# Generalised Wasserstein Dice Score for Imbalanced Multi-class Segmentation using Holistic Convolutional Networks

Lucas Fidon[1], Wenqi Li[1], Luis C. Garcia-Peraza-Herrera[1],
Jinendra Ekanayake[2,3], Neil Kitchen[2],
Sébastien Ourselin[1,3], and Tom Vercauteren[1,3]
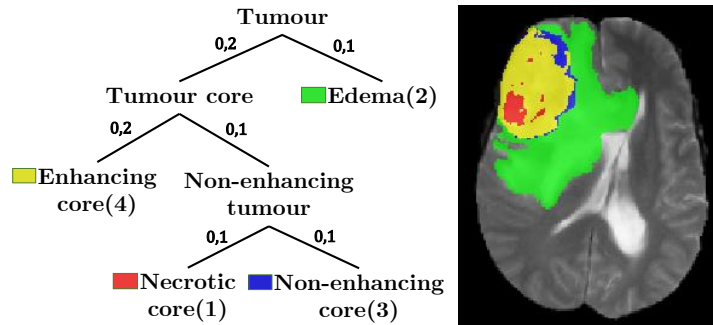
[1] TIG, CMIC, University College London, London, UK
[2] NHNN, University College London Hospitals, London UK
[3] Wellcome / EPSRC Centre for Interventional and Surgical Sciences,
UCL, London, UK

**Abstract.** The Dice score is widely used for binary segmentation due to its robustness to class imbalance. Soft generalisations of the Dice score allow it to be used as a loss function for training convolutional neural networks (CNN). Although CNNs trained using mean-class Dice score achieve state-of-the-art results on multi-class segmentation, this loss function does neither take advantage of inter-class relationships nor multi-scale information. We argue that an improved loss function should balance misclassifications to favour predictions that are semantically meaningful. This paper investigates these issues in the context of multi-class brain tumour segmentation. Our contribution is threefold. 1) We propose a semantically-informed generalisation of the Dice score for multi-class segmentation based on the Wasserstein distance on the probabilistic label space. 2) We propose a holistic CNN that embeds spatial information at multiple scales with deep supervision. 3) We show that the joint use of holistic CNNs and generalised Wasserstein Dice score achieves segmentations that are more semantically meaningful for brain tumour segmentation.

## 1 Introduction

Automatic brain tumour segmentation is an active research area. Learning-based methods using convolutional neural networks (CNNs) have recently emerged as the state of the art [9,11]. One of the challenges is the severe class imbalance. Two complementary ways have traditionally been used when training CNNs to tackle imbalance: 1) using a sampling strategy that imposes constraints on the selection of image patches; and 2) using pixel-wise weighting to balance the contribution of each class in the objective function. For CNN-based segmentation, samples should ideally be entire subject volumes to support the use of fully convolutional network and maximise the computational efficiency of convolution operations within GPUs. As a result, weighted loss functions appear more promising to improve CNN-based automatic brain tumour segmentation. Using

**Fig. 1.** Left: tree on BraTS label space. Edge weights have been manually selected to reflect the distance between labels. Right: illustration on a T2 scan from BraTS'15 [14].

soft generalisations of the Dice score (a popular overlap measure for binary segmentation) directly as a loss function has recently been proposed [15,18]. By introducing global spatial information into the loss function, the Dice loss has been shown to be more robust to class imbalance. However at least two sources of information are not fully utilised in this formulation: 1) the structure of the label space; and 2) the spatial information across scales. Considering the class imbalance and the hierarchical label structure illustrated in Fig.1, both of them are likely to play an important role for multi-class brain tumour segmentation.

In this paper, we propose two complementary contributions that leverage prior knowledge about brain tumour structure. First, we exploit the Wasserstein distance [7,17], which can naturally embed semantic relationships between classes for the comparison of label probability vectors, to generalise the Dice score for multi-class segmentation. Second, we propose a new holistic CNN architecture inspired by [8,19] that embeds spatial information at different scales and introduces deep supervision during the CNN training. We show that the combination of the proposed generalised Wasserstein Dice score and our Holistic CNN achieves better generalisation compared to both mean soft Dice score training and classic CNN architectures for multi-class brain tumour segmentation.

## 2    A Wasserstein approach for multi-class soft Dice score

### 2.1    Dice score for crisp binary segmentation

The Dice score is a widely used overlap measure for pairwise comparison of binary segmentations $S$ and $G$. It can be expressed both in terms of set operations or statistical measures as:

$$D = \frac{2|S \cap G|}{|S| + |G|} = \frac{2\Theta_{TP}}{2\Theta_{TP} + \Theta_{FP} + \Theta_{FN}} = \frac{2\Theta_{TP}}{2\Theta_{TP} + \Theta_{AE}} \tag{1}$$

with $\Theta_{TP}$ the number of true positives, $\Theta_{FP}/\Theta_{FN}$ the number of false positives/false negatives, and $\Theta_{AE} = \Theta_{FP} + \Theta_{FN}$ the number of all errors.

