# USE-Evaluator: Performance Metrics for Medical Image Segmentation Models Supervised by Uncertain, Small or Empty Reference Annotations in Neuroimaging

Sophie Ostmeier[a],[**], Brian Axelrod[a], Fabian Isensee[b], Jeroen Bertels[c], Michael Mlynash[a], Soren Christensen[d], Maarten G. Lansberg[a], Gregory W. Albers[a], Rajen Sheth[e], Benjamin F.J. Verhaaren[c], Abdelkader Mahammedi[a], Li-Jia Li[a], Greg Zaharchuk[a], Jeremy J. Heit[a]

[a]*Stanford University, Center of Academic Medicine, 453 Quarry Rd, Palo Alto, CA 94304*
[b]*Division of Medical Image Computing, German Cancer Research Center (DKFZ), Im Neuenheimer Feld 280, 69120 Heidelberg, Germany*
[c]*KU Leuven, Herestraat 49, 3000 Leuven, Belgium*
[d]*gray number analytics, Lomma, Sweden*
[e]*in personal capacity*

## ABSTRACT

Performance metrics for medical image segmentation models are used to measure the agreement between the reference annotation and the predicted segmentation. Usually, overlap metrics, such as the Dice, are used as a metric to evaluate the performance of these models in order for results to be comparable.

However, there is a mismatch between the distributions of cases and the difficulty level of segmentation tasks in public data sets compared to clinical practice. Common metrics used to assess performance fail to capture the impact of this mismatch, particularly when dealing with datasets in clinical settings that involve challenging segmentation tasks, pathologies with low signal, and reference annotations that are uncertain, small, or empty. Limitations of common metrics may result in ineffective machine learning research in designing and optimizing models. To effectively evaluate the clinical value of such models, it is essential to consider factors such as the uncertainty associated with reference annotations, the ability to accurately measure performance regardless of the size of the reference annotation volume, and the classification of cases where reference annotations are empty.

We study how uncertain, small, and empty reference annotations influence the value of metrics on a stroke in-house data set regardless of the model. We examine metrics behavior on the predictions of a standard deep learning framework in order to identify suitable metrics in such a setting. We compare our results to the BRATS 2019 and Spinal Cord public data sets. We show how uncertain, small, or empty reference annotations require a rethinking of the evaluation. The evaluation code was released to encourage further analysis of this topic https://github.com/SophieOstmeier/UncertainSmallEmpty.git

## 1. Introduction

The performance of machine learning algorithms is assessed by metrics. The optimal choice of metrics depends on the data set and the machine learning task to guarantee that the predictions accurately describe the intended phenomenon (Taha and Hanbury, 2015). Metrics can be used in two different ways. First, as the criteria that the models try to optimize as a loss function. Second, as a way of validating and evaluating the performance of the model. This work focuses on the latter, referred to as performance metrics.

Performance metrics differ in their characteristics. The correlations between them determine the additional information revealed. Therefore, the appropriate selection of a performance metric for a specific task ensures consistency in model performance between development and deployment. For example, physicians that potentially use model predictions for treatment decisions of patients rely on an optimization and evaluation process of the models towards reliable and meaningful clinical information.

For data sets with uncertain (inter-expert variability), small ( e.g. $< 1\%$ of organ), or empty reference annotations, common metrics may penalize or misinterpret clinically meaningful information. Prior studies have described the importance of quantifying uncertainty in the reference annotation (Mehta et al.,

[**]Corresponding author:
*e-mail:* `sostm@stanford.edu` (Sophie Ostmeier)

1

**Fig. 1(a) Row 1: Uncertainty: Non-contrast Computed Tomography from an acute stroke patient within 16h. Row 2: Segmentation of all experts. The segmentations of all experts do not completely overlap.**



**Fig. 1(b) Small Volumes: Boxplot with volume distribution of all data sets**



**Fig. 1(c) Empty reference annotations: Row 1, Segmentation of all experts is "empty". Row 2, predicted voxel probabilities (softmax output values) of the models (low to high probability indicated by blue to red colors)). All colored pixels may be "false positives".**

2022), the dependency of metric values on the segmentation size and degree of class imbalance (Taha and Hanbury, 2015; Liu et al., 2021; Commowick et al., 2018), the equal weighting of all regions of misplaced delineation independently of their distance from the surface (Nikolov et al., 2018) or missing definition for empty reference annotations (Commowick et al., 2018; Maier-Hein et al., 2022).

The failure to describe uncertain, small, or empty segmentations may lead to irrelevant and misleading optimization and evaluation procedures in segmentation models.

Here, we determine how to implement clinically meaningful metrics for medical segmentation models with the UncertainSmallEmpty (USE)-Evaluator. We analyze the behavior of established metrics on benchmark deep learning models trained on four data sets with and without uncertain, small, and empty reference annotations (in-house and public).

