boldsymbol{d_{adv}}}{\|\boldsymbol{d_{adv}}\|_2}, \boldsymbol{y_i}), \quad (14)$$

$$\text{with} \quad \boldsymbol{d_{adv}} = \underset{\boldsymbol{d}; \|\boldsymbol{d}\| < \epsilon}{\arg\max} L(\boldsymbol{x_i} + \boldsymbol{d}, \boldsymbol{y_i}). \tag{15}$$

Here, $\boldsymbol{d_{adv}}$ is the direction of the generated adversarial samples, $l$ and $\epsilon$ are the magnitude and the range of the adversarial perturbations, respectively. $L$ is the chosen loss function, which can be $L_{CE}$ and / or $L_{DSC}$. Similar to the large margin loss, symmetric adversarial training preserves the decision boundary and may cause difficulties for unseen foreground samples, which tends to shift towards background class. Our proposed asymmetric adversarial training aims to produce a larger space between the foreground class and the decision boundary. Specifically, we generate samples by considering more from the rare classes:

$$\hat{\boldsymbol{d}}_{\boldsymbol{adv}} = \underset{\boldsymbol{d}; \|\boldsymbol{d}\| < \epsilon}{\arg\max} L(\boldsymbol{x_i} + \boldsymbol{d}, \boldsymbol{y_i} \odot \boldsymbol{r}) \Big|_{\boldsymbol{y_i} \cdot \boldsymbol{r} > 0}, \tag{16}$$

where "$\odot$" refers to the element product and "$\cdot$" refers to the dot product.

### D. Asymmetric mixup

Mixup is a simple yet effective data augmentation algorithm to improve generalization by generating extra training samples by using the linear combination of pairs of images and their labels [47]:

$$L_{mixup}(\boldsymbol{x_i}, \boldsymbol{y_i}, \boldsymbol{x_k}, \boldsymbol{y_k}) = L(\boldsymbol{x_i}, \boldsymbol{y_i}) + L(\tilde{\boldsymbol{x}}_i, \tilde{\boldsymbol{y}}_i), \quad (17)$$

where $(\tilde{\boldsymbol{x}}_i, \tilde{\boldsymbol{y}}_i)$ is the generated training sample:

$$\tilde{\boldsymbol{x}}_i = \lambda \boldsymbol{x_i} + (1-\lambda)\boldsymbol{x_k}, \quad \tilde{\boldsymbol{y}}_i = \lambda \boldsymbol{y_i} + (1-\lambda)\boldsymbol{y_k}. \tag{18}$$

Here, $\lambda$ is randomly selected based on a beta distribution, $(\boldsymbol{x_k}, \boldsymbol{y_k})$ is another random training sample. Mixup regularizes the model by centering the decision boundary between classes which helps very little in our setting. Different from the original mixup, which generates samples with soft labels, our modification generates hard labels by considered augmented samples near to the foreground samples as foreground class. Asymmetric mixup can keep the decision boundary away from the foreground class and increase the area of the foreground logit distribution. This prevents unseen under-presented samples from shifting across the decision boundary. Specifically, the mixed image $\tilde{\boldsymbol{x}}_i$ which has a certain distance from the background class, is taken as a foreground sample:

$$\hat{\tilde{\boldsymbol{y}}}_i = \begin{cases} \boldsymbol{y_i} & \text{if } (\lambda > m \text{ and } \boldsymbol{y_i} \cdot \boldsymbol{r}(1-\boldsymbol{y_k} \cdot \boldsymbol{r}) == 1) \text{ or } \boldsymbol{y_i} == \boldsymbol{y_k}, \\ \boldsymbol{y_k} & \text{if } (1-\lambda > m \text{ and } \boldsymbol{y_k} \cdot \boldsymbol{r}(1-\boldsymbol{y_i} \cdot \boldsymbol{r}) == 1) \text{ or } \boldsymbol{y_i} == \boldsymbol{y_k}, \\ \boldsymbol{0} & \text{otherwise}, \end{cases} \tag{19}$$

where $m$ is the margin to guarantee that the augmented samples are not getting too close to background samples. In practice, we do not update the model using training samples with $\hat{\tilde{\boldsymbol{y}}}_i = \boldsymbol{0}$.

### E. Asymmetric augmentation

In order to extend the latent space of the foreground class, we also evaluate a simple method to compensate class imbalance by adjusting the magnitude of augmentation for different classes. Standard data augmentation methods would preserve the label $\boldsymbol{y_i}$ and adopt the same set of heuristic transformations such as scaling and rotations to the original

training sample $x_i$ for different classes. The generated training sample $\tilde{x}_i$ can be obtained using:

$$\tilde{x}_i = \mathcal{A}(x_i), \qquad (20)$$

where $\mathcal{A}$ is the chosen transformation with certain probability. When the dataset is highly imbalanced, adding more synthesized background samples is not necessary. Our simple variant of data augmentation reduces the number of transformed samples for the background classes. In this asymmetric setting, the generated sample $\tilde{x}_i$ is obtained using:

$$\hat{\tilde{x}}_i = \begin{cases} \mathcal{A}(x_i) & \text{if} \quad y_i \cdot r == 1, \\ \mathcal{A}_{small}(x_i) & \text{otherwise}, \end{cases} \qquad (21)$$

where $\mathcal{A}_{small}$ is transformations with smaller probability.

### F. The combination of asymmetric techniques

The above-mentioned modifications would introduce more variances for the under-represented classes in the latent space or the image space by adding a bias for the foreground class from different perspectives. In practice, some or all of the techniques can be integrated into a single model to combat overfitting under class imbalance.

Specifically, we can first generate different sets of the augmented samples following the asymmetric adversarial training, the asymmetric mixup and the asymmetric augmentation following equation 16, 19 and 21, separately. The network can then be optimized using the extended training set with the loss functions combined with the asymmetric large margin loss and asymmetric focal loss:

$$
\begin{aligned}
\hat{L}_{CE_{combine}}(x_i, y_i) = \sum_{j=1}^{c} \Big( & -r_j y_{ij} \log(\hat{q}_{ij}) \\
& - (1-r_j)(1-\hat{q}_{ij})^{\gamma} y_{ij} \log(\hat{q}_{ij}) \Big).
\end{aligned}
\qquad (22)
$$

A combined DSC loss can be formulated in a similar way.

### V. EXPERIMENTS

### A. Experimental setup

We demonstrate the effect of our proposed modifications with a variety of medical image segmentation tasks using different models and training scenarios. Here, we summarize the dataset splits and experimental settings, which are kept the same with motivational experiments in Section III. We keep the hyper-parameters of the methods the same for the original baselines and our modified techniques. The hyper-parameters are summarized in Supplementary Table V, VI and VII. Additionally, we conduct a sensitivity analysis of all the hyper-parameters and summarize the results in Supplementary Table VIII. We also provide the source code for our experiments[2].

[2] https://github.com/ZerojumpLine/OverfittingUnderClassImbalance

*1) Brain tumor segmentation:* We first evaluate the asymmetric techniques for the case of binary brain tumor core segmentation using the DeepMedic network architecture, a well performing method for this task [20]. To investigate the behavior under overfitting and to isolate better the effect of the objective functions, we do not use dropout, weight decay and data augmentation in this experiment. We train the network with CE loss, unless otherwise specified. By default, we sample 50% training samples from the foreground class. We conduct experiments using the training dataset of BRATS2017 dataset [2] which contains 285 four modalities Magnetic Resonance (MR) images. The MR images all have the same voxel space of $1.0 \times 1.0 \times 1.0$ mm. We test on 95 cases and train separate models using 190 (100%), 95 (50%), 38 (20%), 19 (10%) and 10 cases (5% of full training set).

*2) Brain stroke lesion segmentation:* We also evaluate the asymmetric techniques for the case of brain stroke lesion segmentation [28] again using DeepMedic. Here, we use a more realistic setting, employing standard regularization techniques including dropout, weight decay and data augmentation, as in the original work where the model achieved high performance for stroke lesion segmentation [20]. We implement our asymmetric techniques with the default training setting and default network architecture. The augmentation includes small intensity shifts and flipping in the sagittal plane wFith probability 0.5. The network is always trained with CE loss. We conduct experiments using ATLAS dataset [28] which contains 220 T1-weighted MR images. The MR images have the same voxel space of $1.0 \times 1.0 \times 1.0$ mm. We test on 75 cases and train separate models using 145 (100%), 73 (50%), 57 (40%) and 43 cases (30% of full training set).

*3) Small organ segmentation:* For organ segmentation in Section III, we use a default DeepMedic network. We conduct experiments using the training datatset of the abdominal organ segmentation challenge [25] which contains 30 computed tomography (CT) scans. We train the network to segment thirteen abdominal organs. We test on 10 cases and train models using 20 (100%) and 5 cases (25% of training set). We resample all the MR images to a common voxel spacing of $2.0 \times 2.0 \times 2.0$ mm. We show the segmentation results of representative small organs including gallbladder, aorta, inferior vena cava as well as portal vein and splenic vein in Fig. 1. Here, we use this dataset to further confirm the observations about the effect on sensitivity when training with varying amounts data.

*4) Kidney tumor segmentation:* In addition, we evaluate the asymmetric techniques for the case of kidney tumor segmentation. We train a well configured 3D U-Net [18] which includes extensive data augmentation with scaling, rotations, brightness, contrast, gamma and Gaussian noise augmentations with a predefined policy [16]. We also train DeepMedic with similar augmentation strategies yielding lower accuracy on this task. Therefore, we evaluate the asymmetric regularization techniques on the U-Net with kidney tumor segmentation. The task includes the segmentation of both the kidney and kidney tumor. As the segmentation of kidney is relatively easy, in this experiment we only focus on tumor segmentation and only take kidney tumor as the foreground class to implement

TABLE II
EVALUATION OF BRAIN TUMOR CORE SEGMENTATION USING DEEPMEDIC WITH DIFFERENT AMOUNTS OF TRAINING DATA AND DIFFERENT TECHNIQUES TO COUNTER OVERFITTING. THE RESULTS ARE CALCULATED WITH POST-PROCESSING. RESULTS WHICH HAVE WORSE DSC THAN THE VANILLA BASELINE ARE HIGHLIGHTED WITH SHADING. BEST AND SECOND BEST RESULTS ARE IN BOLD WITH THE BEST ALSO UNDERLINED.

