



blob loss: Instance Imbalance Aware Loss Functions for Semantic Segmentation

Florian Kofler^{1,2,3,7(✉)}, Suprosanna Shit^{1,2}, Ivan Ezhov^{1,2}, Lucas Fidon⁴,
Izabela Horvath^{1,5}, Rami Al-Maskari^{1,5}, Hongwei Bran Li^{1,13},
Harsharan Bhatia^{5,6}, Timo Loehr^{1,3}, Marie Piraud⁷, Ali Erturk^{5,6,8,9},
Jan Kirschke³, Jan C. Peeken^{10,11,12}, Tom Vercauteren⁴, Claus Zimmer³,
Benedikt Wiestler³, and Bjoern Menze^{1,13}

¹ Department of Informatics, Technical University Munich, Munich, Germany
florian.kofler@tum.de

² TranslaTUM - Central Institute for Translational Cancer Research, Technical University of Munich, Munich, Germany

³ Department of Diagnostic and Interventional Neuroradiology, School of Medicine, Klinikum rechts der Isar, Technical University of Munich, Munich, Germany

⁴ School of Biomedical Engineering and Imaging Sciences, King's College London, London, UK

⁵ Insitute for Tissue Engineering and Regenerative Medicine, Helmholtz Institute Munich (iTERM), Oberschleißheim, Germany

⁶ Institute for Stroke and Dementia research (ISD), University Hospital, LMU Munich, Munich, Germany

⁷ Helmholtz AI, Helmholtz Munich, Neuherberg, Germany

⁸ Graduate School of Neuroscience (GSN), Munich, Germany

⁹ Munich Cluster for Systems Neurology (Synergy), Munich, Germany

¹⁰ Department of Radiation Oncology, Klinikum rechts der Isar, Technical University of Munich, Munich, Germany

¹¹ Institute of Radiation Medicine (IRM), Department of Radiation Sciences (DRS), Helmholtz Zentrum, Munich, Germany

¹² Deutsches Konsortium für Translationale Krebsforschung (DKTK), Partner Site Munich, Munich, Germany

¹³ Department of Quantitative Biomedicine, University of Zurich, Zürich, Switzerland

Abstract. Deep convolutional neural networks (CNN) have proven to be remarkably effective in semantic segmentation tasks. Most popular loss functions were introduced targeting improved volumetric scores, such as the Dice coefficient (DSC). By design, DSC can tackle class imbalance, however, it does not recognize instance imbalance within a class. As a result, a large foreground instance can dominate minor instances and still produce a satisfactory DSC. Nevertheless, detecting tiny instances is crucial for many applications, such as disease monitoring. For example, it is imperative to locate and surveil small-scale lesions in the follow-up of

B. Wiestler and B. Menze—Equal contribution.

Supplementary Information The online version contains supplementary material available at https://doi.org/10.1007/978-3-031-34048-2_58.

multiple sclerosis patients. We propose a novel family of loss functions, *blob loss*, primarily aimed at maximizing instance-level detection metrics, such as *F1* score and *sensitivity*. *Blob loss* is designed for semantic segmentation problems where detecting multiple instances matters. We extensively evaluate a DSC-based *blob loss* in five complex 3D semantic segmentation tasks featuring pronounced instance heterogeneity in terms of texture and morphology. Compared to soft Dice loss, we achieve 5% improvement for MS lesions, 3% improvement for liver tumor, and an average 2% improvement for microscopy segmentation tasks considering *F1* score.

Keywords: semantic segmentation loss function · instance imbalance awareness · multiple sclerosis · lightsheet microscopy

1 Introduction

In recent years convolutional neural networks (CNN) have gained increasing popularity for complex machine learning tasks, such as *semantic segmentation*. In *semantic segmentation*, one segments object from different classes without differentiating multiple instances within a single class. In contrast, *instance segmentation* explicitly takes multiple instances into account, which involves simultaneous localization and segmentation. While U-net variants [23] still represent the state-of-the-art to address semantic segmentation, *Mask-RCNN* and its variants dominate *instance segmentation* [11]. The scarcity of training data often hinders the application of back-bone-dependent Mask RCNNs, while U-Nets have proven to be less data-hungry [5].

However, many semantic segmentation tasks feature relevant instance imbalance, where large instances dominate over smaller ones within a class, as illustrated in Fig. 1. Instances can vary not only with regard to size but also texture and other morphological features. U-nets trained with existing loss functions, such as Soft Dice [6, 18, 19, 24, 28], cannot address this. Instance imbalance is particularly pronounced and significant in medical applications: For example, even a single new multiple sclerosis (MS) lesion can impact the therapy decision. Despite many ways to compensate for class-imbalance [2, 9, 22, 28], there is a notable void in addressing instance imbalance in semantic segmentation settings. Additionally, established metrics have been shown to correlate insufficiently with expert assessment [16].