## 1.1. Uncertain Reference Annotations

While experts annotations may exhibit variations in the volume and location of segmented objects, we assume that each expert possesses the highest level of human ability for the given task and, as a result, their judgments are considered equally valid (Jungo et al., 2018). Identifying a superior expert who can definitively determine the correctness among the experts would require someone with even higher human ability. However, the process of appointing such an overruling expert would necessitate another individual with even greater abilities to make this judgment, leading to an infinite loop. Consequently, in the context of our study, we assume the absence of an overruling expert.

For volume agreement, the reference annotation's classification of a voxel can be true, and the segmentation of another expert or the prediction of the models can be false or vice versa. In practice, the spectrum ranges from a worst-case to a best-case scenario. In the best-case scenario, all false positives ($FP$) are truly positives. In the worst-case scenario, all $FP$ are truly $FP$. For example, in Figure 1(a) the union annotation of an acute stroke from experts A, B, and C (blue, green, and red) is larger than the majority vote (green and red). Some blue voxels at the border of the segmentation might falsely or truly be part of acute stroke ($FP$ or $TP$). Another example is shown in Figure 1(c). Experts reference annotation is empty (first row). However, the prediction (second row) is not empty. Visual investigation shows an ambiguous lesion that was not segmented by the experts making all voxels $FP$ but maybe truly $TP$. The underlying low signal-to-noise ratio of stroke on Non-contrast Computed Tomography (NCCT) and the continuous transition from healthy to pathological brain tissue inevitably prevent a precise membership of these voxels.

For location agreement, the distance between voxels from the reference annotation to another expert or prediction might be longer or shorter. For example, in Figure 1(a) the surface voxels of the green voxels will have a different distance than the blue surface voxels to the surface voxels of a predicted segmentation.

In the BRATS 2019 data set, we reproduce an underlying low signal-to-noise ratio and a more continuous transition by using the non-enhancing tumor segmentation on native T1 as reference annotation. We compare to a high signal-to-noise ratio

setting with a more discrete transition by using the whole tumor segmentation on T1, T1 enhanced, T2-flair, and T2 MRI images and Spinal Cord white matter segmentation on T2 MR images.

We propose the Uncertainty score (U-score) as a quantifying measure of reference annotations uncertainty.

## 1.2. Small Reference Annotations

Depending on the clinical context, small reference annotations may be defined as relative to the total size of the studied body region. (e.g. less than 1%). For the brain, 1% is about 13 ml (Akeret et al., 2021). The distributions of reference annotation volumes vary across medical image data sets and segmentation tasks (Figure 1(b)) (Bakas et al., 2018, 2017; Menze et al., 2015; Prados et al., 2017).

We hypothesize that the distribution of reference annotation volumes influences the value of metrics independently from the model's performance (Figure 1(a))(Maier-Hein et al., 2022). For example, an acute ischemic stroke patient with a suspected large vessel occlusion undergoes emergent imaging to quantify the extent of the irreversible brain injury. The stroke volume is often quite small (1-5 ml in volume (Powers et al., 2018))). Models may segment a 1-2 ml lesion volume that has poor overlap with the segmentation by a neuroradiologist and have a low-performance metric despite properly identifying the volume. A slight difference in volume location within the brain is highly unlikely to influence a physician's decision to treat the patient.

We describe how the distribution of reference annotation volumes produces different metrics values, irrespective of the level of location and volume agreement between the model's predictions and the annotations.

## 1.3. Empty Reference Annotations

Images with empty reference annotations are described as masks in which the object of interest could not be identified by the annotators. The object might have been invisible at the time of the segmentation (Figure 1(c)).

Segmentation of an object within an image is a different task than a classification of an image. A classification task confirms the presence or absence of an object in the image (image-level), while a segmentation task assigns each voxel of the image to an object class (voxel-level) (Maier-Hein et al., 2022). An image-level classification task can also be formulated as a segmentation task by checking if the reference and predicted masks are empty. Therefore, when using a segmentation model in this way, it is important for the performance metrics to capture behavior on empty masks. However, some metrics for image segmentation return "NaN" or 0, if the model correctly predicts an empty mask (e.g. Dice, Specificity, Sensitivity, IoU).

For clinical deployment, images with empty reference annotation are possible and their presence is crucial information. The predictions of segmentation models need to be optimized and evaluated for correct image-level classification (Commowick et al., 2018). For example, it is possible that a stroke lesion in an early time window (0-4h after symptom onset) has a very low signal and cannot be segmented on NCCT. In this case, the reference annotation and the predicted segmentation should both be empty and an image-classification metric should return the optimal value. No visible and no predicted lesion would result in a treatment decision in favor of endovascular therapy (Powers et al., 2018).

We explore a potential solution by setting a volumetric threshold tailored to each clinical context where voxel-wise agreement is expected to go beyond clinical relevance. Below the threshold, the agreement between the reference annotation and prediction is automatically evaluated as an image-level classification task (e.g. stroke present or absent in the image USE-Evaluator).