| Method | 5% training | | | | 10% training | | | | 20% training | | | | 50% training | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | DSC | SEN | PRC | HD | DSC | SEN | PRC | HD | DSC | SEN | PRC | HD | DSC | SEN | PRC | HD |
| Vanilla - CE [20] | 50.4 | 41.0 | 83.5 | 18.0 | 62.5 | 56.0 | 83.1 | 14.3 | 64.9 | 59.8 | 85.7 | 13.8 | 69.4 | 65.4 | 85.3 | 15.7 |
| Vanilla - CE - 80% tumor | 45.5 | 36.0 | 86.7 | 17.8 | 61.5 | 54.2 | 81.7 | 18.5 | 65.3 | 59.6 | 85.0 | 15.1 | 68.6 | 64.1 | 86.1 | 14.8 |
| Vanilla - F1 (DSC) | 47.2 | 37.4 | 86.6 | 15.9 | 58.9 | 51.1 | 83.6 | 20.1 | 64.3 | 58.1 | 83.5 | 16.3 | 67.1 | 62.5 | 86.5 | 15.3 |
| Vanilla - F2 [14] | 45.8 | 36.9 | 81.9 | 17.9 | 59.3 | 52.2 | 84.9 | 18.0 | 66.4 | 61.1 | 83.4 | 14.1 | 68.8 | 66.0 | 83.4 | 13.7 |
| Vanilla - F4 [14] | 51.6 | 42.5 | 83.8 | 18.1 | 59.6 | 53.0 | 82.9 | 18.4 | 65.9 | 61.9 | 85.4 | 14.2 | 67.5 | 64.5 | 84.9 | 13.7 |
| Vanilla - F8 [14] | 47.4 | 38.7 | 83.1 | 19.6 | 59.8 | 52.4 | 87.0 | 15.4 | 64.5 | 60.3 | 85.2 | 14.7 | 67.9 | 65.4 | 81.6 | 14.9 |
| Large margin loss [31] | 44.5 | 35.9 | 82.8 | 20.2 | 60.9 | 53.5 | 84.0 | 17.6 | 67.0 | 61.6 | 86.1 | 14.4 | 66.5 | 62.2 | 88.1 | 13.7 |
| Asymmetric large margin loss | 56.8 | 48.9 | 83.4 | **15.0** | 64.0 | 56.8 | 87.0 | 13.9 | 67.4 | 62.9 | 84.1 | 15.9 | 68.9 | 64.9 | 86.5 | 14.1 |
| Focal loss [29] | 54.0 | 44.8 | 82.6 | 16.0 | 62.6 | 55.1 | 84.3 | 17.7 | 64.9 | 60.0 | 84.4 | 19.5 | 67.0 | 62.3 | 87.0 | 16.5 |
| Asymmetric focal loss | 58.8 | 51.4 | 81.6 | **15.0** | 66.8 | 62.0 | 83.2 | **13.2** | 68.9 | 66.2 | 83.3 | **12.5** | **71.5** | 70.6 | 83.7 | 12.1 |
| Adversarial training [12] | 53.2 | 44.6 | 85.0 | 19.2 | 62.0 | 55.0 | 84.8 | 20.6 | 64.6 | 59.4 | 84.6 | 17.3 | 65.6 | 61.2 | 86.0 | 19.4 |
| Asymmetric adversarial training | 58.5 | 50.8 | 80.1 | 16.2 | 63.9 | 58.2 | 83.1 | 17.2 | 67.7 | 63.7 | 84.2 | 17.0 | 70.5 | 68.4 | 83.0 | 14.8 |
| Mixup [47] | 49.7 | 40.9 | 83.0 | 19.6 | 60.3 | 53.9 | 83.1 | 21.2 | 63.9 | 58.5 | 84.1 | 18.2 | 66.4 | 61.5 | 86.8 | 19.0 |
| Asymmetric mixup | **59.8** | 56.8 | 74.7 | 17.7 | **68.5** | 65.1 | 80.7 | 15.3 | **70.8** | 67.9 | 85.1 | **11.6** | 70.7 | 67.9 | 85.4 | **11.8** |
| Symmetric combination | 50.0 | 42.0 | 84.6 | 21.1 | 60.3 | 53.1 | 84.7 | 25.1 | 64.1 | 58.3 | 86.6 | 19.1 | 67.2 | 63.1 | 86.6 | 15.1 |
| Asymmetric combination | **63.4** | 63.1 | 75.9 | **15.1** | **72.4** | 72.9 | 78.3 | **10.8** | **71.6** | 72.0 | 80.1 | 13.7 | **74.1** | 76.0 | 82.4 | **10.7** |

TABLE III
EVALUATION OF BRAIN STROKE LESION SEGMENTATION ON ATLAS BASED ON DEEPMEDIC WITH DIFFERENT AMOUNTS OF TRAINING DATA AND DIFFERENT TECHNIQUES TO COUNTER OVERFITTING. THE RESULTS ARE CALCULATED WITH POST-PROCESSING. RESULTS WHICH HAVE WORSE DSC THAN THE VANILLA BASELINE ARE HIGHLIGHTED WITH SHADING. BEST AND SECOND BEST RESULTS ARE IN BOLD WITH THE BEST ALSO UNDERLINED.

| Method | 30% training | | | | 50% training | | | | 100% training | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | DSC | SEN | PRC | HD | DSC | SEN | PRC | HD | DSC | SEN | PRC | HD |
| Vanilla - w/ augmentation [20] | 22.2 | 18.3 | 60.9 | 48.6 | 45.2 | 40.9 | 59.7 | 31.1 | 54.5 | 49.5 | 67.3 | 32.2 |
| Vanilla - w/o augmentation | 15.0 | 11.7 | 59.1 | 51.9 | 40.3 | 35.9 | 53.0 | 40.2 | 51.7 | 48.0 | 62.3 | 31.9 |
| Vanilla - asymmetric augmentation | 22.4 | 18.8 | 58.0 | 50.2 | 47.3 | 43.4 | 57.5 | 32.1 | 56.9 | 51.9 | 69.8 | 28.2 |
| Large margin loss [31] | 18.9 | 14.8 | 64.4 | 48.8 | 45.3 | 40.7 | 60.5 | 36.8 | 55.1 | 49.4 | 70.0 | 28.0 |
| Asymmetric large margin loss | 23.5 | 19.8 | 58.6 | 45.9 | 47.7 | 44.3 | 58.4 | 33.8 | **57.6** | 54.1 | 67.6 | **27.8** |
| Focal loss [29] | 20.4 | 16.7 | 62.7 | 47.9 | 46.9 | 41.8 | 61.7 | 31.4 | 56.0 | 50.8 | 69.1 | 30.9 |
| Asymmetric focal loss | 26.3 | 22.2 | 59.0 | 46.4 | 49.0 | 47.8 | 56.3 | 31.7 | 56.6 | 63.2 | 55.6 | 27.9 |
| Adversarial training [12] | 20.1 | 16.7 | 57.3 | 56.9 | 47.2 | 41.6 | 62.6 | 35.0 | 54.0 | 48.5 | 69.7 | 34.9 |
| Asymmetric adversarial training | **28.1** | 23.3 | 63.9 | **43.4** | **50.2** | 46.5 | 64.2 | **29.6** | 55.5 | 51.6 | 69.5 | 33.6 |
| Mixup [47] | 14.5 | 12.0 | 52.7 | 54.1 | 45.7 | 41.7 | 59.5 | 30.3 | 53.8 | 50.0 | 67.8 | 29.2 |
| Asymmetric mixup | 22.8 | 20.9 | 50.9 | 45.9 | 49.0 | 49.2 | 53.9 | 31.5 | 57.0 | 51.0 | 74.6 | 31.7 |
| Symmetric combination | 22.2 | 17.7 | 66.1 | 49.4 | 48.5 | 44.5 | 61.0 | 31.6 | 56.0 | 50.4 | 71.5 | 29.9 |
| Asymmetric combination | **31.1** | 27.9 | 57.6 | **42.1** | **52.2** | 50.6 | 59.5 | **27.8** | **58.5** | 54.9 | 70.9 | **27.3** |

asymmetric techniques. To be specific, we always set $r$ as $[0, 0, 1]^{\mathsf{T}}$. The network is always trained with both CE and sample-wise DSC loss. The two losses have the same weight. We conduct experiments using the training dataset of KiTS19 dataset [16] which contains 210 CT images. We resample all the CT images to a common voxel spacing of $1.6\times1.6\times3.2$ mm following [18]. The original and asymmetric techniques are implemented on both loss functions with the same hyper-parameters. We test on 70 cases and train separate models using 140 (100%), 70 (50%), 28 (20%), 14 (10%) and 7 cases (5% of full training set). Note, we do not manually clean the training dataset or employ deep supervision as done in [18], therefore the absolute performance of our U-Net is not directly comparable with the reported challenge results.

### B. Quantitative results

Taking the provided manual segmentations as the ground truth, we calculate DSC, sensitivity (SEN), precision (PRC) and 95% Hausdorff distance (HD) (mm) to evaluate the segmentation accuracy. The initial segmentation results of our method always have higher sensitivity and DSC, but sometimes would cause more false positive predictions and

therefore lead to worse distance-based metrics such as HD. We argue that in practice this problem can be addressed by taking advantage of some connected component-based post-processing, which is widely adopted in many segmentation methods [18]. Specifically, we assume there is only one target component and suppress all but the largest region. We report both results with or without these post-processing operations. The quantitative segmentation results on BRATS, ATLAS and KiTS datasets using different amounts of training data with post-processing are summarized in Table II, III and IV, respectively. The corresponding quantitative segmentation results without post-processing are summarized in Supplementary Table X, XI and XII. We also evaluate one of the proposed methods, asymmetric focal loss, with abdominal organ segmentation in which multiple classes are considered under-represented. The experiments are summarized in Supplementary Table IX.

Class imbalance affects the segmentation sensitivity of the under-represented class, as shown in Section III. We find that previous attempts to tackle class imbalance do not improve sensitivity, while our asymmetric methods do lead to better results with higher sensitivity across different tasks. This indicates that the proposed methods may effectively mitigate

TABLE IV

EVALUATION OF KIDNEY AND KIDNEY TUMOR SEGMENTATION BASED ON 3D U-NET WITH DIFFERENT AMOUNTS OF TRAINING DATA AND DIFFERENT TECHNIQUES TO COUNTER OVERFITTING. THE RESULTS ARE CALCULATED WITH POST-PROCESSING. RESULTS WHICH HAVE WORSE DSC THAN THE VANILLA BASELINE ARE HIGHLIGHTED WITH SHADING. BEST AND SECOND BEST RESULTS ARE IN BOLD WITH THE BEST ALSO UNDERLINED.