## 2.2   Dice score for soft binary segmentation

Extensions to soft binary segmentations [1,2] rely on the concept of disagreement for pairs of probabilistic classifications. The classes $S_i$ and $G_i$ of each voxel $i \in \mathbf{X}$ can be defined as random variables on the label space $\mathbf{L} = \{0, 1\}$ and the probabilistic segmentations can be represented as label probability maps: $p = \{p^i := P(S_i = 1)\}_{i \in \mathbf{X}}$ and $g = \{g^i := P(G_i = 1)\}_{i \in \mathbf{X}}$. We denote $P(\mathbf{L})$ the set of label probability vectors. We can now generalise $\Theta_{TP}$ and $\Theta_{AE}$ to soft segmentations:

$$\Theta_{AE} = \sum_{i \in \mathbf{X}} |p^i - g^i|, \quad \Theta_{TP} = \sum_{i \in \mathbf{X}} g^i(1 - |p^i - g^i|) \tag{2}$$

In the common case of a crisp segmentation $g$ (i.e. $\forall i \in \mathbf{X}, g^i \in \{0, 1\}$), the associated soft Dice score can be expressed as:

$$D(p, g) = \frac{2 \sum_i g^i p^i}{\sum_i (g^i + p^i)} \tag{3}$$

A second variant has been used in [15], with a quadratic term in the denominator.

## 2.3   Previous work on multi-class Dice score

The easiest way to derive a unique criterion from the soft binary Dice score for multi-class segmentation is to consider the mean Dice score:

$$D_{mean}(p, g) = \frac{1}{|\mathbf{L}|} \sum_{l \in \mathbf{L}} \frac{2 \sum_i g_l^i p_l^i}{\sum_i (g_l^i + p_l^i)} \tag{4}$$

where $\{g_l^i\}_{i \in \mathbf{X}, l \in \mathbf{L}}$, $\{p_l^i\}_{i \in \mathbf{X}, l \in \mathbf{L}}$ are the set label probability vectors for all voxels for the ground truth and the prediction.

A generalised soft multi-class Dice score has also been proposed in [4,18] by generalising the set theory definition of the Dice score (1):

$$D_{FM}(p, g) = \frac{2 \sum_l \alpha_l \sum_i \min(p_l^i, g_l^i)}{\sum_l \alpha_l \sum_i (p_l^i + g_l^i)} \tag{5}$$

where $\{\alpha_l\}_{l \in \mathbf{L}}$ allows to weight the contribution of each class. However, those definitions are still based only on pairwise comparisons of probabilities associated with the same label and don't take into account inter-class relationships.

## 2.4   Wasserstein distance between label probability vectors

The Wasserstein distance (also sometimes called the *Earth Mover's Distance*) represents the minimal cost to transform a probability vector $p$ into another one $q$ when for all $l, l' \in \mathbf{L}$, the cost to move a unit from $l$ to $l'$ is defined as the distance $M_{l,l'}$ between $l$ and $l'$. This is a way to map a distance matrix $M$ (often

referred to as the *ground distance matrix*) on $\mathbf{L}$, into a distance on $P(\mathbf{L})$ that leverages prior knowledges about $\mathbf{L}$. In the case of a finite set $\mathbf{L}$, for $p, q \in P(\mathbf{L})$, the Wasserstein distance between $p$ and $q$ derived from $M$ can be defined as the solution of a linear programming problem [17]:

$$W^M(p,q) = \min_{T_{l,l'}} \sum_{l,l' \in \mathbf{L}} T_{l,l'} M_{l,l'},$$

$$\text{subject to } \forall l \in \mathbf{L}, \sum_{l' \in \mathbf{L}} T_{l,l'} = p_l, \text{ and } \quad \forall l' \in \mathbf{L}, \sum_{l \in \mathbf{L}} T_{l,l'} = q_{l'}. \tag{6}$$

where $T = (T_{l,l'})_{l,l' \in \mathbf{L}}$ is a joint probability distribution for $(p,q)$ with marginal distributions $p$ and $q$. A value $\hat{T}$ that minimises (6) is called an *optimal transport* between $p$ and $q$ for the distance matrix $M$.