## 1.4. Clinical Value

For a successful transition to clinical translatable challenge-winning segmentation models, the focus on clinically meaningful optimization and performance metrics for each clinical context is crucial. Clinical value includes:

- Robustness toward uncertainty in the reference annotation

- Independence from the reference volume

- Reward of volumetric and location agreement between the reference annotation and predictions

- Evaluation of correct classification of empty reference annotations and predictions

## 2. Metrics

### 2.1. Fundamentals

A 3D image consists of a voxel grid with width $w$, height $h$, and depth $d$. We refer to the set of voxels as $X$ with $|X| = w \times h \times d = n$.

A segmentation mask is a grid with the same shape as the image. Pixels/voxels are assigned integer values indicating the semantic class (e.g. organ, pathology) they belong to. In the context of this publication, segmentation masks are either created manually by human experts or, with an automatic algorithm from an image.

A mask can be evaluated by the volume and location agreement of the segmented object. On a voxel level, the agreement between the reference mask, $M$, and the predicted mask, $\hat{M}$, can be measured with (i) voxel class agreement or (ii) spatial distances between corresponding voxels.

For voxel class agreement, we use the assignment of voxels to classes, in the reference mask, as the true classes. The model's classification for each voxel results in a predicted mask. Let $K$ be the set of classes. We note that $K$ completely partitions the mask. That is, $M = \bigcup_{k \in K} M^k$ and $\hat{M} = \bigcup_{k \in K} \hat{M}^k$.

For a binary classification task ($k \in \{0, 1\}$) a confusion matrix of four cardinalities namely $TP$, $FP$, false negatives ($FN$), and true negatives ($TN$) can be defined, where $TP+FP+TN+FN = |X|$ (Table 2).

Table 1: Definitions of Performance Metrics for Medical Image Segmentation

| Category | Metric | Abbr | Usage | Definition |
|---|---|---|---|---|
| **Volume** | Volumetric Similarity | VS | (Caradu et al., 2021; de Vos et al., 2021; Tiulpin et al., 2020; Dewey et al., 2019; Vania et al., 2019) | $1 - \frac{|\hat{V}-V|}{\hat{V}+V+\epsilon}$ |
| | Absolute Volume Difference | AVD | (Amukotuwa et al., 2019; Brosch et al., 2018) | $\frac{1}{m}\sum_{j=1}^{m}\left|V_j - \hat{V}_j\right|$ |
| **Overlap** | Dice Similarity Coefficient | Dice | (Becker et al., 2019; Vania et al., 2019; Brosch et al., 2018) | $\frac{2\times TP}{2\times TP+FN+FP}$ |
| | Jaccard Index, Intersection over Union | IoU | (Bertels et al., 2019) | $\frac{TP}{TP+FN+FP}$ |
| | Recall = Sensitivity | Recall | (Vania et al., 2019) | $\frac{TP}{TP+FN}$ |
| | Precision | Precision | (Vania et al., 2019) | $\frac{TP}{TP+FP}$ |
| **Distance** | Hausdorff Distance, q = 95th percentile | HD 95 | (Huttenlocher et al., 1993; Kuijf et al., 2019; Litjens et al., 2014) | $\max\left(h_p(A,B), h_p(B,A)\right)$ with $h_p(A,B) = P^{th}_{a\in A} \min_{b\in B}\|b-a\|$ and $P = 95$ |
| | Average Symmetric Surface Distance | ASSD | (Heimann et al., 2009; Janssens et al., 2018; Styner et al., 2008) | $\frac{\sum_{x\in S}d(x,\hat{S})+\sum_{y\in\hat{S}}d(S,y)}{|S|+|\hat{S}|}$ |
| | Surface Dice at Tolerance | SDT | (Nikolov et al., 2018; Elguindi et al., 2019; Shusharina et al., 2020) | $\frac{|\hat{S}\cap B^t|+|S\cap\hat{B}^t|}{|\hat{S}|+|S|}$ |
| | Boundary IoU | BIoU | (Cheng et al., 2021) | $\frac{|(S\cap B^d)\cap(\hat{S}\cap\hat{B}^d)|}{|(S\cap B^d)\cup(\hat{S}\cap\hat{B}^d)|}$ |
| **Image-level classification[1]** | Accuracy[2] | ACC | (Gautam and Raman, 2021; Maier-Hein et al., 2022) | $\frac{TP_i+TN_i}{TP_i+TN_i+FP_i+FN_i}$ |
| | $F_1$-score (equivalent to Dice)[2] | $F_1$-score | (Gautam and Raman, 2021; Maier-Hein et al., 2022) | $\frac{2TP_i}{2TP_i+FP_i+FN_i}$ |
| | Sensitivity[2] | Sensitivity | (Maier-Hein et al., 2022) | $\frac{TP_i}{TP_i+FN_i}$ |
| | Specificity[2] | Specificity | (Maier-Hein et al., 2022) | $\frac{TN_i}{TN_i+FP_i}$ |
| | Area Under the Curve[2] | AUC | (Gautam and Raman, 2021; Maier-Hein et al., 2022) | $\int_0^{V_{max}}\frac{TP_i(V_{threshold})}{TP_i+TN_i+FP_i+FN_i}dV_{threshold}$ |