| Method | Kidney | | | | | | | | | | | |
| | 10% training | | | | 50% training | | | | 100% training | | | |
| | DSC | SEN | PRC | HD | DSC | SEN | PRC | HD | DSC | SEN | PRC | HD |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Vanilla - w/ augmentation [18] | 93.3 | 91.2 | 96.9 | 5.4 | 96.4 | 95.8 | 97.1 | 2.7 | 96.6 | 96.1 | 97.3 | 2.4 |
| Vanilla - w/o augmentation | 92.3 | 89.3 | 96.8 | 12.1 | 96.1 | 95.6 | 96.7 | 2.8 | 96.3 | 95.8 | 96.9 | 2.7 |
| Vanilla - asymmetric augmentation | 94.3 | 92.2 | 97.0 | 5.2 | 94.9 | 94.5 | 95.5 | 5.9 | 96.1 | 95.8 | 96.4 | 3.8 |
| Large margin loss [31] | **94.6** | 92.7 | 97.1 | 4.8 | 96.4 | 95.9 | 97.0 | 2.8 | 96.1 | 95.9 | 96.3 | 3.2 |
| Asymmetric large margin loss | 93.8 | 91.4 | 97.2 | 5.3 | 96.1 | 95.5 | 96.9 | 2.9 | **96.8** | 96.6 | 97.1 | **2.2** |
| Focal loss [29] | 91.4 | 85.9 | 99.2 | 10.6 | 94.1 | 89.6 | 99.2 | 4.2 | 94.3 | 90.0 | 99.1 | 4.2 |
| Asymmetric focal loss | 92.0 | 86.7 | 99.0 | 6.0 | 94.7 | 90.9 | 98.9 | 3.5 | 94.8 | 90.9 | 99.1 | 3.1 |
| Adversarial training [12] | 94.1 | 91.9 | 97.3 | 9.1 | 96.3 | 95.7 | 97.1 | 2.6 | 96.6 | 96.2 | 97.2 | **2.3** |
| Asymmetric adversarial training | 94.4 | 92.5 | 97.2 | 5.7 | **96.6** | 96.0 | 97.3 | **2.5** | **96.8** | 96.4 | 97.3 | **2.3** |
| Mixup [47] | **95.0** | 93.2 | 97.3 | **4.2** | **96.8** | 96.2 | 97.5 | **2.3** | **96.9** | 96.4 | 97.5 | **2.2** |
| Asymmetric mixup | **94.6** | 92.6 | 97.3 | **4.5** | 96.0 | 95.2 | 97.0 | 3.1 | 96.4 | 95.7 | 97.3 | 2.7 |
| Symmetric combination | 94.1 | 91.4 | 97.5 | 5.1 | 94.6 | 91.0 | 98.7 | 4.3 | 96.7 | 96.2 | 97.2 | **2.2** |
| Asymmetric combination | 93.5 | 89.7 | 98.5 | 5.2 | 93.9 | 90.0 | 98.3 | 5.3 | 96.7 | 95.6 | 97.9 | **2.2** |

| Method | Kidney tumor | | | | | | | | | | | |
| | 10% training | | | | 50% training | | | | 100% training | | | |
| | DSC | SEN | PRC | HD | DSC | SEN | PRC | HD | DSC | SEN | PRC | HD |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Vanilla - w/ augmentation [18] | 54.6 | 46.0 | 80.0 | **53.2** | 76.0 | 72.8 | 86.1 | 25.1 | 79.2 | 77.0 | 86.2 | **17.8** |
| Vanilla - w/o augmentation | 37.4 | 31.5 | 65.6 | 96.0 | 62.8 | 58.7 | 75.9 | 47.8 | 73.0 | 69.1 | 83.4 | 18.9 |
| Vanilla - asymmetric augmentation | 55.9 | 48.2 | 76.4 | 71.5 | 74.3 | 70.3 | 85.2 | 33.3 | 78.4 | 76.9 | 85.7 | 19.8 |
| Large margin loss [31] | 52.2 | 44.3 | 77.2 | 68.5 | 78.2 | 74.3 | 87.8 | 26.6 | 80.2 | 79.1 | 84.5 | 25.5 |
| Asymmetric large margin loss | 55.5 | 48.3 | 77.4 | 71.6 | **78.4** | 74.9 | 87.5 | 24.1 | **82.3** | 81.4 | 86.0 | **16.9** |
| Focal loss [29] | 47.1 | 37.5 | 78.2 | 74.5 | 73.0 | 66.0 | 87.6 | 40.2 | 79.0 | 73.2 | 90.0 | 20.3 |
| Asymmetric focal loss | **57.9** | 48.9 | 78.4 | 61.4 | 77.4 | 74.4 | 85.0 | **20.2** | 81.5 | 80.6 | 86.7 | 19.4 |
| Adversarial training [12] | 50.9 | 42.5 | 81.3 | 62.0 | 73.2 | 69.6 | 83.9 | 44.1 | 81.9 | 81.1 | 85.8 | 27.6 |
| Asymmetric adversarial training | 55.2 | 47.8 | 79.6 | 66.7 | 78.3 | 74.9 | 87.9 | 23.7 | 82.1 | 81.1 | 87.4 | 19.7 |
| Mixup [47] | 53.3 | 45.2 | 81.6 | 57.8 | 77.0 | 72.9 | 87.3 | 32.1 | 80.3 | 78.5 | 85.9 | 34.1 |
| Asymmetric mixup | 56.8 | 48.1 | 84.6 | 66.5 | 77.9 | 74.0 | 89.2 | 22.0 | 79.7 | 78.1 | 87.3 | 19.3 |
| Symmetric combination | 53.9 | 45.1 | 81.3 | 70.2 | 73.9 | 67.1 | 87.7 | 39.6 | 80.9 | 79.3 | 86.5 | 19.6 |
| Asymmetric combination | **59.2** | 52.2 | 80.3 | **49.5** | **79.4** | 77.0 | 86.7 | **15.5** | **82.7** | 82.1 | 87.0 | 18.8 |

overfitting under class imbalance. These results in turn support our previous analysis of logit shift under class imbalance and indicate that considering this implication would help build unbiased network.

*1) Baseline experiments:* We perform baseline experiments with binary brain tumor core segmentation, as shown in Table II. We show that increasing the weight of tumor samples (from 50% to 80%) decreases performance when the dataset is highly imbalanced. This is because increasing the weight encourages the network to memorize the under-represented samples and may actually lead to more overfitting, thus being counter-productive. Simply changing the objective function to F-score (defined as $F_\beta = (1 + \beta^2)\frac{\text{sensitivity}\cdot\text{precision}}{\beta^2\text{sensitivity}+\text{precision}}$), which is a balancing loss and weights sensitivity $\beta$-times more than precision [14], [33], only shows little improvements increasing the sensitivity slightly. Changing the sampling weights or training with a loss function using F-scores seems to have little impact when foreground training samples are limited, with training accuracy close to 100%, as shown in Fig. 1.

A common approach to alleviate under-segmentation is to adjust thresholds of decision boundaries based on validation sets. In this work, however, we observe distribution shift on unseen test data, which can significantly differ from the training/validation sets. Hence, a threshold selected on validation data may not be optimal for new test data. Due to the lack of ground truth, it is practically not possible to optimise the decision thresholds for a specific test set.

*2) Asymmetric large margin loss:* The original large margin loss decreases performance in some cases, while our modification yields improvements over the symmetric version in all cases.

*3) Asymmetric focal loss:* The original focal loss also decreases the sensitivity in some cases and leads to worse performance. It is because the focal term would decrease the weight of foreground samples and push its logit closer to the decision boundary, making it easier to cause false negative predictions. Our modification removes the loss attenuation for the under-represented class and improve the performance in all cases. We notice that the asymmetric focal loss would make the performance of other class (kidney) overfit more, but it is not the focus of this study and can be easily addressed by just keeping focal term for the background class.

*4) Asymmetric adversarial training:* When the network is trained without data augmentation (as shown in Table II), the original adversarial training seems to be effective when little training data is available while our modifications can further improve the sensitivity and boost the performance substantially.

When the network is trained with data augmentation (as shown in Table III and IV), we find the original adversarial training does not improve the performance when training data is limited. It indicates that in this case the augmented samples by adversarial training might not add anything on top of intensity augmentation. In contrast, our proposed modifi-
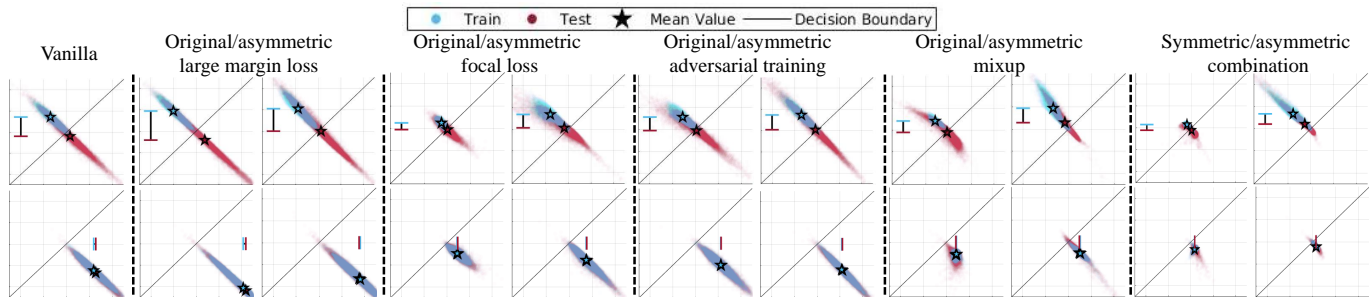
Fig. 7. Activations of the classification layer when processing tumor (top) and background (bottom) samples of BRATS with DeepMedic, using 5% training data. Asymmetric modifications lead to better separation of the logits of unseen tumor samples.
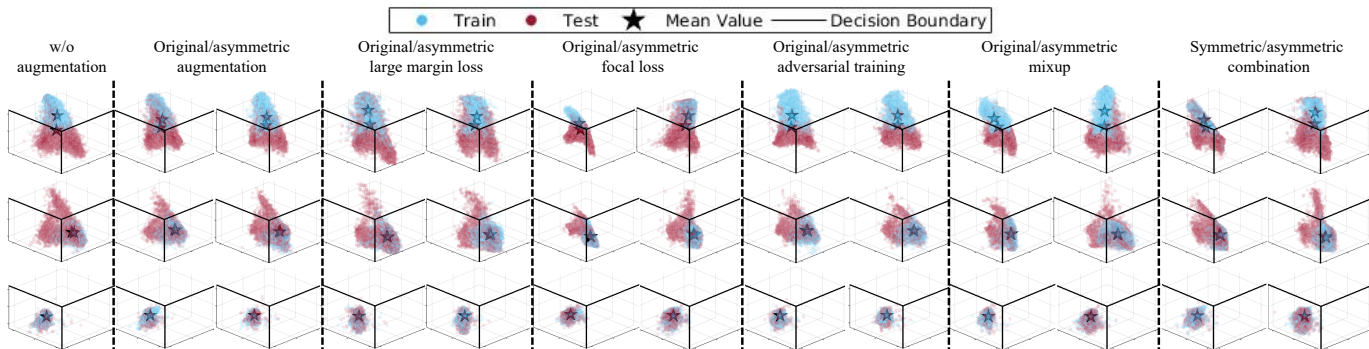


Fig. 8. Activations of the classification layer for tumor (top), kidney (middle) and background (bottom) samples of KiTS with 3D U-Net, using 10% training data. Asymmetric modifications also lead to better separation of the logits of unseen tumor samples and is complementary to standard data augmentation.

cations seem to help always leading to better segmentation performance.

*5) Asymmetric mixup:* We find the original mixup can be effective for the well-represented class. For example, it always improve the segmentation performance of kidney, as shown in Table IV. However, the original mixup leads to lower sensitivity for the under-represented class. We find that the asymmetric mixup can improve the performance for BRATS to a large extent, as shown in Table II. However, we also find it to be less effective for ATLAS and KiTS which only have one image channel, as shown in Table III and IV. This may be because the intensity distributions of healthy and lesion regions in ATLAS and KiTS overlap too much, as shown in the Supplementary Fig. 10. In this case, the mixed samples, which are very similar to the background samples in intensity but are taken as the foreground samples, may confuse the network. For BRATS, however, four image channels are available with the healthy and tumor regions show larger differences in T2 and fluid-attenuated inversion recovery (FLAIR) sequences. The mixed samples seem to take good advantage of the intensity relationship.