### 2.5   Soft multi-class Wasserstein Dice score

The Wasserstein distance $W^M$ in (6) yields a natural way to compare two label probability vectors in a semantically meaningful manner by supplying a distance matrix $M$ on $\mathbf{L}$. Hence we propose using it to generalise the measure of disagreement between a pair of label probability vectors and provide the following generalisations:

$$\Theta_{AE} = \sum_{i \in \mathbf{X}} W^M(p^i, g^i) \tag{7}$$

$$\Theta_{TP}^l = \sum_{i \in \mathbf{X}} g_l^i (W^M(l,b) - W^M(p^i, g^i)), \quad \forall l \in \mathbf{L} \setminus \{b\} \tag{8}$$

where $W^M(l,b)$ is shorthand for $M_{l,b}$ and $M$ is chosen such that the background class $b$ is always the furthest away from the other classes. To generalise $\Theta_{TP}$, we propose to weight the contribution of the classes similarly to (5):

$$\Theta_{TP} = \sum_{l \in \mathbf{L}} \alpha_l \, \Theta_{TP}^l \tag{9}$$

We chose $\alpha_l = W^M(l,b)$ to make sure that background voxels do not contribute to $\Theta_{TP}$. The Wasserstein Dice score with respect to $M$ can then be defined as:

$$D^M(p,g) = \frac{2 \sum_l W^M(l,b) \sum_i g_l^i (W^M(l,b) - W^M(p^i, g^i))}{2 \sum_l [W^M(l,b) \sum_i g_l^i (W^M(l,b) - W^M(p^i, g^i))] + \sum_i W^M(p^i, g^i)} \tag{10}$$

In the binary case, setting $M = \left[\begin{smallmatrix} 0 & 1 \\ 1 & 0 \end{smallmatrix}\right]$ leads to $W^M(p^i, g^i) = |p^i - g^i|$ and reduces the proposed Wasserstein Dice score to the soft binary Dice score (2).

### 2.6   Wasserstein Dice loss with crisp ground truth

Previous work on Wasserstein distance-based loss functions for deep learning have been limited because of the computational burden [17]. However, in the

case of a crisp ground-truth $\{g^i\}_i$, and for any prediction $\{p^i\}_i$, a closed-form solution exists for (6). An optimal transport is $\forall l, l' \in \mathbf{L}, \hat{T}_{l,l'} = p_l^i g_{l'}^i$ and the Wasserstein distance becomes:

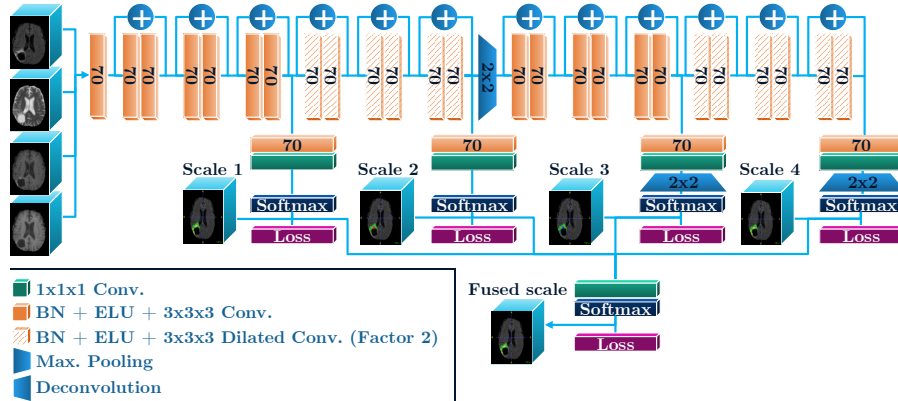$$W^M(p^i, g^i) = \sum_{l,l' \in \mathbf{L}} M_{l,l'} p_l^i g_{l'}^i \tag{11}$$

We define the Wasserstein Dice loss derived from $M$ as $\mathcal{L}_{D^M} := 1 - D^M$.

## 3 Holistic convolutional networks for multi-scale fusion

We now describe a holistically-nested convolutional neural network (HCNN) for imbalanced multi-class brain tumour segmentation inspired by the holistically-nested edge detection (HED) introduced in [19]. HCNN has been used successfully for some imbalanced learning tasks such as edge detection in natural images [19] and surgical tool segmentation [8]. The HCNN features multi-scale prediction and intermediate supervision. It can produce a unified output using a fusion layer while implicitly embedding spatial information in the loss. We further improve on the ability of HCNNs to deal with imbalanced datasets by leveraging the proposed generalised Wasserstein Dice loss. To keep up with state-of-the-art CNNs, we also employ ELU as activation function [3] and use residual connections [10]. Residual blocks include a pair of $3^3$ convolutional filters and Batch Normalisation [20]. The proposed architecture is illustrated in Fig.2.