[1] with *threshold* [2] subscript $i$ indicates the image-level cardinalities

**Table 2: Voxel-Level Cardinalities for a binary classification task in which each voxel is assigned to one of the cardinalities depending on its mapping in the reference mask, $M^1$, and predicted mask, $\hat{M}$**

| | $\hat{M}^1$ | $\hat{M}^0$ |
|---|---|---|
| $M^1$ | $TP$ | $FN$ |
| $M^0$ | $FP$ | $TN$ |

The volume, $V$, of the target object in the reference mask and volume, $\hat{V}$, in the predicted mask is defined as

$$V = |M^1| \times v = (TP + FN) \times v \qquad (1)$$
$$\hat{V} = |\hat{M}^1| \times v = (TP + FP) \times v \qquad (2)$$

where $v$ refers to the physical volume of each voxel.

For distance agreement, the distance between a voxel $x$ and a set of surface voxels $S^k$ is defined as

$$d(x, S^k) = min_{s^k \in S^k} d(x, s^k).$$

For a binary segmentation task with $k \in \{0, 1\}$, the set of voxels, $S^1 \subseteq M$, is defined as the surface voxels of the target object in the reference mask, and $\hat{S}^1 \subseteq \hat{M}$ as the surface voxels of the target object in the predicted mask.

Metric definitions for common volume, overlap, and distance metrics, that were used for the experiments, can be found with their implementations on GitHub and in Table 1.

We note that the frequent class imbalance of 3D medical image segmentation ($|M^1| << |M^0|$) limits meaningful performance evaluation by any metric that includes $TN$ (Specificity, ROC, Accuracy, Kappa, etc.). Therefore, metrics that include

$TN$ in their function should be avoided. Overlap metrics measure $TP$ relative to a combination of $TP$, $FP$, and $FN$. Overlap metrics that exclude $TN$ are Dice, Recall, and Precision. Volume metrics without $TN$ are VS (Volumetric Similarity) and AVD (Absolute volume Difference).

We also note that the Jaccard index $J$ (Intersection over Union, IoU) and Dice $D$ are equivalent and one can be derived from one to the other using the following formula (Bertels et al., 2019).

$$J(M, \hat{M}) = \frac{D(M, \hat{M})}{2 - D(M, \hat{M})} \qquad (3)$$

$$D(M, \hat{M}) = \frac{2 \times J(M, \hat{M})}{1 + J(M, \hat{M})} \qquad (4)$$

The concrete choice for either one of these metrics depends on the user or community preference(Maier-Hein et al., 2022).

## 2.2. Surface Dice at Tolerance

Surface Dice is an evaluation metric introduced by (Nikolov et al., 2018). It describes which portion of voxels on the surface of the target object in the predicted mask have the same spatial location as the surface voxels in the reference mask. For that, it classifies the surface voxels into $TP$, $FP$, and $FN$ depending on their distance to the closest surface voxel in the reference/predicted mask. The contribution of individual voxels to these terms is weighted by the estimated surface area that it represents. The tolerable distance $t$ from the surface at which a voxel still counts as a TP establishes a set of border voxels $B^t$. $t$ is a variable that needs to be set according to the clinical context and (estimated) inter-rater variability for the given segmentation task.

$$SDT = \frac{\left|\hat{S}^k \cap B^{t,k}\right| + \left|S^k \cap \hat{B}^{t,k}\right|}{\left|\hat{S}^k\right| + \left|S^k\right|}.$$

$k$ is the target object class. The tolerated distance can be set depending on the task. A possible method to choose the tolerated distance is to compute the distance between different experts as an acceptable variability (e.g. ASSD Table 1). This might lead to an optimization procedure with a realistic tolerance in which uncertainty within the voxel classification of the reference mask is expected and acceptable.

The Boundary IoU with a distance $d$ proposed by (Cheng et al., 2021) can be converted to the Surface Dice at Tolerance with tolerated distance $t$ by using Equ. 3, where $d$ and $t$ are equivalent.

## 2.3. Uncertainty Score

We develop a score to estimate the uncertainty across a set $E$ of experts across the set $C$ of cases. This score may be used as an indicator of the uncertainty in the data set. Our score is built around evaluating entropy, a measure of information contained in samples, on a target region of each image.

We index cases using $c \in C$. We index the reference masks and membership functions by expert $e$ as $M_e^k$ and membership functions as $f_e^k(x)$. We consider the case where $k = \{0, 1\}$.