*6) Asymmetric augmentation:* Despite its simplicity, we find asymmetric augmentation to be an effective method to improve segmentation performance of the under-represented classes in most cases in terms of DSC and sensitivity, as summarized in Table III. However, we also notice that it could decrease the performance when data augmentation is strong and training data is sufficient, as shown in Table IV. It might be because the strong asymmetric augmentation would drive the model to focus too much on the foreground samples, making the background samples under-represented.

*7) The combination of asymmetric techniques:* We also combine the asymmetric techniques which are found to improve the segmentation accuracy. The combination of the asymmetric techniques is a safe choice leading to the overall best segmentation results with improved sensitivity in all cases. In contrast, the combination of the symmetric counterparts is unable to mitigate overfitting and often decreases sensitivity.

### C. Logit distribution changes

The effects of all the techniques on the logit distributions of BRATS using 5% training data is presented in Fig. 7. Asymmetric techniques would increase the variances of the foreground class and expand its logit distribution. The original large margin loss and adversarial training try to push samples from different classes far from each other, however, the logits of unseen data remain in the center around the decision boundary and thus the predictions are not improved. The original large margin loss results in even larger shifts for the foreground samples. For our asymmetric modifications only the logits of foreground samples are pushed away and the unseen foreground logits tend to remain on the correct side of the decision boundary. The original focal loss encourages the network to prevent the logits of each class from staying too far from the decision boundary. However, it allows foreground logits to remain near the decision boundary which can result in false negative predictions on unseen samples. Our asymmetric focal loss removes the constraints for foreground samples. Original mixup encourages the symmetric distributions of different

classes but does not consider class imbalance. Asymmetric mixup exploits the latent space based on the relationship between samples to generate foreground samples and make the decision boundary stay near the background class. This leads to the largest improvement by increasing the region for the foreground logit distribution and reduce logit shift of unseen foreground samples. The combination of the four asymmetric techniques can stabilize the logits further more.

The effects of all the techniques on the logit distribution of KiTS using 10% training data is summarized in Fig. 8. The original data augmentation can reduce the logit shift of both unseen kidney and kidney tumor samples although it is not specifically designed to regularize the logit distribution. The asymmetric augmentation can further reduce the logit shift of the unseen tumor samples. The asymmetric large margin loss reduces the tumor logit shift towards the kidney class. Although the logit distribution is already regularized by the strong augmentation, the proposed asymmetric techniques provide further benefits in stabilizing the logits.

## VI. CONCLUSION

We study overfitting of neural networks under class imbalance by inspecting network behavior. We observe that when processing unseen under-represented samples, the logit activations tend to shift towards the decision boundary and the sensitivity decreases. This phenomenon is confirmed across a variety of different tasks and two popular different network architectures. We derive simple yet effective asymmetric variants of existing loss functions and regularization techniques to prevent overfitting. We show that our proposed methods can substantially improve segmentation performance under class imbalance in terms of DSC and increased sensitivity, outperforming previous solutions. We believe more regularization methods can be derived to alleviate this problem by considering the biased network behavior. We also believe that the plotting logit distributions may be useful network inspection tool and help to gain a better understanding network behavior under different training scenarios. In future work, we will investigate if monitoring of intermediate activations may provide further insights for other challenging settings such as domain shift or self-supervised learning.

## REFERENCES

[1] N. Abraham and N. M. Khan. A novel focal tversky loss function with improved attention u-net for lesion segmentation. In *2019 IEEE 16th International Symposium on Biomedical Imaging (ISBI 2019)*, pages 683–687. IEEE, 2019.

[2] S. Bakas, H. Akbari, A. Sotiras, M. Bilello, M. Rozycki, J. S. Kirby, J. B. Freymann, K. Farahani, and C. Davatzikos. Advancing the cancer genome atlas glioma mri collections with expert segmentation labels and radiomic features. *Sci. Data*, 4:170117, 2017.

[3] T. Brosch, L. Y. Tang, Y. Yoo, D. K. Li, A. Traboulsee, and R. Tam. Deep 3d convolutional encoder networks with shortcuts for multiscale feature integration applied to multiple sclerosis lesion segmentation. *IEEE transactions on medical imaging*, 35(5):1229–1239, 2016.

[4] M. Buda, A. Maki, and M. A. Mazurowski. A systematic study of the class imbalance problem in convolutional neural networks. *Neural Networks*, 106:249–259, 2018.

[5] K. Cao, C. Wei, A. Gaidon, N. Arechiga, and T. Ma. Learning imbalanced datasets with label-distribution-aware margin loss. In *Advances in Neural Information Processing Systems*, pages 1567–1578, 2019.

[6] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer. Smote: synthetic minority over-sampling technique. *Journal of artificial intelligence research*, 16:321–357, 2002.

[7] E. D. Cubuk, B. Zoph, D. Mane, V. Vasudevan, and Q. V. Le. Autoaugment: Learning augmentation strategies from data. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 113–123, 2019.

[8] E. D. Cubuk, B. Zoph, J. Shlens, and Q. V. Le. Randaugment: Practical data augmentation with no separate search. *arXiv preprint arXiv:1909.13719*, 2019.

[9] Y. Cui, M. Jia, T.-Y. Lin, Y. Song, and S. Belongie. Class-balanced loss based on effective number of samples. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9268–9277, 2019.

[10] J. Deng, J. Guo, N. Xue, and S. Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4690–4699, 2019.

[11] Q. Dong, S. Gong, and X. Zhu. Imbalanced deep learning by minority class incremental rectification. *IEEE transactions on pattern analysis and machine intelligence*, 41(6):1367–1381, 2018.

[12] I. J. Goodfellow, J. Shlens, and C. Szegedy. Explaining and harnessing adversarial examples. In *International Conference on Learning Representations (ICLR)*, 2015.

[13] H. Guo and H. L. Viktor. Learning from imbalanced data sets with boosting and data generation: the databoost-im approach. *ACM Sigkdd Explorations Newsletter*, 6(1):30–39, 2004.

[14] S. R. Hashemi, S. S. M. Salehi, D. Erdogmus, S. P. Prabhu, S. K. Warfield, and A. Gholipour. Asymmetric loss functions and deep densely-connected networks for highly-imbalanced medical image segmentation: Application to multiple sclerosis lesion detection. *IEEE Access*, 7:1721–1735, 2018.

[15] M. Havaei, A. Davy, D. Warde-Farley, A. Biard, A. Courville, Y. Bengio, C. Pal, P.-M. Jodoin, and H. Larochelle. Brain tumor segmentation with deep neural networks. *Medical image analysis*, 35:18–31, 2017.

[16] N. Heller, N. Sathianathen, A. Kalapara, E. Walczak, K. Moore, H. Kaluzniak, J. Rosenberg, P. Blake, Z. Rengel, M. Oestreich, et al. The kits19 challenge data: 300 kidney tumor cases with clinical context, ct semantic segmentations, and surgical outcomes. *arXiv preprint arXiv:1904.00445*, 2019.

[17] C. Huang, Y. Li, C. C. Loy, and X. Tang. Learning deep representation for imbalanced classification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5375–5384, 2016.

[18] F. Isensee, P. F. Jäger, S. A. Kohl, J. Petersen, and K. H. Maier-Hein. Automated design of deep learning methods for biomedical image segmentation. *arXiv preprint arXiv:1904.08128*, 2019.

[19] J. M. Johnson and T. M. Khoshgoftaar. Survey on deep learning with class imbalance. *Journal of Big Data*, 6(1):27, 2019.

[20] K. Kamnitsas, C. Ledig, V. F. Newcombe, J. P. Simpson, A. D. Kane, D. K. Menon, D. Rueckert, and B. Glocker. Efficient multi-scale 3d cnn with fully connected crf for accurate brain lesion segmentation. *Med. Image Anal.*, 36:61–78, 2017.

[21] B. Kang, S. Xie, M. Rohrbach, Z. Yan, A. Gordo, J. Feng, and Y. Kalantidis. Decoupling representation and classifier for long-tailed recognition. In *International Conference on Learning Representations (ICLR)*, 2020.

[22] H. Kervadec, J. Bouchtiba, C. Desrosiers, E. Granger, J. Dolz, and I. B. Ayed. Boundary loss for highly unbalanced segmentation. In *International conference on medical imaging with deep learning*, pages 285–296, 2019.

[23] S. Khan, M. Hayat, S. W. Zamir, J. Shen, and L. Shao. Striking the right balance with uncertainty. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 103–112, 2019.

[24] A. Krogh and J. A. Hertz. A simple weight decay can improve generalization. In *Advances in neural information processing systems*,

pages 950–957, 1992.

[25] B. A. Landman, Z. Xu, J. E. Igelsias, M. Styner, T. R. Langerak, and A. Klein. 2015 miccai multi-atlas labeling beyond the cranial vault – workshop and challenge. Accessed Dec. 2020. [Online]. Available: https://www.synapse.org/#!Synapse:syn3193805, doi: 10.7303/syn3193805.

[26] Y. Li, H. Zaragoza, R. Herbrich, J. Shawe-Taylor, and J. Kandola. The perceptron algorithm with uneven margins. In *ICML*, volume 2, pages 379–386, 2002.

[27] Z. Li, K. Kamnitsas, and B. Glocker. Overfitting of neural nets under class imbalance: Analysis and improvements for segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 402–410. Springer, 2019.

[28] S.-L. Liew, J. M. Anglin, N. W. Banks, M. Sondag, K. L. Ito, H. Kim, J. Chan, J. Ito, C. Jung, N. Khoshab, et al. A large, open source dataset of stroke anatomical brain images and manual lesion segmentations. *Scientific data*, 5:180011, 2018.

[29] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017.

[30] G. Litjens, T. Kooi, B. E. Bejnordi, A. A. A. Setio, F. Ciompi, M. Ghafoorian, J. A. Van Der Laak, B. Van Ginneken, and C. I. Sánchez. A survey on deep learning in medical image analysis. *Medical image analysis*, 42:60–88, 2017.

[31] W. Liu, Y. Wen, Z. Yu, and M. Yang. Large-margin softmax loss for convolutional neural networks. In *International Conference on Machine Leanring (ICML)*, pages 507–516, 2016.

[32] Z. Liu, Z. Miao, X. Zhan, J. Wang, B. Gong, and S. X. Yu. Large-scale long-tailed recognition in an open world. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2537–2546, 2019.

[33] F. Milletari, N. Navab, and S.-A. Ahmadi. V-net: Fully convolutional neural networks for volumetric medical image segmentation. In *2016 Fourth International Conference on 3D Vision (3DV)*, pages 565–571. IEEE, 2016.

[34] T. Miyato, S.-i. Maeda, M. Koyama, and S. Ishii. Virtual adversarial training: a regularization method for supervised and semi-supervised learning. *IEEE transactions on pattern analysis and machine intelligence*, 41(8):1979–1993, 2018.

[35] M. H. Savenije, M. Maspero, G. G. Sikkes, J. R. van der Voort van Zyp, A. N. TJ Kotte, G. H. Bol, and C. A. T. van den Berg. Clinical implementation of mri-based organs-at-risk auto-segmentation with convolutional networks for prostate radiotherapy. *Radiation Oncology*, 15:1–12, 2020.