### 3.1 Multi-scale prediction to leverage spatial consistency in the loss

As the receptive field increases across successive layers, predictions computed at different layers embed spatial information at different scales. Especially for imbalanced multi-class segmentation, different scales can contain complementary



**Fig. 2.** Proposed holistically-nested CNN for multi-class labelling of brain tumours.

information. In this paper, to increase the receptive field and avoid redundancy between successive scale predictions, max pooling and dilated convolutions (with a factor of 2 similar to [13]) have been used. As predictions are computed regularly at intermediate scales along the network (Fig.2), we chose to increase the number of features before the first prediction is made. For simplicity reasons, we then selected the same value for all hidden layers (fixed to 70 given memory constraints).

### 3.2   Multi-scale fusion and deep supervision for multi-class segmentation

While classic CNNs provide only one output, HCNNs provide outputs $\hat{y}^s$ at $S$ different layers of the network, and combine them to provide a final output $\hat{y}^{fuse}$:

$$(\hat{y}_l^{fuse})_{l \in \mathbf{L}} = \text{Softmax}\Big( \big( \sum_{s=1}^{S} w_{l,s} \hat{y}_l^s \big)_{l \in \mathbf{L}} \Big).$$

As different scales can be of different importance for different classes we learn class-specific fusion weights $w_{l,s}$. This transformation can also be represented by a convolution layer with kernels of size $1^3$ where the multi-scale predictions are fused in separated branches for each class, as illustrated in Fig. 2 similarly to the scalable layers introduced in [5]. In addition to applying the loss function $\mathcal{L}$ to the fused prediction, $\mathcal{L}$ is also applied to each scale-specific prediction thereby providing deep supervision (coefficients $\bar{\lambda}$ and $\lambda_s$ are set to $1/(S+1)$ for simplicity):

$$\mathcal{L}_{Total}((\hat{y}^s)_{s=1}^S, \hat{y}^{fuse}, y) = \bar{\lambda}\mathcal{L}(\hat{y}^{fuse}, y) + \sum_{s=1}^{S} \lambda_s \mathcal{L}(\hat{y}^s, y)$$

## 4   Implementation details

### 4.1   Brain tumour segmentation

We evaluate our HCNN model and Wasserstein Dice loss functions on the task of brain tumour segmentation using BraTS'15 training set that provides multimodal images (T1, T1c, T2 and Flair) for 220 high-grade gliomas subjects and 54 low-grade gliomas subjects. We divide it randomly into 80% for training, 10% for validation and 10% for testing so that the proportion of high-grade and low-grade gliomas subjects is the same in each fold. The scans are labelled with five classes (Fig. 1): (0) background, (1) necrotic core, (2) edema, (3) non-enhancing core and (4) enhancing tumour. The most common evaluation criteria for BraTS is to use the Dice scores for the whole tumour (labels 1,2,3,4), the core tumour (labels 1,3,4) and the enhanced tumour (label 4). All the scans of BraTS dataset are skull stripped, resampled to a 1mm isotropic grid and co-registered to the T1-weighted volume of each patient. Additionally, we applied histogram standardisation to each imaging modality independently [16].

**Table 1.** Evaluation of different multi-class Dice scores for training and testing. $\mathcal{L}_{D^{M_{tree}-PT}}$ stands for pre-training the HCNN with mean Dice score (4 epochs) and retraining it with $\mathcal{L}_{D^{M_{tree}}}$ (85 epochs).

| Loss function | Evaluation: Mean(std) Dice scores (%) | | | | | |
|---|---|---|---|---|---|---|
| | Whole | Core | Enh. | Mean Dice | $D^{M_{0-1}}$ | $D^{M_{tree}}$ |
| Mean Dice | 83(13) | 70(21) | 68(26) | **60(12)** | 77(11) | 80(12) |
| $\mathcal{L}_{D^{M_{0-1}}}$ | 86(12) | 59(29) | 69(23) | 48(5) | 82(6) | 85(5) |
| $\mathcal{L}_{D^{M_{tree}}}$ | 88(8) | 73(23) | 70(25) | 54(7) | 84(5) | 86(5) |
| $\mathcal{L}_{D^{M_{tree}-PT}}$ | **89(6)** | **73(22)** | **74(23)** | 59(10) | **84(4)** | **87(4)** |

### 4.2   Implementation details

We train the networks using ADAM [12] with a learning rate $lr = 0.01$, $\beta_1 = 0.9$ and $\beta_2 = 0.999$. To regularise the network, we use early stopping on the validation set and dropout in all residual blocks before the last activation (as proposed in [20]), with a probability of 0.6. We use multi-modal volumes of size $80^3$ from one subject concatenated as input during training and a sampling strategy to maximise the number of classes in each patch. Experiments have been performed using Tensorflow 1.1 [4] and a Nvidia GeForce GTX Titan X GPU.

