



**Universiteit
Leiden**
The Netherlands

Deep learning for automatic segmentation of tumors on MRI

Rodríguez Outeiral, R.

Citation

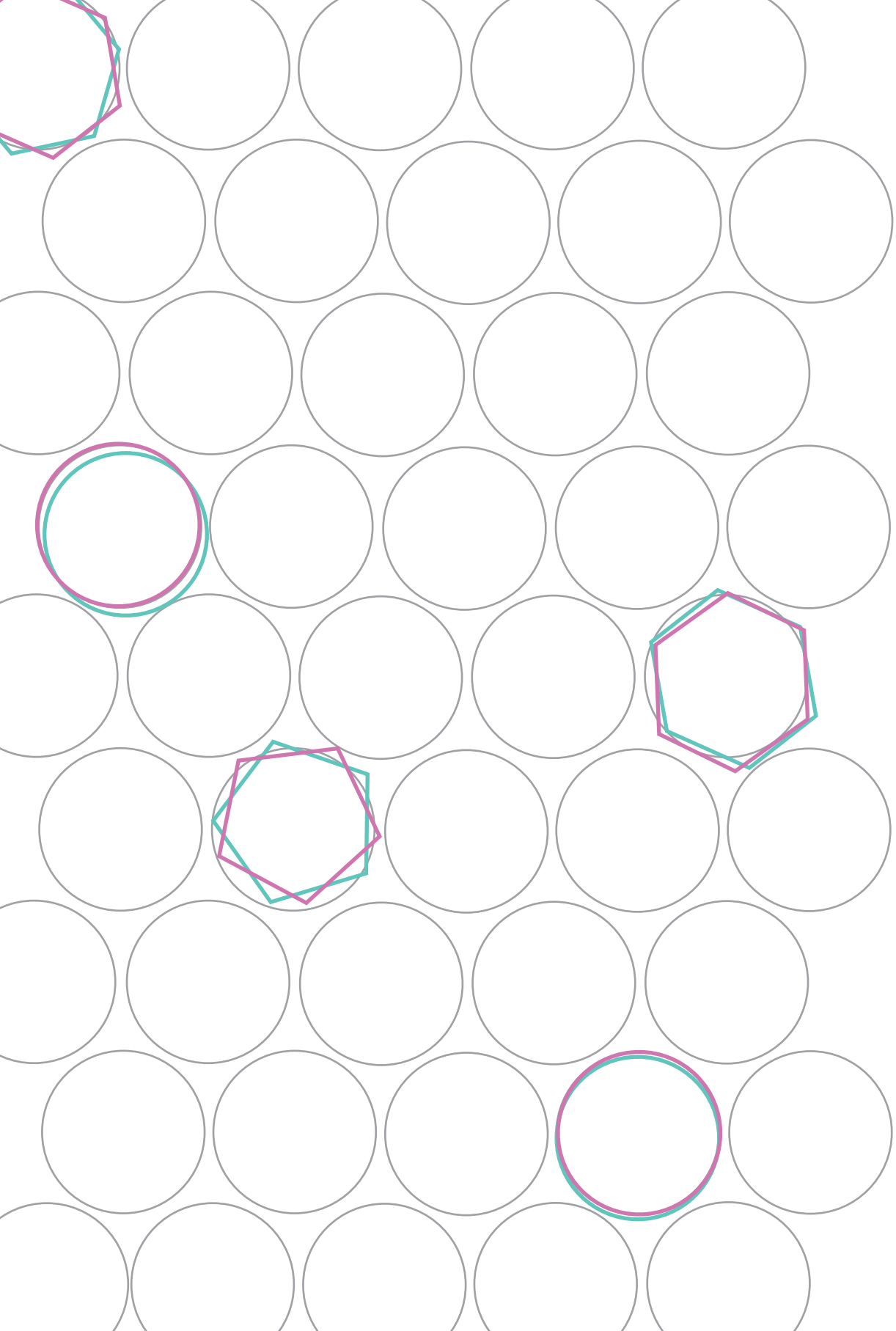
Rodríguez Outeiral, R. (2024, June 25). *Deep learning for automatic segmentation of tumors on MRI*. Retrieved from <https://hdl.handle.net/1887/3765390>

Version: Publisher's Version

License: [Licence agreement concerning inclusion of doctoral thesis in the Institutional Repository of the University of Leiden](#)

Downloaded from: <https://hdl.handle.net/1887/3765390>

Note: To cite this publication please use the final published version (if applicable).