Contribution: We propose *blob loss*, a novel framework to equip semantic segmentation models with instance imbalance awareness. This is achieved by dedicating a specific loss term to each instance without the necessity of instance-wise prediction. *Blob loss* represents a method to convert any loss function into a novel instance imbalance aware loss function for semantic segmentation problems designed to optimize detection metrics. We evaluate its performance on five complex three-dimensional (3D) semantic segmentation tasks, for which the discovery of miniature structures matters. We demonstrate that extending soft

Dice loss to a *blob loss* improves detection performance in these multi-instance semantic segmentation tasks significantly. Furthermore, we also achieve volumetric improvements in some cases.

Related Work: Sirinukunwattana et al. [27] suggested an instance-based Dice metric for evaluating segmentation performance. Salehi et al. [24] were among the first to propose a loss function, called *Tversky loss*, for semantic segmentation of multiple sclerosis lesions in magnetic resonance imaging (MR), trying to improve detection metrics. Similarly, Zhu et al. [32] introduced Focal Loss, initially designed for object detection tasks [17], into medical semantic segmentation tasks.

There have been few recent attempts aiming for a solution to instance imbalance. Zhang et al. [30] propose an auxiliary lesion-level sphere prediction task. However, they do not explicitly consider each instance separately. Shirokikh et al. [25] propose an instance-weighted loss function where a global weight map is inversely proportional to the size of the instances. However, unlike size, not all types of imbalance, such as morphology or texture, can be quantified easily, limiting the method’s applicability.

2 Methods

First, we introduce the problem of instance imbalance in semantic segmentation tasks. Then we present our proposed *blob loss* functions.

Problem Statement: Large foreground areas dominate the calculation of established volumetric metrics (or losses); see Fig. 1. This is because the volumetric measures only accumulate true or false predictions on a voxel level but not at the instance level. Therefore, training models with volumetry-based loss functions, such as soft Dice loss (*dice*), often leads to unsatisfactory instance detection performance. To achieve a better instance detection performance, it is necessary to take instance imbalance into account. Instance imbalance can be of many categories, such as morphology and texture. Importantly, instance imbalance often cannot be easily specified and quantified for use in CNN training, for example, as instance weights in the loss function. Thus, using conventional methods, it is difficult to incorporate instance imbalance in CNN training. Our objective is to design loss functions to compensate for the instance imbalance while being agnostic to the instance imbalance type. Therefore, we aim to dissect the image domain in an instance-wise fashion:

blob loss Formulation: Consider a generic volumetric loss function \mathcal{L} and image domain Ω and foreground domain P . Formally our objective is to find an instance-specific subdomain $\Omega_n \subseteq \Omega$ corresponding to the n^{th} instance such that \mathcal{L} acting on Ω_n is aware of instance imbalance. The criteria to obtain these subsets $\{\Omega_n\}_{n=1}^N$ are such that $\Omega_i \cap \Omega_j \cap P = \phi; \forall(i, j), s.t. 1 \leq i, j \leq N, i \neq j$ and $\cup_{n=1}^N \Omega_n = \Omega$. In simple terms, the subsets $\{\Omega_n\}_{n=1}^N$ need to be mutually exclusive regarding foreground and collectively exhaustive with regard to the whole image domain.

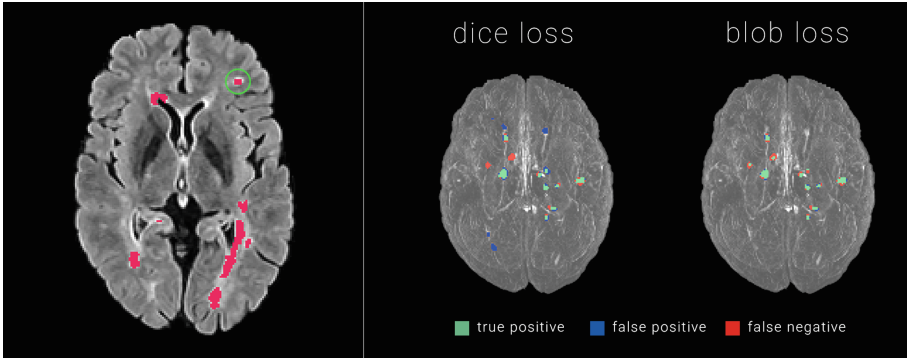


Fig. 1. *Problem statement (left):* The Dice coefficient (DSC) for the segmentation with vs. without a lesion, circled in green, is: 0.9806 . Therefore, the segmentations are hardly distinguishable in terms of DSC. However, from a clinical perspective, the difference is important as the detection of a single lesion can affect treatment decisions. *Comparison of segmentation performance (right):* Maximum intensity projections of the FLAIR images overlaid with segmentations for *dice* and *blob dice*. Lesions are colored according to their detection status: Green for *true positive*; Blue for *false positive*; Red for *false negative*. For this particular patient, applying the transformation to a *blob loss* improves $F1$ from 0.74 to 1.0 and the volumetric Dice coefficient from 0.56 to 0.70 and the latter is caused by an increase in *volumetric precision* from 0.48 to 0.75 , while the *volumetric sensitivity* remains constant at 0.66 .