## 5   Results

We evaluate the usefulness of the proposed soft multi-class Wasserstein Dice loss and the proposed HCNN with deep supervision. We compare the soft multi-class Wasserstein Dice loss to the state-of-the-art mean Dice score [5,13] for the training of our HCNN in Table 1 and 2. We also evaluate the segmentation at the different scales of the HCNN in Table 3.
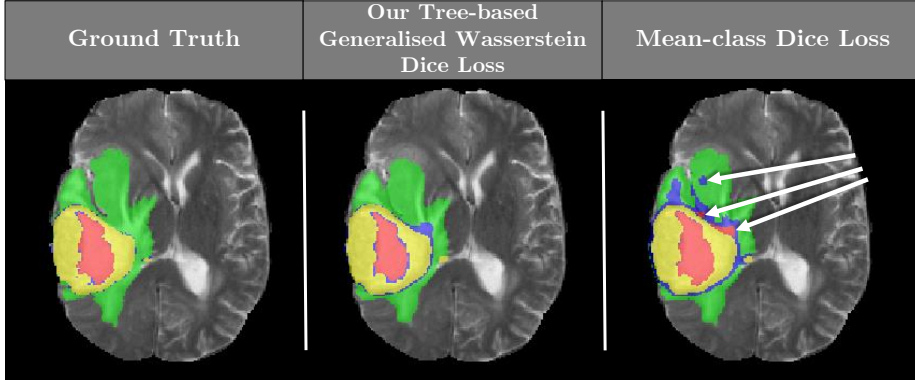
### 5.1   Examples of distance metrics on BraTS label space

To illustrate the flexibility of the proposed generalised Wasserstein Dice score, we evaluate two semantically driven choices for the distance matrix $M$ on $\mathbf{L}$:

$$M_{0-1} = \begin{pmatrix} 0 & 1 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 & 1 \\ 1 & 1 & 0 & 1 & 1 \\ 1 & 1 & 1 & 0 & 1 \\ 1 & 1 & 1 & 1 & 0 \end{pmatrix}, \quad \text{and} \quad M_{tree} = \begin{pmatrix} 0 & 1 & 1 & 1 & 1 \\ 1 & 0 & 0.6 & 0.2 & 0.5 \\ 1 & 0.6 & 0 & 0.6 & 0.7 \\ 1 & 0.2 & 0.6 & 0 & 0.5 \\ 1 & 0.5 & 0.7 & 0.5 & 0 \end{pmatrix}.$$

$M_{0-1}$ is associated with the discrete distance on $\mathbf{L}$ with no inter-class relationship. $M_{tree}$ is derived from the tree structure of $\mathbf{L}$ illustrated in Fig. 1. This tree is based on the tumour hierarchical structure: whole, core and enhancing tumour. We set branch weights to 0.1 for contiguous nodes and 0.2 otherwise.

---

[4] The code is publicly available as part of NiftyNet (http://niftynet.io)

**Fig. 3.** Qualitative comparison of HCNN predictions at testing after training with the proposed Generalised Wasserstein Dice loss ($\mathcal{L}_{D^{M_{tree}-PT}}$) or mean-class Dice loss. Training with $\mathcal{L}_{D^{M_{tree}-PT}}$ allows avoiding implausible misclassifications encountered in predictions after training with mean-class Dice loss (emphasized by white arrows).

### 5.2   Evaluation and training with multi-class Dice score

The mean Dice corresponds to the mean of soft Dice scores for each class as used in [5,13]. Results in Table 1 confirm that training with mean Dice score, $D^{M_{0-1}}$ or $D^{M_{tree}}$ allow maximising results for the associated multi-class Dice score during inference.

While $D^{M_{tree}}$ takes advantage of prior information about the hierarchical structure of the tumour classes it makes the optimisation more complex by adding more constraints. To relax those constraints, we propose to pretrain the network using the mean Dice score during a few epochs (4 in our experiment) and then retrain it using $D^{M_{tree}}$. This approach leads to the best results for all criteria, as illustrated in the last line of Table 1. Moreover, it produces segmentations that are more semantically plausible compared to the HCNN trained with mean Dice only as illustrated by Fig. 3.