Chapter 3

Strategies for tackling the class imbalance problem of oropharyngeal primary tumor segmentation on magnetic resonance imaging

Rodríguez Outeiral R, Bos P, van der Hulst HJ, Al-Mamgani A, Jasperse B, Simões R, et al. *Phys Imaging Radiat Oncol.* 23;144–9 (2022)

ABSTRACT

Background and purpose

Contouring of the oropharyngeal primary tumor is currently done manually which is time-consuming. Autocontouring techniques based on deep learning methods are a desirable alternative, but these methods can render suboptimal results when the structure to segment is considerably smaller than the rest of the image. The purpose of this work was to investigate different strategies to tackle the class imbalance problem in this tumor site.

Materials and methods

A cohort of 230 oropharyngeal cancer patients treated between 2010 and 2018 was retrospectively collected. The following magnetic resonance imaging (MRI) sequences were available: T1-weighted, T2-weighted, 3D T1-weighted after gadolinium injection. Two strategies to tackle the class imbalance problem were studied: training with different loss functions (namely: Dice loss, Generalized Dice loss, Focal Tversky loss and Unified Focal loss) and implementing a two-stage approach (i.e. splitting the task in detection and segmentation). Segmentation performance was measured with Sørensen–Dice coefficient (Dice), 95th Hausdorff distance (HD) and Mean Surface Distance (MSD).

Results

The network trained with the Generalized Dice Loss yielded a median Dice of 0.54, median 95th HD of 10.6 mm and median MSD of 2.4 mm but no significant differences were observed among the different loss functions (p -value > 0.7). The two-stage approach resulted in a median Dice of 0.64, median HD of 8.7 mm and median MSD of 2.1 mm, significantly outperforming the end-to-end 3D U-Net (p -value < 0.05).

Conclusions

No significant differences were observed when training with different loss functions. The two-stage approach outperformed the end-to-end 3D U-Net.

INTRODUCTION

Radiotherapy is one of the common treatment options for head and neck cancer patients [4,41]. One key step of the radiotherapy workflow is tumor contouring. While contouring of organs at risk is increasingly being automated in clinical practice, tumor contouring is still done manually. This is time consuming and suffers from high interobserver variability [11].

Deep learning methods, particularly Convolutional Neural Networks (CNNs), are the current state of the art for automatic segmentation of medical images. Several review papers have been published on deep learning applied to radiotherapy and automatic segmentation is often discussed as one of the main applications [31,60,62,63]. For the particular case of head and neck cancer, various works have focused on the automatic segmentation of organs at risk with deep learning [64], some of them achieving clinically acceptable performance and being commercially available [18]. For the case of tumor contouring, the literature is more scarce and those algorithms are still not implemented in the clinic.

In our previous work [65], we segmented the oropharyngeal primary tumor on magnetic resonance imaging (MRI) and showed that combining multiple anatomical MRI sequences improved the segmentation performance compared to single-sequence. We also proposed a semi-automatic approach that improved the segmentation performance by splitting the segmentation task in manual detection and segmentation. To the best of our knowledge, there is only one other work where the authors segmented the oropharyngeal primary tumor on MRI [42]. The authors studied the impact of combining different anatomical (T1 weighted and T2 weighted) and quantitative images (ADC, Ktrans and ve) as input channels to a CNN and showed that combining anatomical sequences significantly improved the performance.

A known issue in the field of deep learning for medical image segmentation is class imbalance, meaning that the structure to be segmented is present in a smaller amount of voxels compared to the rest of the image. Class imbalance can result in suboptimal solutions because the network is exposed to proportionally less relevant information during the training process. Several works in the field of medical image segmentation have focused on this problem, either by modifying the input data to the network [66,67] or by defining different loss functions [68–70]. This problem is even more critical in the case of tumor segmentation, given that tumors tend to be smaller than other structures and they are heterogeneous in their location, shape and size. This is also the case for the oropharyngeal primary tumor.

Several loss functions have been designed with the aim of tackling class imbalance, such as the Generalized Dice loss [58], the Focal loss [68], the (focal) Tversky loss [69,71] and the Unified Focal loss [70]. Although the choice of the loss function can be critical for

the training of a CNN, comprehensive loss function comparisons for specific tumor sites or anatomies are not commonly performed. Ma et al. [30] showed that the influence in performance of the loss function varies greatly depending on the segmentation task. To the best of our knowledge, this has not been studied yet in the particular case of oropharyngeal cancer segmentation.