To formalize *blob loss*, we address instance imbalance within a binary semantic segmentation framework. At the same time, we remain agnostic towards particular instance attributes and do not incorporate these in the loss function. To this extent, we propose to leverage the existing reference annotations and formally propose a novel family of instance-aware loss functions.

Consider a segmentation problem with N instances; for different input images, N can vary from few to many. Specifically, we propose to compute the instance-specific domain Ω_n by excluding all but the n^{th} foreground from the whole image domain Ω , see Eq. (1):

$$\Omega_n = \Omega \setminus \bigcup_{j=1, j \neq n}^N P_j \tag{1}$$

where P_j is the foreground domain for j^{th} instances of P . This masking process is illustrated by Fig. 2. It is worth noting that the background voxels are included in every Ω_n .

We propose to convert any loss function \mathcal{L} for binary semantic segmentation into an instance-aware loss function \mathcal{L}_{blob} defined as:

$$\mathcal{L}_{blob}((p_i)_{i \in \Omega}, (g_i)_{i \in \Omega}) = \frac{1}{N} \sum_{n=1}^N \mathcal{L}((p_i)_{i \in \Omega_n}, (g_i)_{i \in \Omega_n}) \tag{2}$$

where $\{g_i\}_{i \in \Omega}$ is the ground-truth segmentation, $\{p_i\}_{i \in \Omega}$ is the predicted segmentation, N is the number of instances in the foreground.

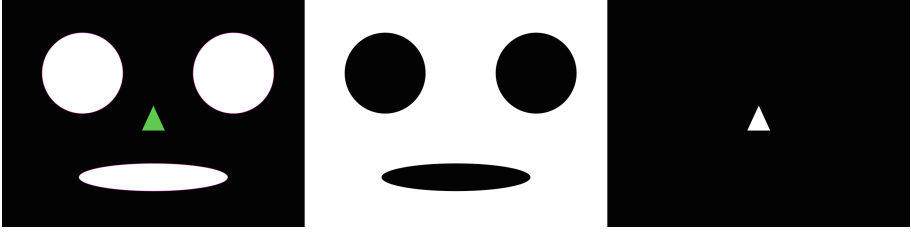


Fig. 2. Masking process described in Eq. (2). *Left:* the global ground truth label (GT), with the n^{th} instance highlighted in green. *Middle:* The loss mask Ω_n for the n^{th} instance (MASK) for multiplication with the network outputs. *Right:* the label used for the computation of the local *blob loss* for the n^{th} instance. This process is repeated for every instance.

As our goal is to assign equal importance to all instances irrespective of their size, shape, texture, and other topological attributes, we average over all instances.

To compute the total loss for a volume, we combine the instance-wise Loss component from Eq. (2) with a global component to obtain the final Loss:

$$\mathcal{L}_{total} = \alpha \mathcal{L}_{global} + \beta \mathcal{L}_{blob} \quad (3)$$

where α and β denote the weights for the global and instance constraint \mathcal{L}_{blob} . We (anonymously) provide a sample Pytorch implementation of a *dice*-based *blob loss* on [GitHub](#). In order to accelerate our training, we precompute the instances, here defined as connected components using *cc3d* [26], version 3.2.1.

Model Training: For all our experiments, we use a basic 3D U-Net implemented via [MONAI](#) inspired by [8] and further depicted in supplementary materials. Furthermore, we use a dropout ratio of 0.1 and employ *mish* as activation function [20]. Otherwise, we stick to the default parameters of the U-Net implementation.

Loss Functions for Comparison: As baselines we use the MONAI implementations of soft Dice loss (*dice*) and Tversky loss (*tversky*) [24]. For *tversky*, we always use the standard parameters of $\alpha = 0.3$ and $\beta = 0.7$ suggested by the authors in the original publication [24]. For comparison we create *blob dice*, by transforming the standard *dice* into a *blob loss* using our conversion method Eq. (2). The final loss is obtained by employing *dice* in the \mathcal{L}_{global} and \mathcal{L}_{blob} terms of the proposed total loss Eq. (3). In analog fashion, we derive *blob tversky*. Furthermore, we compare against *inverse weighting (iw)*, the globally weighted loss function of Shirokikh et al. [25]. For this, we use the official [GitHub implementation](#) to compute the weight maps and loss and deploy these in our training pipelines.