### 5.3   Impact of the Wasserstein Dice loss on class confusion

Evaluating brain tumour segmentation using Dice scores of label subsets like whole, core and enhancing tumour doesn't allow measuring the ability of a model to learn inter-class relationships and to favour voxel classifications, be it correct or not, that are semantically as close as possible to the ground truth. We propose to measure class confusion using pairwise comparisons of all labels pair between the predicted segmentation and the ground truth (Table 2). Mathematically, for all $l, l' \in \mathbf{L}$, the quantity in row $l$ and colomn $l'$ stands for the soft binary Dice score:

$$D_{l,l'} = \frac{2 \sum_i g_l^i p_{l'}^i}{\sum_i (g_l^i + p_{l'}^i)} \tag{12}$$

**Table 2.** Dice score evaluation of the confusion after training the HCNN using different loss functions. Each line (resp. column) corresponds to the mean(standard deviation) Dice scores (%) of a region of the ground truth (resp. prediction) with all regions of the prediction (resp. ground truth) computed on the testing set.

| Mean Dice | Prediction | | | | |
|---|---|---|---|---|---|
| **Ground truth** | Background | Necrotic core | Edema | Non-enh. | Enh. |
| Background | 99.6(0) | 0(0) | 0.8(0) | 0.1(0) | 0.1(0) |
| Necrotic core | 0.(0) | 36.8(30) | 1.2(3) | 8.3(9) | 1(1) |
| Edema | 0.3(0) | 0.9(1) | 62.9(18) | 21.7(13) | 4.3(6) |
| Non-enh. | 0.1(0) | 8.7(9) | 6.5(8) | 33(15) | 14.8(11) |
| Enh. | 0(0) | 0.9(1) | 0.3(0) | 6.9(7) | 67.6(25) |

| $\mathcal{L}_{D^{M_{tree}}}$ | Prediction | | | | |
|---|---|---|---|---|---|
| **Ground truth** | Background | Necrotic core | Edema | Non-enh. | Enh. |
| Background | 99.7(0) | 0(0) | 0.3(0) | 0(0) | 0(0) |
| Necrotic core | 0(0) | 0(0) | 2.4(5) | 28.2(22) | 1.4(1) |
| Edema | 0.6(0) | 0(0) | 71.3(12) | 8.5(7) | 3.5(5) |
| Non-enh. | 0.1(0) | 0(0) | 15.4(13) | 28.9(14) | 14.2(10) |
| Enh. | 0(0) | 0(0) | 1.7(1) | 6.9(7) | 70.5(25) |

| $\mathcal{L}_{D^{M_{tree}-PT}}$ | Prediction | | | | |
|---|---|---|---|---|---|
| **Ground truth** | Background | Necrotic core | Edema | Non-enh. | Enh. |
| Background | 99.7(0) | 0(0) | 0.2(0) | 0(0) | 0(0) |
| Necrotic core | 0(0) | 20.2(27) | 2.3(5) | 23.2(18) | 1(1) |
| Edema | 0.6(0) | 0.3(0) | 73.3(11) | 5.7(5) | 3.1(4) |
| Non-enh. | 0.1(0) | 2.1(7) | 16.1(13) | 30(17) | 13(8) |
| Enh. | 0(0) | 0(0) | 2.4(2) | 3.8(4) | 73.5(22) |

Results in Table 2 compare class confusion of the proposed HCNN after being trained either using mean Dice loss, tree-based Wasserstein Dice loss ($\mathcal{L}_{D^{M_{tree}}}$) or tree-based Wasserstein Dice loss pre-trained with mean Dice loss($\mathcal{L}_{D^{M_{tree}-PT}}$). The first one aims only at maximising the true positives (diagonal) while the two other additionally aim at balancing the misclassifications to produce semnatically meaningful segmentations.

The network trained with mean Dice loss segments correctly most of the voxels (diagonal in Table 2) but makes misclassifications that are not semantically meaningful. For example, it makes poor differentiation between the edema and the core tumour as can be seen in the line corresponding to edema in Table 2 and in Fig. 3.

In contrast, the network trained with $\mathcal{L}_{D^{M_{tree}}}$ makes more meaningful confusion but it is not able to differentiate necrotic core and non-enhancing tumour at all (columns 2 and 4). It illustrates the difficulty to train the network with $\mathcal{L}_{D^{M_{tree}}}$ starting from a random initialisation because $\mathcal{L}_{D^{M_{tree}}}$ embeds more constraints than the mean Dice loss.

$\mathcal{L}_{D^{M_{tree}-PT}}$ allows combining advantages of both loss function: pre-training the network using the mean Dice loss allows initialising it so that it produces

**Table 3.** Evaluation of scale-specific and fused predictions of the HCNN with Dice score of whole, core, enhancing tumour and $D^{M_{tree}}$ after being pre-trained with mean Dice score (4 epochs) and retrained with $\mathcal{L}_{D^{M_{tree}}}$ (85 epochs).