[36] A. A. A. Setio, F. Ciompi, G. Litjens, P. Gerke, C. Jacobs, S. J. Van Riel, M. M. W. Wille, M. Naqibullah, C. I. Sánchez, and B. van Ginneken. Pulmonary nodule detection in ct images: false positive reduction using multi-view convolutional networks. *IEEE Trans. Med. Imaging*, 35(5):1160–1169, 2016.

[37] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1):1929–1958, 2014.

[38] K. Tang, J. Huang, and H. Zhang. Long-tailed classification by keeping the good and removing the bad momentum causal effect. *arXiv preprint arXiv:2009.12991*, 2020.

[39] V. V. Valindria, I. Lavdas, J. Cerrolaza, E. O. Aboagye, A. G. Rockall, D. Rueckert, and B. Glocker. Small organ segmentation in whole-body mri using a two-stage fcn and weighting schemes. In *MICCAI-MLMI*, pages 346–354. Springer, 2018.

[40] F. Wang, J. Cheng, W. Liu, and H. Liu. Additive margin softmax for face verification. *IEEE Signal Processing Letters*, 25(7):926–930, 2018.

[41] S. Wang, W. Liu, J. Wu, L. Cao, Q. Meng, and P. J. Kennedy. Training deep neural networks on imbalanced data sets. In *2016 international joint conference on neural networks (IJCNN)*, pages 4368–4374. IEEE, 2016.

[42] Y.-X. Wang, D. Ramanan, and M. Hebert. Learning to model the tail. In *Advances in Neural Information Processing Systems*, pages 7029–7039, 2017.

[43] K. C. Wong, M. Moradi, H. Tang, and T. Syeda-Mahmood. 3d segmentation with exponential logarithmic loss for highly unbalanced object sizes. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 612–619. Springer, 2018.

[44] C. Xie, M. Tan, B. Gong, J. Wang, A. Yuille, and Q. V. Le. Adversarial examples improve image recognition. *arXiv preprint arXiv:1911.09665*, 2019.

[45] S. Yun, D. Han, S. J. Oh, S. Chun, J. Choe, and Y. Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 6023–6032, 2019.

[46] B. Zadrozny, J. Langford, and N. Abe. Cost-sensitive learning by cost-proportionate example weighting. In *Third IEEE international conference on data mining*, pages 435–442. IEEE, 2003.

[47] H. Zhang, M. Cisse, Y. N. Dauphin, and D. Lopez-Paz. mixup: Beyond empirical risk minimization. In *International Conference on Learning Representations (ICLR)*, 2018.

[48] B. Zhou, Q. Cui, X.-S. Wei, and Z.-M. Chen. Bbn: Bilateral-branch network with cumulative learning for long-tailed visual recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9719–9728, 2020.

[49] P. Zhou and J. Feng. Understanding generalization and optimization performance of deep cnns. In *International Conference on Machine Learning*, 2018.

## Supplementary Material

### VII. The analysis of large margin loss from a sample re-weighting perspective

Some of the proposed methods are based on two well known loss functions, noted as focal loss and large margin loss. We argue that these two loss functions can be both seen as sample-level re-weighting methods, which change the magnitude of gradient of the network output by multiplying a scalar. Specifically, focal loss would decrease the weights of well-classified samples, while large margin loss would increase the weights of all samples, especially for well-classified samples. Here, we provide an analysis of these loss functions from sample re-weighting perspective.

Following the formulation in the main text, we first consider a CNN trained using cross-entropy loss with one sample $x_i$ and its one-hot label $y_i$:

$$L_{CE}(x_i, y_i) = -\sum_{j=1}^{c} y_{ij} \log(p_{ij}) = -y_i \cdot \log(p_i), \quad (23)$$

where $p_i$ is the calculated probability, which is a normalized term from the network output $z_i$:

$$p_i = \frac{e^{z_i}}{e^{z_i} \cdot 1}. \quad (24)$$

Let's look at the gradient of a general loss function $L(x_i, y_i)$ with respect to the network parameters $\theta$:

$$\frac{\partial L(x_i, y_i)}{\partial \theta} = \frac{\partial L(x_i, y_i)}{\partial z_i} \frac{\partial z_i}{\partial \theta}, \quad (25)$$

where the former term is associated with the design of loss functions and the latter is related to the network architecture. Considering a cross entropy loss $L_{CE}(x_i, y_i)$, we can have:

$$\begin{aligned}
\frac{\partial L_{CE}(x_i, y_i)}{\partial z_i} &= \frac{\partial L_{CE}(x_i, y_i)}{\partial p_i} \frac{\partial p_i}{\partial z_i} \\
&= -\frac{1}{p_i \cdot y_i} p_i \cdot y_i (y_i - p_i) = p_i - y_i.
\end{aligned} \quad (26)$$

In instance-level, we can assign different weights for different samples $x_i$ with different scalars $w_i$. The weights could be derive based on the frequency of samples or other heuristic rules. Typically, we multiply $L_{CE}(x_i, y_i)$ by $w_i$, and the gradient after re-weighting would become:

$$\frac{\partial \Big( w_i L_{CE}(x_i, y_i) \Big)}{\partial z_i} = w_i (p_i - y_i). \quad (27)$$

It can be seen that re-weighting would change the gradient of the network output for different samples and make the model fit better the chosen samples, which are assigned a larger weight $w_i$.

Similarly, we can also derive the gradient of the focal loss as:

$$\frac{\partial\Big(L_{CE_{focal}}(\boldsymbol{x_i},\boldsymbol{y_i})\Big)}{\partial \boldsymbol{z_i}} = \frac{\partial\Big(L_{CE_{focal}}(\boldsymbol{x_i},\boldsymbol{y_i})\Big)}{\partial \boldsymbol{p_i}}\frac{\partial \boldsymbol{p_i}}{\partial \boldsymbol{z_i}}$$
$$= \Big((1-\boldsymbol{p_i}\cdot\boldsymbol{y_i})^\gamma - \gamma\boldsymbol{p_i}\cdot\boldsymbol{y_i}\log(\boldsymbol{p_i}\cdot\boldsymbol{y_i})(1-\boldsymbol{p_i}\cdot\boldsymbol{y_i})^{\gamma-1}\Big)(\boldsymbol{p_i}-\boldsymbol{y_i})$$
$$= w_{i_{focal}}(\boldsymbol{p_i}-\boldsymbol{y_i}), \tag{28}$$

The weight term of focal loss $w_{i_{focal}}$ is a scalar and related to the sample probability $\boldsymbol{p_i}\cdot\boldsymbol{y_i}$. Generally speaking, $w_{i_{focal}}$ would decrease when $\boldsymbol{p_i}\cdot\boldsymbol{y_i}$ is large, therefore focal loss would make the model fit the easy cases less.

More interestingly, we next look into the effect of large margin loss on the gradient. Large margin loss would change the calculation of probability for the training process. Specifically, we substitute $\boldsymbol{p_i}$ with $\boldsymbol{q_i}$ to calculate the loss function, where we require:

$$\boldsymbol{q_i} = \frac{\mathrm{e}^{\boldsymbol{z_i}-\boldsymbol{y_i}m}}{\mathrm{e}^{\boldsymbol{z_i}-\boldsymbol{y_i}m}\cdot\mathbf{1}}, \tag{29}$$

where $m$ is the hyper-parameter for the margin. In this case, the gradient of large margin loss can be derived as:

$$\frac{\partial L_{CE_M}(\boldsymbol{x_i},\boldsymbol{y_i})}{\partial \boldsymbol{z_i}} = \frac{\partial L_{CE}(\boldsymbol{x_i},\boldsymbol{y_i})}{\partial \boldsymbol{q_i}}\frac{\partial \boldsymbol{q_i}}{\partial \boldsymbol{z_i}}$$
$$= \frac{\mathrm{e}^{\boldsymbol{z_i}}\cdot\mathbf{1}}{\mathrm{e}^{\boldsymbol{z_i}-\boldsymbol{y_i}m}\cdot\mathbf{1}}(\boldsymbol{p_i}-\boldsymbol{y_i}) = w_{i_M}(\boldsymbol{p_i}-\boldsymbol{y_i}). \tag{30}$$

The weight term of large margin loss $w_{i_M}$ is also a scalar and related to the network output $\boldsymbol{z_i}\cdot\boldsymbol{y_i}$. It can be seen that the existence of a margin $m$ would increase the gradient of sample $\boldsymbol{x_i}$. Moreover, $w_{i_M}$ would be larger as $\boldsymbol{z_i}\cdot\boldsymbol{y_i}$ becomes larger, therefore large margin loss would make the model fit the easy cases more, and keep the distribution of $\boldsymbol{z_i}$ away from the decision boundary.

The analysis for $L_{DSC}$ can be done in a similar way.

## VIII. THE MAGNITUDE OF FOCAL DSC LOSS

The proposed focal DSC loss has similar behaviour with the original of focal loss, as shown in Figure 9. In addition, it does not change the magnitude of loss too much compared with existing solutions [1], [43], making it easier to be combined with other losses. We find it is particularly important for our experiments with 3D U-net [18] because this framework adopts a loss function which is a combination of cross entropy and DSC loss.

## IX. HYPER-PARAMETERS OF THE REGULARIZATION TECHNIQUES

We summarize the hyper-parameters in Table V, Table VI and Table VII as a reference for practitioners. We find when the asymmetric regularization techniques are combined together, the network could be regularized too much. In this case,



Fig. 9. The comparison of focal loss with cross entropy and DSC loss. The behavior of focal loss for cross entropy and DSC loss are similar using the formulation in equation 10 and 12.

the model would not converge and even perform poorly on the training data. Therefore, we always choose hyper-parameters with smaller regularization magnitude for the experiments with the combined regularization. Empirically, we find decreasing the hyper-parameters of large margin loss and/or focal loss is a sensible choice.