Other works have implemented two-stage approaches (i.e. detection and segmentation) that resulted in more accurate segmentations than their one-stage counterparts [72–74]. By locating the tumor first, the context around the tumor is reduced. Consequently, two-stage approaches are a possible way of tackling class imbalance. The semi-automatic approach from our previous work [65] consisted of having human observers outlining a box around the tumor to provide a first approximation of the tumor location and consequently ease the segmentation task. However, the semi-automatic approach still needed manual intervention. The implementation of a two-stage approach will also allow us to fully automate the semi-automatic approach proposed in our previous work [65].

The aim of this study was to investigate two different strategies for tackling the class imbalance problem for oropharyngeal primary tumor segmentation: training with different loss functions and implementing a fully automatic two-stage approach.

MATERIALS AND METHODS

Data

A cohort of 230 patients treated at our institute between January 2010 and May 2018 was used for this project. The mean age of the patients was 61 years (standard deviation ± 7 years) and 66 % of the patients were male. Further details on tumor stage and HPV status can be found in the Supplemental Material (table S.1). All patients had histologically proven primary oropharyngeal squamous cell carcinoma and received a pre-treatment MRI for primary staging. The institutional review board approved the study (IRBd18047). Informed consent was waived by the institutional review board considering the retrospective design. The cohort was extended from our previous work [65]. A total of 59 new patients were included.

The scans were acquired on 1.5T (n=108) or 3.0T (n=122) MRI scanners (Philips Medical System, Best, The Netherlands). The imaging protocol included: 2D T1-weighted fast spin-echo, 2D T2-weighted fast spin-echo with fat suppression, 3D T1-weighted high-resolution isotropic volume excitation after gadolinium injection with fat suppression. Further details on the MRI protocols are given in the Supplemental Material (table S.2). The primary tumors were manually contoured in 3D Slicer (version 4.8.0, <https://www.slicer.org/>) by one observer with 1 year of experience (P.B. or H.H.). Afterwards, they were reviewed and adjusted, if needed, by a radiologist with 7 years of experience (B.J.).

All tumor volumes were delineated on the T1gd but the observers were allowed to consult the other sequences.

For the experimental set-up, the data set was split in three subsets: a training set (n=190), a validation set (n=20) and a test set (n=20). The test set was not used for training or hyper-parameter tuning. We stratified the three subsets for tumor volume, subsite, and aspect ratio since these features are likely relevant for segmentation. Subsites were defined as tonsillar tissue, soft palate, base of tongue and posterior wall. The aspect ratio was defined as the ratio between the shortest and the longest axis of the tumor. All images were resampled to a voxel size of 0.8 mm x 0.8 mm x 0.8 mm.

Baseline model architecture

The 3D U-Net architecture [53,57] was used as the basis for our experiments. The Adam optimizer [59] and early stopping were used for training. Dropout and data augmentation were used for regularization. Further details on the training procedure can be found in table S.4. and in the code which is publicly available in: https://github.com/RoqueRouteiral/oroph_segms_ts.

Training with different loss functions.

We trained the 3D U-Net with four different loss functions: Dice loss [75], Generalized Dice loss[58], Focal Tversky loss [71] and Unified Focal loss [70]. For the particular case of the Unified Focal loss, Yeung et al. [70] showed that the choice of the γ hyperparameter can affect the performance. Consequently, we trained four networks with the Unified Focal loss for different values of its hyperparameter γ ($\gamma = [0.2, 0.4, 0.6, 0.8]$). We compared the segmentation performance of all the networks among each other.

Two-stage approach

In our previous work, we demonstrated that the segmentation of the oropharyngeal primary tumor was more accurate when the input image was manually cropped with a clipbox around the tumor before being fed to a segmentation network.

In this work, we fully automated this two-stage approach (figure 1). The first stage consisted of roughly detecting the tumor by automatically selecting a clipbox around it. In the second stage, this clipbox was used to crop the image which was then used as input to a segmentation network. The loss function chosen for both stages was the Generalized Dice loss function. The loss was backpropagated through each network separately.