Training Procedure: Our CNNs are trained on multiple cuboid-shaped crops per batch element, with higher resolution in the axial dimension, enabling the learning of contextual image features. The crops are randomly sampled around a center voxel that consists of *foreground* with a 95% probability. We consider one

epoch as one full iteration of forward and backward passes through all batches of the training set. For all training, *Ranger21* [29] serves as our optimizer. For each experiment, we keep the initial learning rate (lr) constant between training runs. Depending on the segmentation task, we deploy varying suitable image normalization strategies. For comparability, we keep all training parameters except for the loss functions constant on a segmentation task basis and stick to this standard training procedure.

Training-Test Split and Model Selection: Given the high heterogeneity of our bio-medical datasets and the limited availability of high-quality ground truth annotations due to the very costly labeling procedures requiring domain experts, we do not set aside data for validation and therefore do not conduct model selection. Instead, inspired by [13], we split our data *80:20* into training and test set and evaluate on the last checkpoint of the model training. As an exception, the MS dataset comes with predefined training, validation, and test set splits; therefore, we additionally evaluate the *best* model checkpoint, meaning the model with the lowest loss on the validation set. As we are more interested in *blob loss*' generalization capabilities than exact quantification of improvements on particular datasets, we prioritize a broad validation on multiple datasets over cross-validation.

Technical Details: Our experiments were conducted using NVIDIA RTX8000, RTX6000, RTX3090, and A6000 GPUs using CUDA version *11.4* in conjunction with Pytorch version *1.9.1* and MONAI version *0.7.0*.

2.1 Evaluation Metrics and Interpretation

Metrics: We obtain global, volumetric performance measures from *pymia* [14]. In addition to DSC, we also evaluate *volumetric sensitivity* (S), *volumetric precision* (P), and the *Surface Dice similarity coefficient* ($SDSC$). To compute instance-wise detection metrics, namely instance F1 ($F1$), *instance sensitivity* (IS) and *instance precision* (IP), we employ a proven evaluation pipeline from Pan et al. [21].

Interpretation: By design, human annotators tend to overlook tiny structures. For comparison, human annotators initially missed *29%* of micrometastases when labeling the DeepMACT light-sheet microscopy dataset [21]. Therefore, the likelihood of a structure being correctly labeled in the ground truth is much higher for foreground than for background structures. Additionally, human annotators have a tendency to label a structure's center but do not perfectly trace its contours. Both phenomena are illustrated in Fig. 3. These effects are particularly pronounced for microscopy datasets, which often feature thousands of blobs. These factors are important to keep in mind when interpreting the results. Consequently, volumetric - and instance sensitivity are much more informative than volumetric and instance precision.

3 Experiments

To validate *blob loss*, we train segmentation models on a selection of datasets from different 3D imaging modalities, namely brain MR, thorax CT, and light-sheet microscopy. We select datasets featuring a variety of fragmented semantic segmentation problems. For simplicity, we use the default values $\alpha = 2.0$ and $\beta = 1.0$ across all experiments.

Multiple Sclerosis (MRI): The Multiple Sclerosis (MS) dataset, comprising 521 single timepoint MRI examinations of patients with MS, was collected for internal validation of MS lesion segmentation algorithms. The patients come from a representative, institutional cohort covering all stages (in terms of time from disease onset) and forms (relapsing-remitting, progressive) of MS. A 3D T1w and a 3D FLAIR sequence were acquired on a 3 T *Philips Achieva* scanner. All 3D volumes feature $193 \times 193 \times 229$ voxels in 1mm isotropic resolution. The dataset divides into a fixed training set of 200, a validation set of 21, and a test set of 200 cases. The annotations feature a total of 4791 blobs, with 25.69 ± 23.01 blobs per sample. Expert neuroradiologists annotated the MS lesions manually and ensured pristine ground truth quality with consensus voting.

For all training runs of 500 epochs, we set the initial learning rate to $1e-2$ and the batch size to 4. The networks are trained on a single GPU using 2 random crops with a patch size of $192 \times 192 \times 32$ voxels per batch element after applying a *min/max* normalization. As the MS dataset comes with a predefined validation set of 21 images, we also save the checkpoint with the lowest loss on the validation set and compare it to the respective last checkpoint of the training. In addition to the standard *dice*, we also compare against *tversky*. Furthermore, we conduct an ablation study to find out how the performance metrics are affected by choosing different values for α and β .