TABLE V
HYPER-PARAMETERS OF EXPERIMENTS USING DEEPMEDIC WITH BRATS.

| | BRATS | 5% data | 10% data | 20% data | 50% data |
|---|---|---|---|---|---|
| Individual | large margin $m$ | 1 | 0.2 | 1 | 1 |
| | focal $\gamma$ | 2 | 4 | 4 | 4 |
| | adversarial $\epsilon$ | 1e-5 | 1e-5 | 1e-5 | 1e-5 |
| | adversarial $l$ | 10 | 10 | 10 | 20 |
| | mixup $\lambda$ (symmetric) | 0.2 | 0.2 | 0.2 | 0.2 |
| | mixup $\lambda$ (asymmetric) | 1 | 1 | 1 | 1 |
| | mixup $m$ | 0.2 | 0.2 | 0.2 | 0.1 |
| Combination | large margin $m$ | 1 | 0 | 0 | 0 |
| | focal $\gamma$ | 1.5 | 2 | 2 | 4 |
| | adversarial $\epsilon$ | 1e-5 | 1e-5 | 1e-5 | 1e-5 |
| | adversarial $l$ | 10 | 10 | 10 | 20 |
| | mixup $\lambda$ (symmetric) | 0.2 | 0.2 | 0.2 | 0.2 |
| | mixup $\lambda$ (asymmetric) | 1 | 1 | 1 | 1 |
| | mixup $m$ | 0.2 | 0.2 | 0.2 | 0.1 |

TABLE VI
HYPER-PARAMETERS OF EXPERIMENTS USING DEEPMEDIC WITH ATLAS.

| | ATLAS | 30% data | 50% data | 100% data |
|---|---|---|---|---|
| Individual | large margin $m$ | 0.1 | 3 | 2 |
| | focal $\gamma$ | 4 | 4 | 4 |
| | adversarial $\epsilon$ | 1e-5 | 1e-5 | 1e-5 |
| | adversarial $l$ | 10 | 10 | 10 |
| | mixup $\lambda$ (symmetric) | 0.2 | 0.2 | 0.2 |
| | mixup $\lambda$ (asymmetric) | 1 | 1 | 1 |
| | mixup $m$ | 0.8 | 0.8 | 0.2 |
| | Probability of background samples being augmented | 50% | 50% | 25% |
| Combination | large margin $m$ | 0.1 | 0 | 1 |
| | focal $\gamma$ | 4 | 3 | 2 |
| | adversarial $\epsilon$ | 1e-5 | 1e-5 | 1e-5 |
| | adversarial $l$ | 10 | 10 | 10 |
| | mixup $\lambda$ (symmetric) | 0.2 | 0.2 | 0.2 |
| | mixup $\lambda$ (asymmetric) | —— | —— | 1 |
| | mixup $m$ | —— | —— | 0.2 |
| | Probability of background samples being augmented | 50% | 50% | —— |

## X. SENSITIVITY ANALYSIS

We conduct a series of controlled experiments with different hyper-parameters to provide more practical details of the proposed regularization techniques. Specifically, we use a baseline DeepMedic model for brain tumor core segmentation with 5% training data of BRATS. The experimental details are

TABLE VII
HYPER-PARAMETERS OF EXPERIMENTS USING 3D U-NET WITH KITS.

| | KiTS | 10% data | 50% data | 100% data |
|---|---|---|---|---|
| Individual | large margin $m$ | 0.8 | 0.6 | 0.8 |
| | focal $\gamma$ | 6 | 6 | 6 |
| | adversarial $\epsilon$ | 1e-5 | 1e-5 | 1e-5 |
| | adversarial $l$ | 100 | 50 | 50 |
| | mixup $\lambda$ (symmetric) | 0.2 | 0.2 | 0.2 |
| | mixup $\lambda$ (asymmetric) | 1 | 1 | 1 |
| | mixup $m$ | 0.05 | 0.05 | 0.2 |
| | Probability of background samples being augmented | 0% | 0% | 50% |
| Combination | large margin $m$ | 0.8 | 0.2 | 0.8 |
| | focal $\gamma$ | 4 | 6 | 2 |
| | adversarial $\epsilon$ | —— | —— | —— |
| | adversarial $l$ | —— | —— | —— |
| | mixup $\lambda$ (symmetric) | —— | —— | —— |
| | mixup $\lambda$ (asymmetric) | —— | —— | —— |
| | mixup $m$ | —— | —— | —— |
| | Probability of background samples being augmented | 0% | —— | —— |

consistent with descriptions in Section V. We summarize the results with and without any post-processing in Table VIII. We can see from the results that the proposed methods can improve the baseline segmentation results with varied hyper-parameters in most cases. Specifically, asymmetric large margin loss yields improvements for most cases, however, a specific hyper-parameter may yield unexpected results (i.e. $m = 0.5$). A potential reason is that the model which focuses on a small portion of easy under-represented samples (c.f. equation 30 in Section VII) would overfit more. Asymmetric large margin loss with larger $m$ makes the model emphasize on more under-represented samples and therefore generalize better. Asymmetric adversarial training and asymmetric mixup yields considerable improvements when the perturbation in data augmentation is larger (i.e. $l > 2.5$ for asymmetric adversarial training and $m < 0.8$ for asymmetric mixup). Asymmetric focal loss is robust and can improve the segmentation results with all chosen hyper-parameters. Therefore, we recommend to choose asymmetric focal loss at first for new applications.

## XI. THE INTENSITY HISTOGRAM OF DIFFERENT DATASETS

Empirically, we find the asymmetric mixup is the most effective method for tumor segmentation with BRATS. However, asymmetric mixup show limited improvements for ATLAS and KiTS. We think it is because the multi-channel information in BRATS could create more useful information, as shown in Figure 10.

## XII. THE QUANTITATIVE RESULTS OF ABDOMINAL ORGAN SEGMENTATION

We evaluate one of our proposed techniques, asymmetric focal loss, with the application of abdominal organ segmentation to demonstrate our method can be feasibly applied to multi-class segmentation. Specifically, we train a model of basic DeepMedic using 25% of the training data, with the same setting in empirical experiments in Section III. Considering the class distribution of the dataset, as shown in Figure 11, we take class 4, class 5, class 8, class 9, class 10, class 11, class 12 and class 13 as rare classes. Specifically, we initiate the one-hot vector $r$ as $[0, 0, 0, 0, 1, 1, 0, 0, 1, 1, 1, 1, 1, 1]^{\mathsf{T}}$. We

TABLE VIII
THE SENSITIVITY ANALYSIS OF DIFFERENT HYPER-PARAMETERS. WE CONDUCT EXPERIMENTS WITH DIFFERENT PARAMETERS WITH BRAIN TUMOR CORE SEGMENTATION (5% TRAINING DATA) WITH BRATS USING DEEPMEDIC. RESULTS WHICH HAVE WORSE DSC THAN THE VANILLA BASELINE ARE HIGHLIGHTED WITH GRAY SHADING.

| Method | Parameter | | DSC | SEN | PRC | HD |
|---|---|---|---|---|---|---|
| w/ post-processing | | | | | | |
| Vanilla - CE | —— | | 50.4 | 41.0 | 83.5 | 18.0 |
| Asymmetric large margin loss | m = 0.2 | | 53.6 | 44.8 | 84.8 | 15.6 |
| | m = 0.5 | | 48.4 | 39.4 | 81.5 | 16.9 |
| | m = 1 | | 56.8 | 48.9 | 83.4 | 15.0 |
| | m = 1.5 | | 54.1 | 45.6 | 81.7 | 15.3 |
| | m = 2 | | 51.6 | 42.8 | 84.0 | 16.7 |
| | m = 3 | | 54.4 | 45.7 | 82.3 | 14.3 |
| Asymmetric focal loss | $\gamma$ = 0.5 | | 53.4 | 44.8 | 79.6 | 16.6 |
| | $\gamma$ = 1 | | 53.9 | 45.2 | 81.9 | 17.9 |
| | $\gamma$ = 1.5 | | 56.5 | 48.3 | 87.8 | 13.8 |
| | $\gamma$ = 2 | | 58.8 | 51.4 | 81.6 | 15.0 |
| | $\gamma$ = 3 | | 57.5 | 49.0 | 85.5 | 14.2 |
| | $\gamma$ = 4 | | 55.2 | 48.5 | 78.3 | 15.5 |
| Asymmetric adversarial training | $\epsilon$ = 1e-5 | $l$ = 2.5 | 50.3 | 41.8 | 82.0 | 17.2 |
| | $\epsilon$ = 1e-5 | $l$ = 5 | 58.1 | 50.0 | 84.7 | 14.1 |
| | $\epsilon$ = 1e-5 | $l$ = 10 | 58.5 | 50.8 | 80.1 | 16.2 |
| | $\epsilon$ = 1e-5 | $l$ = 15 | 53.8 | 46.2 | 80.1 | 16.9 |
| | $\epsilon$ = 1e-5 | $l$ = 20 | 56.6 | 50.7 | 76.9 | 18.8 |
| | $\epsilon$ = 1e-4 | $l$ = 10 | 57.6 | 51.1 | 78.9 | 16.1 |
| | $\epsilon$ = 1e-6 | $l$ = 10 | 56.2 | 48.5 | 81.1 | 17.8 |
| Asymmetric mixup | m = 0.1 | | 52.1 | 47.3 | 73.8 | 20.7 |
| | m = 0.15 | | 58.1 | 53.7 | 75.0 | 19.9 |
| | m = 0.2 | | 59.8 | 56.8 | 74.7 | 17.7 |
| | m = 0.25 | | 60.4 | 55.0 | 82.0 | 15.6 |
| | m = 0.3 | | 59.1 | 54.3 | 82.0 | 15.3 |
| | m = 0.4 | | 52.1 | 44.2 | 84.2 | 21.4 |
| | m = 0.8 | | 50.3 | 41.6 | 85.5 | 17.7 |
| w/o post-processing | | | | | | |
| Vanilla - CE | —— | | 51.0 | 42.6 | 78.6 | 17.5 |
| Asymmetric large margin loss | m = 0.2 | | 53.2 | 46.0 | 79.8 | 18.3 |
| | m = 0.5 | | 48.8 | 40.8 | 78.1 | 17.5 |
| | m = 1 | | 55.5 | 50.6 | 76.2 | 23.9 |
| | m = 1.5 | | 52.6 | 47.2 | 73.1 | 25.8 |
| | m = 2 | | 51.4 | 44.2 | 78.2 | 18.8 |
| | m = 3 | | 53.4 | 47.3 | 75.1 | 21.3 |
| Asymmetric focal loss | $\gamma$ = 0.5 | | 54.2 | 48.0 | 76.2 | 22.0 |
| | $\gamma$ = 1 | | 53.7 | 46.8 | 76.0 | 22.8 |
| | $\gamma$ = 1.5 | | 54.3 | 49.6 | 76.3 | 25.9 |
| | $\gamma$ = 2 | | 57.3 | 52.7 | 76.4 | 24.4 |
| | $\gamma$ = 3 | | 55.7 | 50.3 | 75.4 | 24.6 |
| | $\gamma$ = 4 | | 54.4 | 50.3 | 71.1 | 25.4 |
| Asymmetric adversarial training | $\epsilon$ = 1e-5 | $l$ = 2.5 | 50.5 | 43.6 | 76.3 | 21.3 |
| | $\epsilon$ = 1e-5 | $l$ = 5 | 56.6 | 51.3 | 76.1 | 21.9 |
| | $\epsilon$ = 1e-5 | $l$ = 10 | 56.8 | 51.8 | 74.8 | 23.6 |
| | $\epsilon$ = 1e-5 | $l$ = 15 | 53.3 | 47.6 | 74.8 | 21.5 |
| | $\epsilon$ = 1e-5 | $l$ = 20 | 55.4 | 53.2 | 72.0 | 26.2 |
| | $\epsilon$ = 1e-4 | $l$ = 10 | 56.9 | 53.3 | 74.1 | 22.4 |
| | $\epsilon$ = 1e-6 | $l$ = 10 | 55.2 | 50.0 | 76.0 | 23.8 |
| Asymmetric mixup | m = 0.1 | | 52.0 | 48.8 | 68.7 | 32.2 |
| | m = 0.15 | | 58.0 | 55.7 | 70.6 | 31.6 |
| | m = 0.2 | | 59.3 | 57.9 | 70.6 | 27.8 |
| | m = 0.25 | | 60.1 | 55.9 | 78.0 | 23.5 |
| | m = 0.3 | | 59.2 | 55.4 | 77.9 | 17.6 |
| | m = 0.4 | | 52.8 | 45.3 | 80.2 | 21.5 |
| | m = 0.8 | | 51.0 | 43.6 | 79.2 | 18.7 |

use $\gamma = 4$ in this experiments. We adopt post-processing described in Section V separately to the results of every classes. The results are shown in Table IX. The asymmetric focal loss can get better overall segmentation results than cross entropy or its symmetric variant. More importantly. it can get better segmentation results with higher sensitivity for most rare classes. Specifically, asymmetric focal loss can improve

Fig. 10. (a) The intensity histogram of BRATS, (b) ATLAS and (c) KiTS. The intensity of the foreground and background classes overlap a lot for ATLAS and KiTS. This can be a potential factor due to which the asymmetric mixup does not create useful synthetic samples and cannot improve the segmentation performance that much.