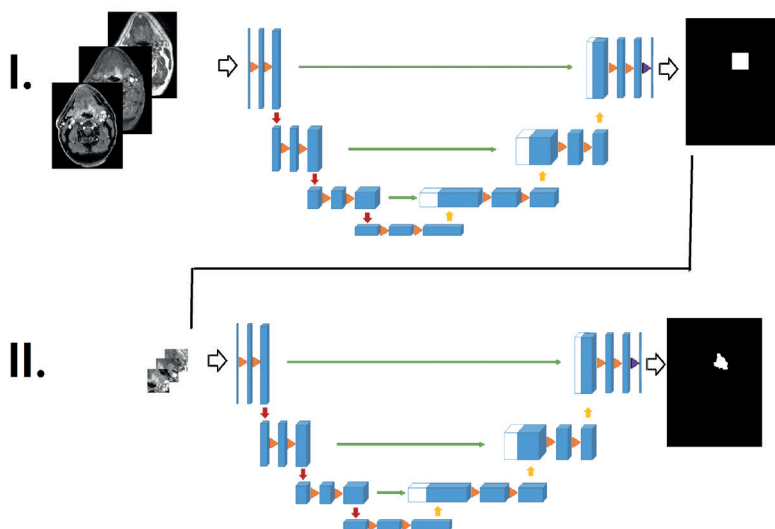


Figure 1. Overview of the two-stage approach.

For the detection stage, a 3D U-Net was trained using the bounding box of the tumor as ground truth. At inference time, the output of the detection was computed as the bounding box of the output.

For the segmentation stage, the same architecture as in our previous work was used [65]. This segmentation network was trained with only the information contained inside the clipboxes. In every training iteration, the clipboxes were randomly shifted by an amount of up to 25 mm in different directions to make the network robust to possible displacements in the detection. At inference time the input images were cropped by the clipboxes defined by the detection network. Similarly to our previous work, the clipboxes were dilated by 5 mm.

Statistics

To confirm that the three subsets were balanced in subsite, volume and aspect ratio, a Kruskal-Wallis test was used for continuous variables (volume and aspect ratio) and a chi-square test for independence for the categorical data (subsite).

Predicted segmentations and the segmentations from the human observers were compared for the patients on the unseen test set. Common segmentation metrics were used: Sørensen–Dice coefficient (Dice), 95th Hausdorff Distance (HD) and Mean surface distance (MSD). The metrics were implemented using the Python package from DeepMind (<https://github.com>).

com/deepmind/surface-distance). For the two-stage approach, the detection was evaluated by measuring the absolute mean shift in all 6 directions between the tumor bounding box and the detected clipbox for the patients on the unseen test set. The average shift of the boxes for the observers from our previous work was used for comparison [65]. Differences among the loss function experiments were assessed by the Friedman test whereas the two-stage approach experiments were assessed by the Wilcoxon signed-ranked test. P-values below 0.05 were considered statistically significant.

All networks were retrained four times. Reported results are the mean of the results of the four versions of each network. We opted for this approach over N-fold cross-validation to account for the random initialization of the network while ensuring proper stratification in the three sets for all the folds.

RESULTS

Summary of tumor characteristics

Table S.3 shows the tumor characteristics (location, volume and aspect ratio) of our cohort. No significant differences were found in the distributions of subsite, volume and aspect ratio between the training, validation and test sets.

Training with different loss functions

When comparing the performance of the networks trained with different loss functions no significant differences were found (p -value > 0.25 for the three metrics). Lower variance in the MSD and Dice can be observed for the network trained with the Generalized Dice loss (figure 2). The network achieved a median Dice of 0.54, median 95th HD of 10.6 mm and median MSD of 2.4 mm. Non-significant differences were observed when training the network with different γ values for the Unified Focal loss (Figure S.1).