Our score will require counting over experts, classes, and voxels. Let $\beta^k(c, x)$ be the function that returns the number of experts that puts voxel $x$ of case $c$ in class $k$. Formally, $\beta^k(c, x) = \left|\{e | e \in E, f_e^k(c, x) = 1\}\right|$.

We compute the U-score over the set of voxels where at least one expert classified the voxel as positive. We denote this set as $U = \left(\bigcup_{e \in E} M_e^1\right)$.

For a case, we compute its U-score as the average, over voxels, of the expert annotation entropy of the voxel. Formally

$$\text{U-score} = \frac{1}{|U|} \sum_{x \in U} \text{entropy}(x). \qquad (5)$$

With entropy computed as

$$\text{entropy}(c, x) = \sum_k \frac{\beta^k(c, x)}{|E|} \log \frac{\beta^k(c, x)}{|E|}.$$

We can compute the U-score of a dataset as the average U-score over cases.

## 2.4. Voxel-level Class Imbalances

The class imbalance ratio ($IR$) is commonly defined as the ratio between the cardinality of the majority class and the cardinality of the minority class (Zhu et al., 2020).

### 2.4.1. Class Imbalances of Segmentation

In this context of image segmentation, the $IR$ is given by $IR = \frac{|M^{majority}|}{|M^{minority}|}$.

An image segmentation task can be a perfectly balanced voxel-level binary classification problem $|M^1| = |M^0|$. However, it is often the case that the target object is small relative to the image, that is $|M^1| << |M^0|$ and $\frac{|M^0|}{|M^1|} >> 1$. This indicates high class imbalance. This can result in even very simple models achieving many true negatives on $|M^0|$ with a low false positive rate (Table 2). Since the number of background voxels in medical images may vary (due to scanner settings, and image processing) we aim to control the considered background voxels in a consistent way. We do so by restricting the region of interest to either an organ or the immediate body cavity. Background voxels in this region of interest are referred to as $M^{0,region}$. For example, for stroke and brain tumor this would be the brain, for the gray matter in the spinal cord this would be the total spinal cord. We then get

$$IR = \frac{|M^{0,region}|}{|M^1|}. \qquad (6)$$

### 2.4.2. Image-level Class Imbalances

In the realm of image classification, we denote the class imbalance ratio as $IR_i$. When considering reference and predicted volumes that fall below a clinically reasonable *threshold* (i.e. 1ml for the NCCT and BRATS models), the significance of segmentation performance diminishes. Images that are correctly or incorrectly classified below this *threshold* are designated as

**Table 3: Data set properties**

| Data set | Target Object | Multiple Labels | USE[1] | Positive Cases[2] | Negative Cases[3] |
|---|---|---|---|---|---|
| NCCT[4] | ischemic core | ✓ | ✓ | ✓ | ✓ |
| BRATS 2019 | non-enhancing tumor | - | ✓ | ✓ | ✓ |
| BRATS 2019 | whole tumor | - | - | ✓ | - |
| Spinal Cord | gray matter | ✓ | - | ✓ | - |

[1] USE= uncertain, small and empty reference annotations,
[2] cases with the target object present in the image,
[3] case without the target object or below a volume threshold present in the image
[4] NCCT= Non-Contrast Computed Tomography

$TN_i$ or $FN_i$ respectively. As a result, we derive the equation:

$$IR_i = \frac{TP_i + FN_i}{TN_i + FP_i} \qquad (7)$$

with an optimal value of 1. Here, the positive cases (i.e. patients with a stroke larger than 1ml), represented by $TP_i + FN_i$, serve as the majority class, while the negative cases (i.e. patients with a stroke smaller than 1ml), represented by $TN_i + FP_i$, serve as the minority class. Visual examples illustrating $TN_i, FN_i, TP_i$, and $FP_i$ can be observed in Figure 2.

## 3. Methods

### 3.1. Data Sets

To evaluate metrics for models trained on uncertain, small, or empty reference annotations we use several data sets (Table 3).

A de-identified dataset of 200 NCCT images of patients with an acute ischemic stroke from the DEFUSE3 trial (Albers et al., 2018) was provided to three neuroradiologists 4, 4, and 5 years of experience in neuroradiology (B.V.,A.M.,J.J.H.) (study design https://clinicaltrials.gov/ct2/show/NCT02586415). The experts were instructed to segment abnormal hypodensity on the NCCT that corresponds to acute ischemic brain injury. Detailed instructions and videos, as well as an oral explanation of the task, were given. Any missed lesions or missed slices were not corrected. The experts' masks were fused by a majority vote to form the reference mask. In addition, 156 institutional NCCT images were added of patients who were scanned with suspicion of stroke but were confirmed on follow-up Diffusion-weighted MR imaging not to have a stroke.