Liver Tumors - LiTS (CT): To develop an understanding of *blob loss* performance on other imaging modalities, we train a model for segmenting liver tumors on CT images of the *LiTS* challenge [4]. The dataset consists of varying high-resolution CT images of the abdomen. The challenge’s original task was segmenting liver and liver tumor tissue. As we are primarily interested in segmenting small fragmented structures, we limit our experiments to the liver area and segment only liver tumor tissue (in contrast to tumors, the liver represents a huge solid structure, and we are interested in blobs). We split the publicly available training set into 104 images for training and 27 for testing. The annotations were created by expert radiologists and feature a total of 908 blobs, with 12.39 ± 14.92 blobs per sample.

For all training runs of 500 epochs, we set the initial learning rate to $1e-2$ and the batch size to 2. The networks are trained on two GPUs in parallel using 2 random crops with a patch size of $192 \times 192 \times 64$ voxels per batch element. We apply normalization based on windowing on the Hounsfield (HU) scale. Therefore, we define a normalization window suitable for liver tumor segmentation around center $30 HU$ with a width of $150 HU$, and 20% added tolerance.

DISCO-MS (Light-Sheet Microscopy). To develop an understanding for *blob loss* performance on other imaging modalities, we train a model for segmenting Amyloid plaques in light-sheet microscopy images of the *DISCO-MS* dataset [3].

The volumes of $300 \times 300 \times 300$ voxels resolution contain cleared tissue of mouse brain. We split the publicly available dataset into 41 volumes for training and six for testing. The annotations feature a total of 988 blobs, with 28.32 ± 24.44 blobs per sample. Even though the label quality is very high, the results should still be interpreted with care following the guidelines in Sect. 2.1.

For all training runs of 800 epochs, we set the initial learning rate to $1e-3$ and the batch size to 6. As our initial model trained with *dice* does not produce satisfactory results, we furthermore try learning rates of $1e-2$, $3e-4$ and $1e-4$, following the heuristics suggested by [1] without success. The networks are trained on two GPUs in parallel using 2 random crops with a patch size of $192 \times 192 \times 64$ per batch element. The images are globally normalized, using a minimum and maximum threshold defined by the 0.5 and 99.5 percentile.

SHANEL (Light-Sheet Microscopy). For further validation, we evaluate neuron segmentation in light-sheet microscopy images of the *SHANEL* dataset [31]. The volumes of $200 \times 200 \times 200$ voxels resolution contain cleared human brain tissue from the primary visual cortex, the primary sensory cortex, the primary motor cortex, and the hippocampus. We split this publicly available dataset into nine volumes for training and three for testing. The annotations feature a total of 20684 blobs, with 992.14 ± 689.39 blobs per sample. As the data is more sparsely annotated than DISCO-MS, F1 and especially DSC should be interpreted with great care, as described in Sect. 2.1.

For all training runs of 1000 epochs, we set the initial learning rate to $1e-3$ and the batch size to 3. The networks are trained on two GPUs in parallel using 6 random crops with a patch size of $128 \times 128 \times 32$ per batch element, with min/max normalization.

DeepMACT (Light-Sheet Microscopy). For further validation, we evaluate the segmentation of micrometastasis in light-sheet microscopy images of

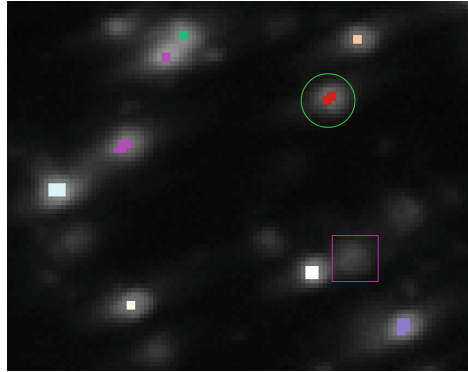


Fig. 3. Zoomed in 2D view on a volume of the SHANEL [31] dataset. The overlaid labels are colored according to a 3D connected component analysis. The expert biologists did not label each foreground object in every slice, e.g., the magenta-colored square. Furthermore, the contours of the structures are imperfectly segmented, for instance, the red label within the bright green circle. These effects can partially be attributed to the ambiguity of the light-sheet microscopy signal [15]. However, they are also observed in the human annotations of the MS and LiTS dataset. (Color figure online)

the *DeepMact* dataset [21]. The volumes of $350 \times 350 \times 350$ resolution contain cleared tissue featuring different body parts of a mouse. We split the publicly available dataset into 115 images for training and 19 for testing. The annotations feature a total of 484 blobs, with 6.99 ± 8.14 blobs per sample. As the data is sparsely annotated, *F1* and especially *DSC* should be interpreted with great care, as described in Sect. 2.1.

For all training runs of 500 epochs, we set the initial learning rate to $1e-2$ and the batch size to 4. The networks are trained on a single GPU using 2 random crops with a patch size of $192 \times 192 \times 48$. The images are globally normalized based using a minimum and maximum threshold defined by the 0.0 and 99.5 percentile.