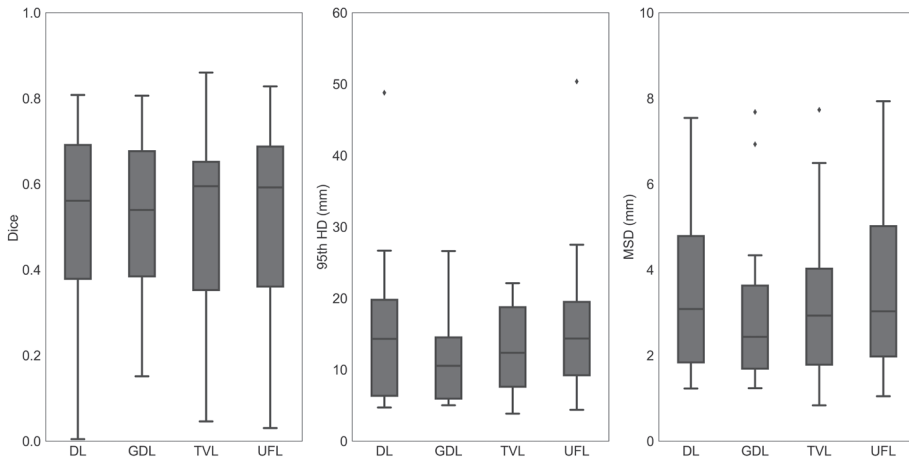


Figure 2. Segmentation performance of the 3D U-Net trained with different loss functions: Dice Loss (DL), Generalized Dice Loss (GDL), Tversky Loss (TVL) and Unified Focal Loss (UFL).

Two-stage approach

The mean shift for the detection network was of 8.9 mm (Table 1) and no significant differences were found when comparing to the detection of observer 2 from our previous work (p -value = 0.40). Significant differences were found when comparing the detection of this work to the detection of the observer 1 from our previous work (p -value < 0.001). When separating the mean shift per direction, we observed a mean shift of 10.0 mm in the cranial-caudal direction, 8.4 mm in the medial-lateral direction and 7.7 mm in the dorsal-ventral direction.

Table 1. Detection and segmentation performance of the two-stage approach and comparison to results of the previous work [65].

	Detection		Segmentation	
	Avg. shift (mm) – [SD]	Dice	HD (mm)	MSD (mm)
This work				
3D end-to-end UNet	--	0.54	10.6	2.4
Two-stage approach	8.7 [8.2]	0.64	8.7	2.1
Previous work				
Semi-automatic approach (Obs. 1)	3.0 [3.9]	0.74	4.6	1.2
Semi-automatic approach (Obs. 2)	8.9 [6.9]	0.67	7.2	1.7

The segmentation results of the two-stage approach were significantly better for Dice (p -value = 0.03) and MSD (p -value = 0.02) than the results of the end-to-end 3D UNet (Table 1). The fully automated two-stage approach yielded a median Dice of 0.64, median HD of 8.7 mm and median MSD of 2.1 mm. One patient was missed in the detection of the two-stage approach for one of the folds, and thus removed from that fold for the analysis.

Qualitative results

Examples of segmentations obtained by the end-to-end 3D U-Net, the two-stage approach and ground truth segmentation are shown in Figure 3. The end-to-end 3D U-Net approach oversegmented (Figure 3a-c) the tumor, where the two-stage approach showed better segmentation comparison to the ground truth. Figure 3b shows cases where the segmentation end-to-end 3D U-Net rendered additional false positive structures on the image.

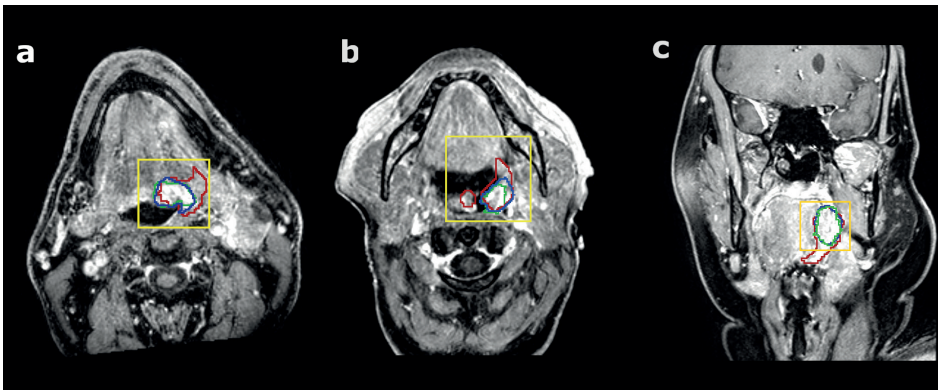


Figure 3. Comparison of the oropharyngeal segmentations in three different patients (a, b, c) trained with the end-to-end 3D U-Net (red), with the two-stage approach (blue) and the manual delineation (green). The yellow boxes are drawn by detection network from the two-stage approach. All the images correspond to the 3D T1gd sequence.

DISCUSSION

This work investigated two different strategies to tackle the class imbalance problem for the task of oropharyngeal primary tumor segmentation: training with the different loss functions and implementing a two-stage approach. Additionally, the proposed two-stage approach fully automated the semi-automatic approach described in our previous work [65].