The BRATS 2019 public data set was used to reproduce and compare results and included 345 MRIs of high and low-grade glioma patients (Bakas et al., 2018, 2017; Menze et al., 2015). One to four experts segmented the brain tumors followed by a

consensus procedure. The reference masks had four target objects; "background", "edema", "non-enhancing", and "enhancing". We used this data set to train two segmentation tasks (i) with the target object "non-enhancing" tumor on only T1 and (ii) with the target object whole tumor, defined as the union over "edema", "non-enhancing", and "enhancing" target objects, on T1, T1 contrast-enhanced, T2-Flair and T2.

The Spinal cord data set is a public data set with 40 annotated MRIs of 40 healthy patients from 4 different hospitals and annotated by 4 experts per case. The annotations include the white and gray matter of the spinal cord on T2(Prados et al., 2017). The experts' masks were fused by a majority vote to form the reference mask.

### 3.2. Data Partition

For each segmentation tasks the cases were randomly divided into 5 folds that consisted of 80% training and 20% test examples. The default self-configured nnUNet was used to train all folds for each segmentation task. All analyses were done on the aggregated 5 test sets for each segmentation task (Supplemental material, Figure 6).

All models shared the same training schedule with 500 epochs, Stochastic gradient descent with Nestov momentum of 0.99, the initial learning rate of 0.01 with linear decay, and oversampling of 33% for the target lesion.

### 3.3. Models

#### 3.3.1. Deep Learning models

We chose the 3D full-resolution nnUNet as our deep learning framework (Isensee et al., 2021). For fairness and ease of comparability, we let all models undergo the same training schedule and did not modify hyperparameters.

The default configured model included a patch size of ($1 \times 28 \times 512 \times 512$) and spacing of (3.00, 0.45, 0.45), Dice and Cross-Entropy loss function, seven stages, two 3D convolutions per stage and a leaky ReLU as activation function.

For the NCCT ischemic core segmentation task, the model input was the NCCT image (356 cases) to output a predicted mask for ischemic brain tissue.

For the first BRATS 2019 segmentation task, the model input encompasses only the T1 image (345 cases) to simulate a lower signal-to-noise ratio and the output was the predicted mask for the non-enhancing tumor. For the second BRATS 2019 segmentation task, the model input included all available MR sequences images (345 cases) to output the predicted mask of all tumor parts.

For the Spinal cord gray matter segmentation task (40 cases), the model input was a T2 image to output a predicted mask of gray matter.

#### 3.3.2. Random model

To demonstrate that this trend is independent of the model, we analyze the behavior of the Dice on a model that randomly labels voxels.

Our objective is to investigate the impact of target object size on the Dice score. As the class imbalance ratio increases, the Dice score tends to decrease, creating difficulties in comparing

**Fig. 2:** Example of true positives ($TP_i$), true negatives ($TN_i$), false positive ($FP_i$) and false negative ($FN_i$) cases for the NCCT and BRATS data set for a threshold of 1ml.

7

model performance across different levels of class imbalance. However, we aim to demonstrate that this relationship is also a general property not tied to a single task.

Consider a binary segmentation model that decides each voxel's membership in the predicted mask randomly with a biased coin toss . We will show this by deriving the expected Dice score $E_D$ for images drawn from a random model and showing that the trend observed empirically matches the trend in this theoretical model (Supplemental material 8.1.). It is cleaner to parameterize this random model using the expected portion of voxels that are positive. We refer to this as $p$ and note that it can be directly computed from the class imbalance ratio $p = \frac{1}{1+IR}$.

### 3.4. Evaluation Tool

All evaluations were performed using the USE-Evaluator inspired by (Nikolov et al., 2018; Isensee et al., 2021) (Table 1). The source code can be applied to folders with reference annotation and prediction mask in .nii.gz format and produces a .xslx file with sheets for all studies, the means, medians, and image-level classification with bootstrapped 95% confidence interval. A threshold flag can be set as a lower volume threshold for the segmentation and image-level classification evaluation. If the reference or predicted volume is below the threshold, a case is excluded from the segmentation evaluation but included as a negative case for the image-level classification evaluation.

### 3.5. Evaluation of Reference Annotations

We analyzed the variability among different experts' annotations masks, available for the NCCT and the Spinal cord data set, with the evaluation tool described in Section 3.4. To estimate uncertainty we compute the U-score (Equ. 5) and the median inter-expert agreement, and the median agreement to the majority vote (majority-expert) with the metrics presented in Table 1.

### 3.6. Evaluation of Model Performance

Performances were measured with the evaluation tool (Section 3.4) with a threshold of 1ml for the NCCT and BRATS 2019 data sets. For other medical applications, this might depend on the clinical task the model is trained on. With the evaluator tool, this threshold can be easily changed. For the Spinal cord data set, we did not set a threshold, because the clinical concern in healthy populations would not be about the non-existence vs. existence of gray matter in the spinal cord.