4 Results

Table 1 summarizes the results of our experiments. Across all datasets, we find that extending *dice* to a *blob loss* helps to improve detection metrics. Furthermore, in some cases, we also observe improvements in volumetric performance measures. While model selection seems not beneficial on this dataset, employing *blob loss* produces more robust results, as both the *dice* and *tversky* models suffer performance drops for the *best* checkpoints. Notably, even though *tversky* was explicitly proposed for MS lesion segmentation, it is clearly outperformed by *dice*, as well as *blob dice* and *blob tversky*. Further, even with the mitigation strategies suggested by the authors, *inverse weighting* produced over-segmentations.

Table 2 summarizes the results of the ablation study on α and β parameters of *blob loss*. We find that assigning higher importance to the global parameter by choosing $\alpha = 2$ and $\beta = 1$ seems to produce the best results. Overall, we find that *blob loss* seems quite robust regarding the choice of hyperparameters as long as the global term remains included by choosing a α greater than 0.

5 Discussion

Contribution: *blob loss* can be employed to provide existing loss functions with instance imbalance awareness. We demonstrate that the application of *blob loss* improves detection- and in some cases, even volumetric segmentation performance across a diverse set of complex 3D bio-medical segmentation tasks. We evaluate *blob loss*' performance in the segmentation of multiple sclerosis (MS) lesions in MR, liver tumors in CT, and segmentation of different biological structures in 3D light-sheet microscopy datasets. Depending on the dataset, it achieves these improvements either due to better detection of foreground objects, better suppression of background objects, or both. We provide an implementation of blob loss leveraging on a precomputed connected component analysis for fast processing times.

Limitations: Certainly, the biggest disadvantage of *blob loss* is the dependency on instance segmentation labels; however, in many cases, these can be simply

Table 1. Experimental results for five datasets. For all training runs with *blob loss* we use $\alpha = 1$ and $\beta = 2$. Note that the results for LiTS are based on a different, more challenging test set and are therefore not comparable with the public leaderboard of the LiTS challenge. For DISCO-MS, the *dice* model completely over-segments and produces dissatisfactory results. Therefore, we try two additional training runs with reduced learning rates following the heuristics suggested by [1], resulting in similar over-segmentation. The same problem is observed for *inverse weighting (iw)*. Shirokikh et al. [25] themselves note the stability problems of the method and suggest lowering the learning rate to $1 - \epsilon_3$.

| dataset | loss | lr | DSC | SDSC | F1 | IS | IP |
|----------|--------------|------|--------------|--------------|--------------|--------------|--------------|
| MS | blob dice | 1e-2 | 0.680 | 0.848 | 0.810 | 0.822 | 0.828 |
| | dice | 1e-2 | 0.660 | 0.820 | 0.758 | 0.854 | 0.711 |
| | iw [25] | 1e-2 | 0.153 | 0.167 | 0.278 | 0.801 | 0.188 |
| | iw [25] | 1e-3 | 0.243 | 0.273 | 0.282 | 0.819 | 0.189 |
| | blob tversky | 1e-2 | 0.690 | 0.852 | 0.804 | 0.829 | 0.804 |
| | tversky | 1e-2 | 0.601 | 0.697 | 0.566 | 0.854 | 0.459 |
| LiTS | blob dice | 1e-2 | 0.663 | 0.542 | 0.657 | 0.861 | 0.631 |
| | dice | 1e-2 | 0.660 | 0.546 | 0.623 | 0.801 | 0.599 |
| SHANEL | blob dice | 1e-3 | 0.543 | 0.808 | 0.792 | 0.874 | 0.724 |
| | dice | 1e-3 | 0.539 | 0.794 | 0.783 | 0.854 | 0.723 |
| DISCO-MS | blob dice | 1e-3 | 0.546 | 0.678 | 0.589 | 0.760 | 0.481 |
| | dice | 1e-3 | 0.095 | 0.083 | 0.012 | 0.870 | 0.006 |
| | dice | 3e-4 | 0.016 | 0.036 | 0.379 | 0.896 | 0.240 |
| | dice | 1e-4 | 0.007 | 0.011 | 0.228 | 0.825 | 0.132 |
| DeepMACT | blob dice | 1e-2 | 0.357 | 0.393 | 0.391 | 0.871 | 0.276 |
| | dice | 1e-2 | 0.353 | 0.372 | 0.367 | 0.801 | 0.254 |

Table 2. Ablation analysis on the *blob loss*' hyperparameters α and β for the MS lesions dataset. We observe that *blob loss* seems to be quite robust with regard to hyperparameter choice, as long as the global term remains present, compare Eq. (3). The default parameters $\alpha = 2$ and $\beta = 1$ provide the best results.