When training the networks with different loss functions, no significant improvements were observed in the segmentation metrics. Hyperparameter tuning for the hyperparameter of the Unified Focal loss did not yield significantly better results either. This result is consistent with the work of Ma et al. [30], where they concluded that Dice-related losses are often optimal for medical image segmentation tasks. Additionally, it is also in line with the conclusions described by Isensee et al. and their proposed “no new Net” (nnU-Net) [76]. They showed that a tailored-to-task method configuration is more relevant than specific setup choices when designing a segmentation deep learning pipeline.

The two-stage approach achieved significantly better results compared to the conventional end-to-end approach. The high complexity of the task may make the end-to-end training of the network suboptimal, while focusing on two simpler tasks can render better results. In our previous work [10], a semi-automatic approach in which an observer selected a clipbox around the tumor was implemented. When comparing the current detection results to the semi-automatic approach of our previous work, we noted that one of the observers (Obs. 1) selected a tighter box (although all the tumors were included inside the clipboxes) compared to that of our two-stage approach which resulted in significantly different detection performance. However, we did not observe significant differences with the detection performance of the semi-automatic approach for the other observer (Obs. 2), showing that a fully automatic two-stage approach can be a feasible alternative to a semi-automatic approach. Also, the time spent on delineating in the clinical practice is aimed to be as low as possible. We reported in our previous work that the time spent on drawing the boxes was lower for observer 2 than for observer 1, making the delineations of observer 2 a more realistic representation of what is expected in the clinic. In the present work, the whole pipeline is automated, which can save time in the clinic. That said, further efforts in improving the detection are of interest to improve the segmentation performance of the two-stage approach.

The literature on automatic segmentation for the oropharyngeal tumor on MRI is scarce and its aims are heterogeneous. Besides our previous work [65], only Wahid et al. [42] have focused on the segmentation of this tumor site on MRI. Their work focused on studying the value of multiparametric MRI on the segmentation performance, both for qualitative and quantitative imaging. Other works focused on the automatic segmentation on multiparametric MRI of the head and neck cancer in general, rather than on the particular subset of oropharyngeal cancer: Bielak et al. [77] used diffusion weighted imaging while Schouten et al. [78] proposed a multiview CNN architecture. To the best of our knowledge, only our work is focused on tackling the class imbalance problem for head and neck cancer segmentation on MRI, and particularly for the oropharyngeal subsite.

In 2020, the first head and neck tumor segmentation challenge, known as HECKTOR challenge, was launched [79]. The main subsite of the challenge was the oropharyngeal

tumor and the winner of the challenge achieved a mean Dice of 0.76, but the image modalities used were PET/CT. Additionally, Ren et al. [80] compared the use of PET/CT/MRI as different input image combinations for the automatic segmentation of head and neck GTV and observed that, when including PET, the segmentation performance improved. Considering all the above, it is possible that PET is a useful modality for the task of head and neck tumor segmentation. However, the differences in resolution between imaging modalities may be reflected in the detail of the manual ground truth delineations used for training and evaluation. Potentially, this can also explain the difference in performance of the MRI-based task. That said, we argue that the strategies to tackle class imbalance in this work can be useful in the development of autocontouring tools for the case of oropharyngeal cancer.

This study has limitations. Firstly, there is a high interobserver variability on this tumor subsite, especially in case of tonsillar fossa and base of tongue tumor which are rich in lymphatic tissue, so it is possible that the ground truth delineations used in this work are partially biased. However, one observer corrected the other's delineation, reducing this observer variation. Secondly, validation of our results is still needed with an independent cohort in a multi-center study. Thirdly, the performance could also be improved by making different decisions on the training setup, such as using larger batch sizes or non downsampled data, but other strategies to mitigate memory limitations would be needed. Finally, there is a certain variability in the scan protocols. However, variability in the training set can be desirable as it makes the network robust to protocol differences.

In conclusion, the loss functions designed to tackle class imbalance performed comparably among each other. The approach of splitting the problem into localization and segmentation outperformed the end-to-end network, proving an effective strategy to mitigate the class imbalance problem in oropharyngeal cancer segmentation.

SUPPLEMENTAL MATERIAL

Table S.1. Tumor stage and HPV status.