### 3.7. Evaluation of Metrics

We evaluate the segmentation metrics by correlation to uncertainty among the expert's masks, independence from reference volume, the reward of volumetric and location expert-model agreement, and evaluation of correct classification of empty reference masks or small reference volumes cases using the R package corrplot (Version 0.92).

To compute $IR$ and $p$ for the stylized model, we defined the *region*s as the entire brain for the NCCT and BRATS datasets, and the entire spinal cord for the Spinal Cord dataset (Section 2.4.1). The BET_CT was used to extract the brain on NCCT according to (Schell et al., 2019). For the extraction of the brain

on MRI, the HD_BET was applied (Isensee et al., 2019). For extraction of the spinal cord, the union of the gray and white matter in the majority vote reference mask was used.

For the evaluation of empty reference and predicted masks, we explore possible image-classification metrics and their relationship to $IR_i$, where we refer to $p_i = \frac{1}{1+IR_i}$.

## 4. Results and Discussion

In this section, we will examine the relationship between metric values and varying prevalence of uncertain, small, or empty reference annotations.

In Section 4.1 we measure uncertainty in reference annotations. We conduct empirical validation of the U-score across data sets and its correlation with inter-expert variability and consensus among the majority of experts.

In Section 4.2 we analyze all models' performances with each metric across data sets in order to provide a first indication of trends between dataset properties and metric values that we explore in further detail in the following section.

In Section 4.3 we use the correlation of metric values to provide empirical evidence of the link between the uncertain, small, and empty reference annotations and the metric values.

For the Dice metric, we demonstrate that the link is even more general by illustrating that the relationship found empirically is present in the evaluation of a stylized theoretical model (Section 4.3.2).

Finally, we explore trends in image-classification metrics in section 4.4.

Upon negative tests for normal distribution, results for each metric are shown as medians with 95% confidence interval (bootstrapped, 1000 repetitions), and the correlations are reported as Spearman's rank correlation coefficient.

### 4.1. Evaluation of Reference Annotations of Experts

Variability in reference annotations can impact the model's segmentation performance and solutions have been discussed (Karimi et al., 2020). However, we focus on a better choice of evaluation techniques to enhance the clinical applicability of segmentation models. In this regard, we first propose the introduction of the U-score as a measure of uncertainty for reference annotations (Section 1.1). We found an overall median U-score is significantly different between the NCCT ischemic core and the Spinal cord gray matter segmentation task (0.87 ± 0.05 vs. 0.39 ± 0.02, respectively). These findings are consistent with common measures such as inter-expert and majority-expert agreement (supplemental material, Table 7)(Yang et al., 2023). Inter-expert and majority-expert agreements use pairwise expert comparison and rely on common segmentation metrics to indirectly estimate uncertainty in reference annotations. The U-score directly measures uncertainty.

We found varying distributions of reference volumes across the studied data sets (median (IQR) volume 6 (2-21)ml, 10 (4-25)ml, 89 (48-146)ml and 0.7 (0.3-1.1)ml for NCCT, BRATS 2019 non-enhancing tumor part segmentation task, BRATS 2019 whole tumor segmentation task and Spinal Cord gray matter segmentation, respectively.

**Fig. 3: Scatter plot with log-scale and confusion matrix with a volume threshold of 1ml dividing $TP$ and $TN$ from $FP$ and $FN$. For the NCCT data set(violet points), almost all incorrectly classified cases are too small, namely $FN$, whereas for the BRATS non-enhancing tumor data set the opposite is the case. None of the cases of BRATS whole tumor are incorrectly classified.**

We further conduct correlation analyses between the U-score and reference volumes to common metrics outlined in Section 4.3.

### 4.2. Evaluation of Segmentation and Image Classification Performance

The performance of the NCCT ischemic core and BRATS 2019 non-enhancing tumor models, trained using uncertain, small, and empty reference annotations, shows similar results across volume, overlap, and distance metrics. However, the BRATS 2019 whole tumor and Spinal cord model, trained on larger and more certain reference annotations, consistently outperforms both the NCCT and BRATS 2019 non-enhancing tumor model (Table 4).

For image classification, the total number of cases with reference volumes < 1ml is 192 cases for the NCCT ischemic core segmentation task, 36 cases for the BRATS 2019 non-enhancing tumor part segmentation task and 0 cases for the BRATS 2019 whole tumor . We visualize these class distributions and report the confusion matrix (Figure 3)(Maier-Hein et al., 2022). Reference Volumes >1ml cluster around the identity line, whereas references <1ml are more spread. Sensitivity, F1-score and ACC of the NCCT model are lower compared to the BRATS non-enhancing tumor models, however, the AUC and Specificity are higher for the NCCT models (Table 5). Further analysis between data set properties and common image classification metrics analysis are summarized in Section 4.4.