| loss | α | β | DSC | S | P | SDSC | F1 | IS | IP |
|-----------|----------|---------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| blob dice | 3 | 1 | 0.674 | 0.629 | 0.765 | 0.833 | 0.790 | 0.796 | 0.815 |
| blob dice | 2 | 1 | 0.680 | 0.626 | 0.782 | 0.848 | 0.810 | 0.822 | 0.828 |
| blob dice | 1 | 1 | 0.658 | 0.580 | 0.802 | 0.839 | 0.804 | 0.840 | 0.801 |
| blob dice | 1 | 2 | 0.630 | 0.552 | 0.803 | 0.819 | 0.792 | 0.832 | 0.786 |
| dice | 1 | 0 | 0.660 | 0.704 | 0.656 | 0.820 | 0.758 | 0.854 | 0.711 |
| blob | 0 | 1 | 0.522 | 0.409 | 0.837 | 0.728 | 0.744 | 0.805 | 0.727 |

obtained by a connected component analysis, as demonstrated in our experiments. Another disadvantage of *blob loss* compared to other loss functions are the more extensive computational requirements. By definition, the user is required to run computations with large patch sizes that feature multiple instances. This results in an increased demand for GPU memory, especially when working with 3D data (as in our experiments). However, larger patch sizes have proven helpful for bio-medical segmentation problems, in general, [12]. Furthermore, according to our formulation, *blob loss* possesses an interesting mathematical property, it penalizes false positives proportionally to the number of instances in the volume. Additionally, even though *blob loss* can easily be reduced to a single hyperparameter, and it proved quite robust in our experiments, it might be sensitive to hyperparameter tuning. Moreover, by design *blob loss* can only improve performance for multi-instance segmentation problems.

Interpretation: One can only speculate why *blob loss* improves performance metrics. CNNs learn features that are very sensitive to texture [10]. Unlike conventional loss functions, *blob loss* adds attention to every single instance in the volume. Thus the network is forced to learn the instance imbalanced features such as, but not limited to morphology and texture, which would not be well represented by optimizing via *dice* and alike. Such instance imbalance was observed in the medical field, as it has been shown that MS lesions change their imaging phenotype over time, with recent lesions looking significantly different from older ones [7]. These aspects might explain the gains in instance sensitivity. Furthermore, adding the multiple instance terms leads to heavy penalization on background, which might explain why we often observe an improvement in precision, see supplementary materials.

Outlook: Future research will have to reveal to which extent transformation to *blob loss* can be beneficial for other segmentation tasks and loss functions. A first and third place in recent public segmentation challenges using a compound-based variant *blob loss* indicate that *blob loss* might possess broad applicability towards other instance imbalanced semantic segmentation problems.

References

1. Bengio, Y.: Practical recommendations for gradient-based training of deep architectures. In: Montavon, G., Orr, G.B., Müller, K.-R. (eds.) *Neural Networks: Tricks of the Trade*. LNCS, vol. 7700, pp. 437–478. Springer, Heidelberg (2012). https://doi.org/10.1007/978-3-642-35289-8_26
2. Berman, M., et al.: The lovász-softmax loss: a tractable surrogate for the optimization of the intersection-over-union measure in neural networks. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4413–4421 (2018)
3. Bhatia, et al.: Proteomics of spatially identified tissues in whole organs. *arXiv* (2021)
4. Bilic, P., et al.: The liver tumor segmentation benchmark (LiTS) (2019)
5. Caicedo, J.C., et al.: Nucleus segmentation across imaging experiments: the 2018 data science bowl. *Nat. Methods* **16**(12), 1247–1253 (2019)