	N patients (percentage in %)
Tumor staging	
T1	38 (16.52)
T2	85 (36.96)
T3	52 (22.61)
T4	55 (23.91)
N-stage	
N0	44 (19.13)
N1	33 (14.35)
N2	149 (64.78)
N3	4 (1.74)
Subsite	
Tonsillar tissue	121 (52.61)
Soft palate	21 (9.13)
Base of tongue	80 (34.78)
Posterior wall	8 (3.48)
HPV status	
Positive	106 (46.09)
Negative	101 (43.91)
Unknown	23 (10.0)

Table S.2. Overview of MRI sequences. The MRI sequences are 2D T1 weighted (T1w), 2D T2 weighted with fat suppression (T2w), 3D T1 weighted after gadolinium injection with fat suppression (T1gd).

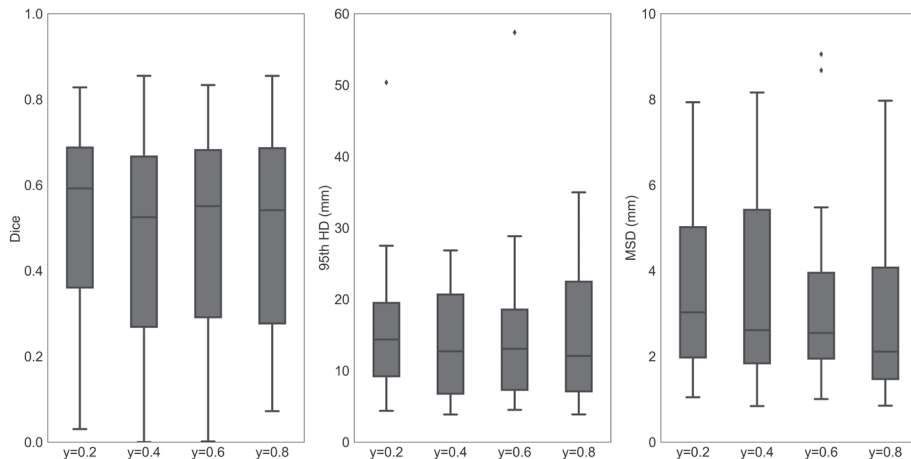
	TR (ms)	TE (ms)	Pixel size (mm)	Slice thickness (mm)
T1w	[180 – 890]	[2 - 10]	[0.4 – 0.9]	[3.0 – 5.0]
T2w	[2500 – 16000]	[20 – 90]	[0.4 – 0.9]	[3.0 – 5.0]
T1gd	[4 – 10]	[2 – 5]	[0.2– 1.0]	[0.8 – 2.0]

Table S.3. Overview of tumor characteristics per data set. *p-value was calculated using chi-square test. ** p-value was calculated using Kruskal-Wallis test

	Training set	Validation set	Testing set	p-value
Patients (n)	190	20	20	
Subsite				0.757*
Tonsillar tissue	93 (53.2 %)	10 (50.0 %)	10 (50.0 %)	
Soft palate	14 (7.90 %)	3 (15.0 %)	3 (15.0 %)	
Base of tongue	63 (34.7 %)	7 (35.0%)	7 (35.0 %)	
Posterior wall	7 (4.20 %)	0 (0 %)	0 (0 %)	
Volume (cm ³)				0.465**
Median	7.54	7.40	6.88	
Range	[0.23,71.54]	[0.51, 41.6]	[0.46, 17.20]	
Aspect ratio (%)				0.350**
Median	52.17	53.54	55.61	
Range	[16.21, 90.67]	[41.53, 64.35]	[47.47, 84.28]	

Table S.4. Training details of the networks of the paper.

Hyper-parameter / set-up	
Optimizer	Adam
Loss function	Dice Loss / Unified Focal Loss/ Generalized Dice Loss/ Focal Tversky Loss
Initial learning rate	0.001
Learning rate scheduler	Multiply by 0.5 if validation loss does not decrease in ten epochs by an amount of 0.001
Batch size	1* (*Batch normalization with running mean and variance during inference time, because of stability issues during training with batch size of 1)
Dropout	0.2 in bottleneck convolutions
Data augmentation (only full context)	Horizontal flip (in coronal view) with a chance of 0.5 in every epoch. Random elastic deformation Using elasticdeform library (https://github.com/gvtulder/elasticdeform) Random rotations between -10 and 10 degrees in every epoch. Downsampled by a factor of 2.5
Shifts (only second stage of two-stage approach)	For second stage of the training of the two stage approach, random shifts of the tumor within the box of up to 25 mm.

**Figure S.1.** Segmentation performance of the network trained with Unified Focal Loss for different values of γ .