## XIII. QUANTITATIVE RESULTS WITHOUT POST-PROCESSING

The quantitative segmentation results without post-processing are summarized in Table X, Table XI and Table XII. Without post-processing, the proposed asymmetric regularization methods can improve DSC but could lead to worse distance-based evaluation metrics such as Hausdorff distance (HD). It is because the regularized model, which is more sensitive for the under-represented classes, would make relatively more false positive predictions. The false positive predictions which are far from the ground truth would increase HD significantly. However, in practice most false positive predictions could be easily removed by some connected component-based post-processing, as described in Section V. In this way, eventually we can get better or similar HD with our methods, as shown in the main text.

the average DSC of rare classes by 4.9%. We also notice that asymmetric focal loss would decrease the segmentation performance of esophagus which is taken as a rare class. It is because esophagus is too small, and post-processing would remove the correct segmentation regions by mistake but leave the false positive predictions. We think more advanced post-processing would help improve the segmentation in this case.



Fig. 11. The class distribution of the abdomen dataset we use in this study. We summarize the total pixel number of different classes. We take class 4, 5, 8, 9, 10, 11, 12 and 13 as rare classes.

TABLE IX

EVALUATION OF ABDOMEN SEGMENTATION WITH 25% OF TRAINING DATA WITH SYMMETRIC (SY.) AND ASYMMETRIC (ASY.) FOCAL LOSS. THE RARE CLASSES ARE MARKED WITH $r$. AVG IS THE AVERAGE PERFORMANCE OF ALL CLASSES. $\text{AVG}_r$ IS THE AVERAGE PERFORMANCE OF ALL RARE CLASSES. BEST RESULTS ARE HIGHLIGHTED IN BOLD.

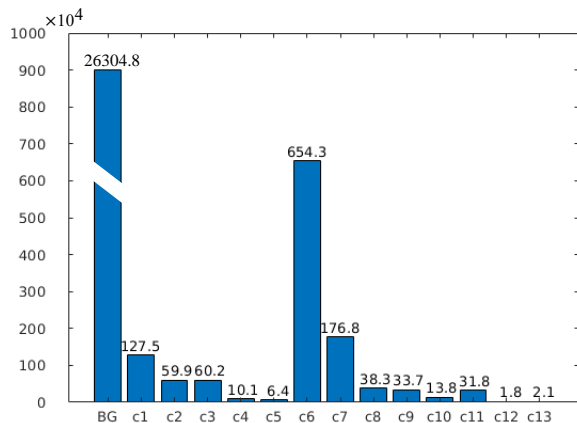| | c1 (spleen) | | | c2 (right kidney) | | | c3 (left kidney) | | | c4 (gallbladder) $r$ | | | c5 (esophagus) $r$ | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | vanilla - CE | sy. focal loss | asy. focal loss | vanilla - CE | sy. focal loss | asy. focal loss | vanilla - CE | sy. focal loss | asy. focal loss | vanilla - CE | sy. focal loss | asy. focal loss | vanilla - CE | sy. focal loss | asy. focal loss |
| DSC | **85.3** | 78.7 | 84.3 | 78.3 | **84.4** | 68.4 | 82.7 | **85.0** | 82.6 | 34.6 | 30.5 | **53.7** | **55.7** | 52.7 | 41.8 |
| Sensitivity | 80.0 | 70.5 | 76.9 | 73.7 | 77.2 | 62.9 | 77.5 | 79.5 | 77.8 | 25.8 | 22.9 | 46.3 | 50.7 | 48.2 | 43.8 |
| Precision | 93.6 | 96.0 | 95.9 | 84.6 | 94.8 | 76.2 | 94.7 | 94.4 | 93.0 | 78.4 | 68.5 | 68.9 | 74.4 | 69.5 | 45.0 |

| | c6 (liver) | | | c7 (stomach) | | | c8 (aorta) $r$ | | | c9 (vena cava) $r$ | | | c10 (vein) $r$ | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | vanilla - CE | sy. focal loss | asy. focal loss | vanilla - CE | sy. focal loss | asy. focal loss | vanilla - CE | sy. focal loss | asy. focal loss | vanilla - CE | sy. focal loss | asy. focal loss | vanilla - CE | sy. focal loss | asy. focal loss |
| DSC | 87.4 | 88.4 | **88.6** | 39.7 | 39.7 | **42.1** | 82.7 | 80.3 | **84.0** | 65.1 | 66.4 | **73.9** | 42.0 | 26.3 | **43.3** |
| Sensitivity | 84.1 | 85.2 | 84.0 | 28.8 | 28.5 | 31.0 | 76.6 | 73.1 | 82.9 | 56.8 | 60.3 | 76.9 | 28.4 | 16.5 | 31.0 |
| Precision | 92.3 | 92.8 | 94.2 | 91.3 | 84.1 | 86.3 | 91.6 | 91.5 | 86.3 | 86.1 | 79.3 | 72.7 | 91.5 | 82.1 | 80.0 |

| | c11 (pancreas) $r$ | | | c12 (right adrenal) $r$ | | | c13 (left adrenal) $r$ | | | AVG | | | $\text{AVG}_r$ | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | vanilla - CE | sy. focal loss | asy. focal loss | vanilla - CE | sy. focal loss | asy. focal loss | vanilla - CE | sy. focal loss | asy. focal loss | vanilla - CE | sy. focal loss | asy. focal loss | vanilla - CE | sy. focal loss | asy. focal loss |
| DSC | 17.2 | 24.3 | **26.4** | 54.3 | 32.8 | **55.9** | 34.8 | 28.9 | **47.0** | 58.5 | 55.3 | **60.9** | 48.3 | 42.8 | **53.2** |
| Sensitivity | 11.1 | 17.3 | 18.3 | 45.3 | 24.3 | 50.3 | 27.2 | 22.2 | 41.6 | 51.2 | 48.1 | 55.7 | 40.2 | 35.6 | 48.9 |
| Precision | 56.6 | 52.0 | 61.7 | 74.5 | 60.8 | 69.4 | 61.4 | 52.7 | 63.6 | 82.4 | 78.3 | 76.4 | 76.8 | 69.6 | 68.5 |

TABLE X

EVALUATION OF BRAIN TUMOR CORE SEGMENTATION USING DEEPMEDIC WITH DIFFERENT AMOUNTS OF TRAINING DATA AND DIFFERENT TECHNIQUES TO COUNTER OVERFITTING. THE RESULTS ARE CALCULATED WITHOUT ANY POST-PROCESSING. RESULTS WHICH HAVE WORSE DSC THAN THE VANILLA BASELINE ARE HIGHLIGHTED WITH GRAY SHADING. BEST AND SECOND BEST RESULTS ARE HIGHLIGHTED IN BOLD WITH THE BEST ALSO BEING UNDERLINED.

| Method | 5% training | | | | 10% training | | | | 20% training | | | | 50% training | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | DSC | SEN | PRC | HD | DSC | SEN | PRC | HD | DSC | SEN | PRC | HD | DSC | SEN | PRC | HD |
| Vanilla - CE [20] | 51.0 | 42.6 | 78.6 | 17.5 | 62.8 | 56.9 | 81.6 | **13.7** | 65.3 | 61.0 | 83.2 | 12.8 | 69.5 | 66.4 | 83.8 | 14.3 |
| Vanilla - CE - 80% tumor | 46.2 | 38.3 | 77.1 | 22.5 | 61.6 | 55.3 | 79.1 | 17.8 | 65.5 | 60.9 | 81.7 | 17.1 | 68.8 | 65.3 | 83.4 | 15.1 |
| Vanilla - F1 (DSC) | 47.7 | 38.7 | 82.0 | **15.4** | 59.4 | 52.2 | 82.7 | **13.7** | 64.6 | 59.0 | 82.9 | 13.6 | 67.4 | 63.4 | 84.7 | 13.4 |
| Vanilla - F2 [14] | 46.9 | 38.6 | 79.2 | **14.9** | 59.9 | 53.6 | 82.6 | 15.4 | 66.5 | 62.1 | 81.8 | **12.7** | 68.8 | 67.0 | 81.1 | 14.2 |
| Vanilla - F4 [14] | 51.6 | 44.0 | 78.6 | 19.8 | 60.1 | 54.1 | 81.7 | 16.5 | 65.9 | 63.1 | 80.9 | 19.1 | 67.8 | 65.7 | 83.1 | 12.2 |
| Vanilla - F8 [14] | 48.6 | 40.2 | 80.2 | 16.8 | 60.2 | 53.6 | 83.2 | 15.4 | 64.7 | 61.3 | 81.6 | 15.4 | 67.9 | 66.4 | 79.6 | 13.7 |
| Large margin loss [31] | 46.4 | 38.3 | 77.8 | 21.3 | 61.2 | 54.3 | 82.2 | 15.3 | 67.0 | 62.7 | 83.3 | **12.5** | 66.8 | 63.4 | 86.1 | **11.0** |
| Asymmetric large margin loss | 55.5 | 50.6 | 76.2 | 23.9 | 64.3 | 57.9 | 84.1 | **14.3** | 67.8 | 63.9 | 82.4 | 13.2 | 69.3 | 66.2 | 84.6 | 12.7 |
| Focal loss [29] | 53.6 | 46.3 | 78.7 | 20.3 | 62.9 | 56.0 | 82.5 | 17.3 | 65.2 | 61.1 | 82.6 | 19.2 | 67.2 | 63.2 | 84.7 | 15.3 |
| Asymmetric focal loss | 57.3 | 52.7 | 74.4 | 24.4 | 66.3 | 62.9 | 79.1 | 16.5 | 68.6 | 67.3 | 78.8 | 15.6 | **71.2** | 71.7 | 79.9 | 12.4 |
| Adversarial training [12] | 53.4 | 45.7 | 81.8 | 21.5 | 62.4 | 55.8 | 83.1 | 19.4 | 65.2 | 60.4 | 83.4 | 15.4 | 66.0 | 62.0 | 84.8 | 17.8 |
| Asymmetric adversarial training | 56.8 | 51.8 | 74.8 | 23.6 | 64.0 | 59.2 | 80.5 | 17.2 | 68.0 | 64.7 | 82.8 | 15.6 | 70.6 | 69.3 | 81.5 | 15.0 |
| Mixup [47] | 50.0 | 42.2 | 77.6 | 21.1 | 60.9 | 55.0 | 81.4 | 19.7 | 64.9 | 60.0 | 82.3 | 14.0 | 67.2 | 62.7 | 86.3 | 17.3 |
| Asymmetric mixup | **59.2** | 57.9 | 70.7 | 27.8 | **68.5** | 66.3 | 79.2 | 16.5 | **70.6** | 69.2 | 81.2 | 16.0 | 70.8 | 69.1 | 83.7 | **11.1** |
| Symmetric combination | 50.6 | 43.0 | 82.2 | 20.3 | 61.0 | 54.2 | 83.4 | 23.3 | 64.9 | 59.5 | 85.9 | 16.8 | 67.4 | 63.9 | 84.4 | 15.7 |
| Asymmetric combination | **62.4** | 64.7 | 71.5 | 27.8 | **71.4** | 73.8 | 74.3 | 20.8 | **71.9** | 74.1 | 79.2 | 20.7 | **72.9** | 77.0 | 77.8 | 19.0 |