6. Eelbode, T., et al.: Optimization for medical image segmentation: theory and practice when evaluating with dice score or Jaccard index. *IEEE Trans. Med. Imaging* **39**(11), 3679–3690 (2020)
7. Elliott, C., et al.: Slowly expanding/evolving lesions as a magnetic resonance imaging marker of chronic active multiple sclerosis lesions. *Mult. Scler. J.* **25**(14), 1915–1925 (2019)
8. Falk, T., et al.: U-net: deep learning for cell counting, detection, and morphometry. *Nat. Methods* **16**(1), 67–70 (2019)
9. Fidon, L., et al.: Generalised wasserstein dice score for imbalanced multi-class segmentation using holistic convolutional networks. In: Crimi, A., Bakas, S., Kuijf, H., Menze, B., Reyes, M. (eds.) *BrainLes 2017*. LNCS, vol. 10670, pp. 64–76. Springer, Cham (2018). https://doi.org/10.1007/978-3-319-75238-9_6
10. Geirhos, R., et al.: ImageNet-trained CNNs are biased towards texture; increasing shape bias improves accuracy and robustness. *arXiv preprint arXiv:1811.12231* (2018)
11. He, K., et al.: Mask R-CNN. In: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 2961–2969 (2017)
12. Isensee, F., Jaeger, P.F., Kohl, S.A., Petersen, J., Maier-Hein, K.H.: nnU-net: a self-configuring method for deep learning-based biomedical image segmentation. *Nat. Methods* **18**(2), 203–211 (2021)
13. Isensee, F., et al.: nnU-net: breaking the spell on successful medical image segmentation. *arXiv preprint arXiv:1904.08128*, vol. 1, pp. 1–8 (2019)
14. Jungo, A., et al.: pymia: a python package for data handling and evaluation in deep learning-based medical image analysis. *Comput. Methods Programs Biomed.* **198**, 105796 (2021)
15. Kofler, F., et al.: Approaching peak ground truth. *arXiv preprint arXiv:2301.00243* (2022)
16. Kofler, F., et al.: Are we using appropriate segmentation metrics? Identifying correlates of human expert perception for CNN training beyond rolling the dice coefficient (2021)
17. Lin, T.Y., et al.: Focal loss for dense object detection. In: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 2980–2988 (2017)
18. Ma, J., et al.: Loss odyssey in medical image segmentation. *Med. Image Anal.* **71**, 102035 (2021)
19. Milletari, F., Navab, N., Ahmadi, S.A.: V-net: fully convolutional neural networks for volumetric medical image segmentation. In: *2016 Fourth International Conference on 3D Vision (3DV)*, pp. 565–571, IEEE (2016)
20. Misra, D.: Mish: a self regularized non-monotonic neural activation function. *arXiv preprint arXiv:1908.08681* (2019)
21. Pan, C., et al.: Deep learning reveals cancer metastasis and therapeutic antibody targeting in the entire body. *Cell* **179**(7), 1661–1676 (2019)
22. Rahman, M.A., Wang, Y.: Optimizing intersection-over-union in deep neural networks for image segmentation. In: Bebis, G., et al. (eds.) *ISVC 2016*. LNCS, vol. 10072, pp. 234–244. Springer, Cham (2016). https://doi.org/10.1007/978-3-319-50835-1_22
23. Ronneberger, O., Fischer, P., Brox, T.: U-net: convolutional networks for biomedical image segmentation. In: Navab, N., Hornegger, J., Wells, W.M., Frangi, A.F. (eds.) *MICCAI 2015*. LNCS, vol. 9351, pp. 234–241. Springer, Cham (2015). https://doi.org/10.1007/978-3-319-24574-4_28

24. Salehi, S.S.M., Erdogmus, D., Gholipour, A.: Tversky loss function for image segmentation using 3D fully convolutional deep networks. In: Wang, Q., Shi, Y., Suk, H.-I., Suzuki, K. (eds.) MLMI 2017. LNCS, vol. 10541, pp. 379–387. Springer, Cham (2017). https://doi.org/10.1007/978-3-319-67389-9_44
25. Shirokikh, B., et al.: Universal loss reweighting to balance lesion size inequality in 3D medical image segmentation. In: Martel, A.L., et al. (eds.) MICCAI 2020. LNCS, vol. 12264, pp. 523–532. Springer, Cham (2020). https://doi.org/10.1007/978-3-030-59719-1_51
26. Silversmith, W.: seung-lab/connected-components-3d: Zenodo release v1. Zenodo (2021). <https://doi.org/10.5281/zenodo.5535251>
27. Sirinukunwattana, K., Snead, D.R., Rajpoot, N.M.: A stochastic polygons model for glandular structures in colon histology images. *IEEE Trans. Med. Imaging* **34**(11), 2366–2378 (2015)
28. Sudre, C.H., Li, W., Vercauteren, T., Ourselin, S., Jorge Cardoso, M.: Generalised dice overlap as a deep learning loss function for highly unbalanced segmentations. In: Cardoso, M.J., et al. (eds.) DLMIA/ML-CDS -2017. LNCS, vol. 10553, pp. 240–248. Springer, Cham (2017). https://doi.org/10.1007/978-3-319-67558-9_28
29. Wright, L., Demeure, N.: Ranger21: a synergistic deep learning optimizer. arXiv preprint [arXiv:2106.13731](https://arxiv.org/abs/2106.13731) (2021)
30. Zhang, H., et al.: All-net: Anatomical information lesion-wise loss function integrated into neural network for multiple sclerosis lesion segmentation. *NeuroImage: Clin.* **32**, 102854 (2021)
31. Zhao, S., et al.: Cellular and molecular probing of intact human organs. *Cell* **180**(4), 796–812 (2020)
32. Zhu, W., et al.: AnatomyNet: deep learning for fast and fully automated whole-volume segmentation of head and neck anatomy. *Med. Phys.* **46**(2), 576–589 (2019)