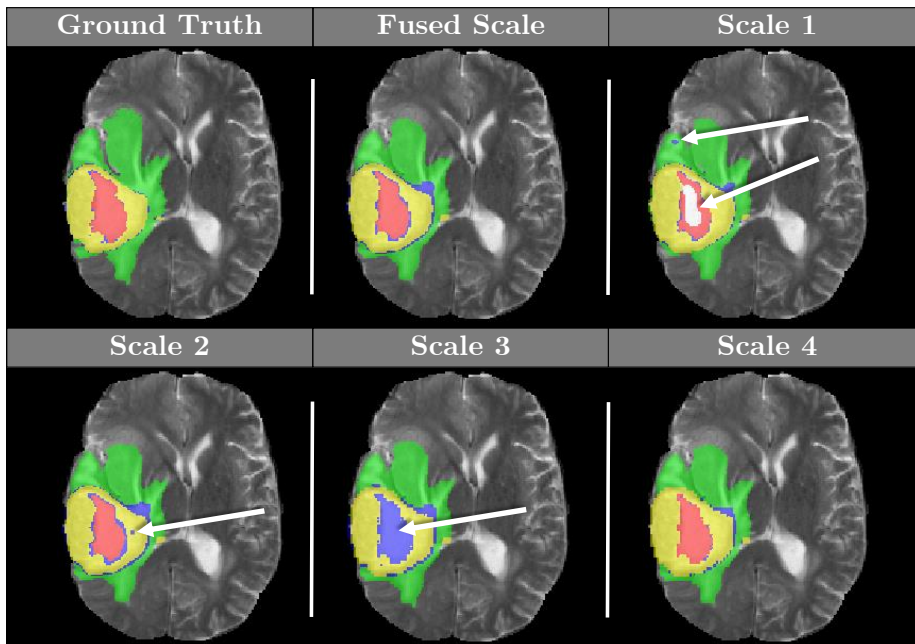| Prediction | Mean(Std) Dice score (%) | | | |
|---|---|---|---|---|
| | Whole tumour | Core tumour | Enh. tumour | $D^{M_{tree}}$ |
| Scale 1 | 84(8) | 68(23) | 70(25) | 84(5) |
| Scale 2 | **89(5)** | **73(22)** | **74(23)** | **87(4)** |
| Scale 3 | 88(6) | 72(23) | 71(22) | 86(4) |
| Scale 4 | **89(5)** | 72(22) | 71(21) | 86(3) |
| Fused | 89(6) | **73(22)** | **74(23)** | **87(4)** |

quickly an approximation of the segmentation, and retraining it with $\mathcal{L}_{D^{M_{tree}}}$ allows reaching a model which provides semantically meaningful segmentations (Fig. 3) with a higher rate of true positives compared to training with $\mathcal{L}_{D^{M_{tree}}}$ or mean Dice loss alone (Table 2).

### 5.4   Evaluation of deep supervision

Results in Table 3 are obtained after pre-training HCNN with mean Dice score during 4 epochs and then training it with $\mathcal{L}_{D^{M_{tree}}}$ during 85 additional epochs. Scales 2 to 4 and fused achieve similar Dice scores for whole, core tumour and the objective function $D^{M_{tree}}$ while scale 1 obtains lower Dice scores. Holes in tumour segmentations produced by scale 1, as illustrated in Fig. 4, suggest an unsufficient receptive field could account for those lower Dice scores. The best result for the enhancing tumour is achievied by both scale 2 and fused, which was expected as this is the smallest region of interest and the full resolution is maintained until scale 2. Moreover, as illustrated in Fig. 4, scales 3 and 4 fail at segmenting the thinest regions of the tumour because of their lower resolution contrary to scales 1 and 2 and fused. However, scales 1 to 3 contained implausible segmentation regions contrary to scale 4 and fused. This suggests trade-offs between high receptive field and high resolution that are class specific. It confirms the usefulness of the multi-scale holistic approach for the multi-class brain tumour segmentation task.

## 6   Conclusion and future work

We proposed a semantically driven generalisation of the Dice score for soft multiclass segmentation based on the Wasserstein distance. This embeds prior knowledge about inter-class relationships represented by a distance matrix on the label space. Additionally, we proposed a holistic convolutional network that uses multiscale predictions and deep supervision to make use of multi-scale information. We successfully used the proposed Wasserstein Dice score as a loss function to train our holistic networks and show the importance of multi-scale and inter-class

**Fig. 4.** Qualitative comparison of fused and scales predictions at testing after training our HCNN with the proposed Generalised Wasserstein Dice loss ($\mathcal{L}_{DM_{tree-PT}}$). White arrows emphasize implausible misclassifications.

relationships for the imbalanced task of multi-class brain tumour segmentation. The proposed distance matrix based on the label space tree structure leads to higher Dice scores compared to the discrete distance. Because the tree-based distance matrix used was heuristically chosen we think that better heuristics or a method to directly learn the matrix from the data could lead to further improvements.