TABLE XI

EVALUATION OF BRAIN STROKE LESION SEGMENTATION ON ATLAS USING DEEPMEDIC WITH DIFFERENT AMOUNTS OF TRAINING DATA AND DIFFERENT TECHNIQUES TO COUNTER OVERFITTING. THE RESULTS ARE CALCULATED WITHOUT POST-PROCESSING. RESULTS WHICH HAVE WORSE DSC THAN THE VANILLA BASELINE ARE HIGHLIGHTED WITH SHADING. BEST AND SECOND BEST RESULTS ARE IN BOLD WITH THE BEST ALSO UNDERLINED.

| Method | 30% training | | | | 50% training | | | | 100% training | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | DSC | SEN | PRC | HD | DSC | SEN | PRC | HD | DSC | SEN | PRC | HD |
| Vanilla - w/ augmentation [20] | 22.9 | 22.4 | 52.3 | 41.7 | 47.7 | 48.5 | 55.3 | **32.8** | 55.7 | 56.8 | 62.2 | 30.2 |
| Vanilla - w/o augmentation | 17.1 | 15.2 | 51.4 | 38.8 | 40.5 | 46.8 | 45.8 | 46.4 | 53.5 | 56.0 | 58.2 | 33.3 |
| Vanilla - asymmetric augmentation | 23.3 | 22.6 | 49.7 | 41.8 | 48.7 | 51.3 | 55.0 | 35.0 | **57.1** | 58.9 | 62.7 | **29.1** |
| Large margin loss [31] | 20.2 | 17.0 | 59.4 | **37.2** | 46.8 | 46.2 | 57.4 | 34.9 | 56.0 | 55.3 | 64.3 | **29.0** |
| Asymmetric large margin loss | 24.0 | 23.8 | 50.4 | 41.4 | 49.2 | 52.6 | 54.7 | 35.4 | 56.9 | 59.9 | 60.8 | 27.7 |
| Focal loss [29] | 21.9 | 20.2 | 55.0 | **37.7** | 47.8 | 49.3 | 55.6 | 33.7 | **57.1** | 59.6 | 63.4 | 32.6 |
| Asymmetric focal loss | 24.7 | 27.2 | 42.8 | 49.0 | 49.6 | 56.3 | 51.5 | 36.6 | 56.6 | 64.7 | 56.9 | 31.0 |
| Adversarial training [12] | 21.1 | 18.3 | 55.1 | 49.2 | 48.3 | 46.1 | 59.1 | 35.9 | 56.4 | 55.2 | 65.2 | 34.0 |
| Asymmetric adversarial training | **27.5** | 28.1 | 52.3 | 42.2 | **50.8** | 52.6 | 57.6 | 33.3 | 56.7 | 58.5 | 64.3 | 33.4 |
| Mixup [47] | 15.9 | 14.9 | 46.0 | 41.0 | 47.6 | 47.5 | 56.8 | **31.6** | 55.9 | 57.4 | 63.6 | 30.3 |
| Asymmetric mixup | 21.5 | 24.9 | 39.4 | 48.0 | 47.8 | 60.3 | 46.3 | 45.8 | 57.0 | 56.3 | 67.1 | 33.8 |
| Symmetric combination | 24.6 | 20.9 | 63.5 | 41.3 | 49.7 | 51.7 | 56.0 | 34.7 | 56.8 | 57.0 | 65.1 | 29.9 |
| Asymmetric combination | **29.9** | 34.2 | 47.1 | 47.4 | **51.1** | 58.6 | 52.2 | 38.6 | **57.9** | 62.4 | 61.5 | 32.1 |

TABLE XII

EVALUATION OF KIDNEY AND KIDNEY TUMOR SEGMENTATION BASED ON 3D U-NET WITH DIFFERENT AMOUNTS OF TRAINING DATA AND DIFFERENT TECHNIQUES TO COUNTER OVERFITTING. THE RESULTS ARE CALCULATED WITHOUT ANY POST-PROCESSING. RESULTS WHICH HAVE WORSE DSC THAN THE VANILLA BASELINE ARE HIGHLIGHTED WITH SHADING. BEST AND SECOND BEST RESULTS ARE IN BOLD WITH THE BEST ALSO UNDERLINED.

| Method | Kidney | | | | | | | | | | | |
| | 10% training | | | | 50% training | | | | 100% training | | | |
| | DSC | SEN | PRC | HD | DSC | SEN | PRC | HD | DSC | SEN | PRC | HD |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Vanilla - w/ augmentation [18] | 93.7 | 91.7 | 96.9 | 5.6 | 96.5 | 96.1 | 97.0 | 3.6 | **96.8** | 96.4 | 97.2 | **2.2** |
| Vanilla - w/o augmentation | 92.8 | 90.2 | 96.6 | 12.4 | 96.3 | 93.1 | 96.6 | 2.5 | 96.5 | 96.4 | 96.8 | 3.8 |
| Vanilla - asymmetric augmentation | 94.7 | 93.0 | 96.9 | 5.2 | 95.6 | 95.9 | 95.8 | 5.9 | 96.5 | 96.6 | 96.6 | 3.7 |
| Large margin loss [31] | **94.9** | 93.1 | 97.0 | 4.7 | 96.4 | 96.3 | 96.7 | 4.0 | 96.3 | 96.7 | 96.1 | 4.6 |
| Asymmetric large margin loss | 94.1 | 91.9 | 97.1 | 5.9 | 96.3 | 95.9 | 96.8 | 3.2 | **96.8** | 96.8 | 96.8 | 4.0 |
| Focal loss [29] | 91.6 | 86.1 | 99.1 | 6.6 | 94.1 | 89.9 | 99.0 | 5.5 | 94.4 | 90.2 | 99.1 | 4.1 |
| Asymmetric focal loss | 92.4 | 87.3 | 98.9 | 5.8 | 94.9 | 91.3 | 98.9 | 3.4 | 94.9 | 91.2 | 99.1 | 3.0 |
| Adversarial training [12] | 94.3 | 92.3 | 97.3 | 6.3 | 96.5 | 96.1 | 97.0 | **2.4** | **96.8** | 96.5 | 97.2 | **2.2** |
| Asymmetric adversarial training | 94.6 | 92.8 | 97.2 | 4.6 | **96.6** | 93.4 | 97.0 | 3.6 | **<u>97.0</u>** | 96.7 | 97.3 | **<u>2.1</u>** |
| Mixup [47] | **<u>95.2</u>** | 93.6 | 97.3 | **4.0** | **<u>96.9</u>** | 96.4 | 97.5 | **2.2** | **<u>97.0</u>** | 96.6 | 97.3 | 2.5 |
| Asymmetric mixup | 94.8 | 92.9 | 97.3 | **4.3** | 96.1 | 95.4 | 97.0 | 3.1 | 96.5 | 95.9 | 97.3 | 2.5 |
| Symmetric combination | 94.3 | 91.9 | 97.4 | 5.9 | 94.7 | 91.2 | 98.7 | 4.1 | **96.8** | 96.4 | 97.2 | **<u>2.1</u>** |
| Asymmetric combination | 94.0 | 90.5 | 98.4 | 4.8 | 94.3 | 90.8 | 95.4 | 5.1 | **96.8** | 95.9 | 97.7 | 3.3 |
| Method | Kidney tumor | | | | | | | | | | | |
| | 10% training | | | | 50% training | | | | 100% training | | | |
| | DSC | SEN | PRC | HD | DSC | SEN | PRC | HD | DSC | SEN | PRC | HD |
| Vanilla - w/ augmentation [18] | 54.9 | 47.9 | 77.5 | 93.8 | 75.6 | 73.7 | 83.9 | 49.8 | 79.3 | 78.3 | 84.8 | 48.3 |
| Vanilla - w/o augmentation | 37.8 | 32.9 | 62.8 | 121.2 | 64.1 | 62.1 | 74.6 | 85.6 | 71.3 | 69.7 | 79.3 | 54.9 |
| Vanilla - asymmetric augmentation | 56.1 | 50.4 | 74.5 | 97.1 | 75.3 | 73.6 | 84.3 | 60.3 | 79.6 | 79.5 | 85.2 | **35.0** |
| Large margin loss [31] | 54.8 | 48.2 | 77.2 | 84.0 | 77.1 | 75.2 | 84.5 | 58.8 | 80.9 | 82.1 | 83.5 | 47.4 |
| Asymmetric large margin loss | 55.7 | 50.1 | 75.4 | 99.6 | 77.9 | 76.0 | 84.9 | 54.6 | **81.9** | 82.3 | 84.0 | 56.2 |
| Focal loss [29] | 48.4 | 39.3 | 78.1 | 80.7 | 73.0 | 66.8 | 86.1 | 63.0 | 78.6 | 73.6 | 88.1 | 52.5 |
| Asymmetric focal loss | **57.0** | 49.9 | 74.9 | 95.9 | **78.2** | 76.6 | 84.6 | **43.8** | 80.8 | 81.1 | 83.5 | 48.9 |
| Adversarial training [12] | 51.6 | 45.0 | 79.5 | **78.3** | 73.6 | 71.6 | 83.2 | 55.3 | 81.4 | 81.6 | 84.0 | 52.2 |
| Asymmetric adversarial training | 56.9 | 51.1 | 79.4 | 87.5 | 77.4 | 75.6 | 85.5 | 58.3 | 81.8 | 81.5 | 86.1 | **<u>30.8</u>** |
| Mixup [47] | 54.5 | 48.3 | 79.6 | **<u>74.3</u>** | 77.1 | 73.8 | 86.2 | 48.3 | 80.6 | 79.5 | 84.9 | 52.0 |
| Asymmetric mixup | 55.1 | 48.6 | 79.9 | 92.3 | 77.9 | 74.4 | 87.9 | **40.8** | 79.8 | 79.0 | 85.9 | 54.9 |
| Symmetric combination | 54.2 | 47.1 | 79.0 | 105.4 | 73.6 | 67.5 | 86.0 | 51.5 | 80.5 | 80.0 | 84.5 | 48.4 |
| Asymmetric combination | **<u>59.2</u>** | 54.1 | 77.1 | 82.8 | **<u>79.4</u>** | 79.0 | 85.1 | 50.6 | **<u>82.2</u>** | 82.7 | 85.0 | 36.7 |