As the memory capacity of GPUs increases, entire multi-modal volumes could be used as input of CNN-based segmentation. However, it will also increase the class imbalance in the patches used as input. We expect this to increase the impact of our contributions. Future work includes extending the use of Wasserstein distance by defining a matrix distance on the entire output space $\mathbf{X} \times \mathbf{L}$ similarly to [6]. This would allow embedding spatial information directly in the loss, but the computation burden of the Wasserstein distance, in that case, remains a challenge [17].

## References

1. Anbeek, P., Vincken, K.L., van Bochove, G.S., van Osch, M.J., van der Grond, J.: Probabilistic segmentation of brain tissue in MR imaging. NeuroImage 27(4), 795 – 804 (2005)
2. Chang, H.H., Zhuang, A.H., Valentino, D.J., Chu, W.C.: Performance measure characterization for evaluating neuroimage segmentation algorithms. Neuroimage (2009)
3. Clevert, D.A., Unterthiner, T., Hochreiter, S.: Fast and accurate deep network learning by exponential linear units (elus). arXiv:1511.07289 (2015)
4. Crum, W.R., Camara, O., Hill, D.L.: Generalized overlap measures for evaluation and validation in medical image analysis. IEEE TMI 25(11), 1451–1461 (2006)
5. Fidon, L., Li, W., Garcia-Peraza-Herrera, L.C., Ekanayake, J., Kitchen, N., Ourselin, S., Vercauteren, T.: Scalable multimodal convolutional networks for brain tumour segmentation. MICCAI (2017)
6. Fitschen, J.H., Laus, F., Schmitzer, B.: Optimal transport for manifold-valued images. In: International Conference on Scale Space and Variational Methods in Computer Vision. pp. 460–472. Springer (2017)
7. Frogner, C., Zhang, C., Mobahi, H., Araya, M., Poggio, T.A.: Learning with a wasserstein loss. In: NIPS. pp. 2053–2061 (2015)
8. Garcia-Peraza-Herrera, L.C., Li, W., Fidon, L., Gruijthuijsen, C., Devreker, A., Attilakos, G., Deprest, J., Vander Poorten, E., Stoyanov, D., Vercauteren, T., Ourselin, S.: ToolNet : Holistically-Nested Real-Time Segmentation of Robotic Surgical Tools. In: IROS (2017)
9. Havaei, M., Davy, A., Warde-Farley, D., Biard, A., Courville, A., Bengio, Y., Pal, C., Jodoin, P.M., Larochelle, H.: Brain tumor segmentation with deep neural networks. Med. Image Anal. 35, 18–31 (2017)
10. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: IEEE CVPR (2016)
11. Kamnitsas, K., Ledig, C., Newcombe, V.F., Simpson, J.P., Kane, A.D., Menon, D.K., Rueckert, D., Glocker, B.: Efficient multi-scale 3D CNN with fully connected CRF for accurate brain lesion segmentation. Med. Image Anal. 36, 61–78 (2017)
12. Kingma, D., Ba, J.: Adam: A method for stochastic optimization. arXiv:1412.6980 (2014)
13. Li, W., Wang, G., Fidon, L., Ourselin, S., Cardoso, M.J., Vercauteren, T.: On the compactness, efficiency, and representation of 3D convolutional networks: Brain parcellation as a pretext task. IPMI (2017)
14. Menze, B.H., Jakab, A., Bauer, S., Kalpathy-Cramer, J., Farahani, K., Kirby, J., Burren, Y., Porz, N., Slotboom, J., Wiest, R., et al.: The multimodal brain tumor image segmentation benchmark (BraTS). IEEE TMI 34(10), 1993–2024 (2015)
15. Milletari, F., Navab, N., Ahmadi, S.A.: V-Net: Fully convolutional neural networks for volumetric medical image segmentation. In: Proc. 3DV'16. pp. 565–571 (2016)
16. Nyul, L.G., Udupa, J.K., Zhang, X.: New variants of a method of MRI scale standardization. IEEE TMI 19(2), 143–150 (Feb 2000)
17. Pele, O., Werman, M.: Fast and robust earth mover's distances. ICCV (2009)
18. Sudre, C.H., Li, W., Vercauteren, T., Ourselin, S., Cardoso, M.J.: Generalised Dice overlap as a deep learning loss function for highly unbalanced segmentations. arXiv:1707.03237 (2017)
19. Xie, S., Tu, Z.: Holistically-nested edge detection. In: ICCV (2015)
20. Zagoruyko, S., Komodakis, N.: Wide residual networks. arXiv:1605.07146 (2016)