### 4.3. Evaluation of Segmentation Metrics

We use the relationship between metrics and dataset properties to identify evaluation strategies robust of the presence to uncertain, small or empty reference annotations. These recommendations are backed by the empirical data (Figure5) and we



**Fig. 4: Dot plot with regression lines for the Dice over class imbalance $p$ for all segmentation models, where $p = \frac{1}{1+IR}$. The gray areas represent 95% confidence intervals. The dark red dots and line represent the random model with the expected Dice $E_D$ defined here. The dashed line indicates the expected Dice $E_D$ for a balanced reference mask.**

provide an intuition of how the given formula provides the observed effect (Table1).

We categorize the segmentation metrics according to volume, overlap, or distance agreement. We then analyze every segmentation metric based on the characteristics outlined in section 1.4. If not otherwise specified, all numbers in this section refer to the Spearman correlation coefficients presented in Figure 5. A summary of the core results and guidelines for the choice of metrics are provided in Table 6.

#### 4.3.1. Volume Agreement
**VS:**

*Robustness toward Uncertainty and Independence from Reference Volume:*

Conceptually, VS allows location variability of reference volumes (Table 6), because $FN$ and $FP$ voxels can be anywhere in the image without an influence on the value of VS. This characteristic becomes particularly valuable when dealing with uncertain reference annotations. Assuming that a source of $FP$ and $FN$ is uncertainty (see Section 1.1); VS does not penalize uncertainty as long as their difference has a linear relationship to reference volumes. This is because VS normalizes to the sum of the reference and predicted volume. Our findings support this with a low correlation to uncertainty (-0.17 and -0.32 for NCCT and Spinal Cord) and reference volume (across all data set below 0.25). We conclude that VS value is less driven by uncertainty or reference volume.

*Reward of Volume and Location Agreement:*

VS does not reward location agreement since VS measures the relative relationship between $FP$ and $FN$ rather than their

9

**Fig. 5(a)** Uncertain, small and empty reference annotations (NCCT stroke)

| | V | V̂ | VS | AVD | Dice | Precision | Recall | ASSD | HD 95 | SDT_small | SDT_large | Uncertainty |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| V | | 0.82 | 0.13 | 0.84 | 0.56 | 0.60 | 0.63 | X | X | X | X | −0.72 |
| V̂ | 0.82 | | 0.47 | 0.49 | 0.85 | 0.54 | 0.89 | −0.49 | −0.43 | 0.58 | 0.57 | −0.66 |
| VS | 0.13 | 0.47 | | X | 0.84 | | 0.58 | −0.75 | −0.73 | 0.82 | 0.80 | −0.17 |
| AVD | 0.84 | 0.49 | X | | X | 0.49 | X | X | 0.31 | X | X | −0.53 |
| Dice | 0.56 | 0.85 | 0.84 | X | | 0.37 | 0.96 | −0.82 | −0.74 | 0.86 | 0.84 | −0.63 |
| Precision | 0.60 | 0.54 | | 0.49 | 0.37 | | 0.51 | −0.24 | −0.15 | 0.19 | 0.20 | −0.59 |
| Recall | 0.63 | 0.89 | 0.58 | X | 0.96 | 0.51 | | −0.81 | −0.74 | 0.86 | 0.84 | −0.60 |
| ASSD | X | −0.49 | −0.75 | X | −0.82 | −0.24 | −0.81 | | 0.95 | −0.94 | −0.96 | X |
| HD 95 | X | −0.43 | −0.73 | 0.31 | −0.74 | −0.15 | −0.74 | 0.95 | | −0.85 | −0.92 | X |
| SDT_small | X | 0.58 | 0.82 | X | 0.86 | 0.19 | 0.86 | −0.94 | −0.85 | | 0.96 | −0.37 |
| SDT_large | X | 0.57 | 0.80 | X | 0.84 | 0.20 | 0.84 | −0.96 | −0.92 | 0.96 | | −0.37 |
| Uncertainty | −0.72 | −0.66 | −0.17 | −0.53 | −0.63 | −0.59 | −0.60 | X | X | −0.37 | −0.37 | |

**Fig. 5(a) Uncertain, small and empty reference annotations (NCCT stroke): The overlap metrics (Dice, Recall and Precision) and AVD have a strong negative correlation to Uncertainty and positive correlation to Reference Volume, Distance metrics (ASSD, HD 95, SDT) show insignificant correlation, respectively.**

**Fig. 5(b)** Uncertain, small and empty reference annotations (BRATS 2019 non-enhancing tumor)

| | V | V̂ | VS | AVD | Dice | Precision | Recall | ASSD | HD 95 | SDT_small | SDT_large | Uncertainty |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| V | | 0.76 | 0.32 | 0.62 | 0.47 | 0.62 | 0.22 | X | X | X | X | ? |
| V̂ | 0.76 | | 0.39 | 0.56 | 0.46 | 0.31 | 0.55 | X | X | X | 0.12 | ? |
| VS | 0.32 | 0.39 | | −0.42 | 0.76 | 0.38 | 0.56 | −0.69 | −0.58 | 0.75 | 0.66 | ? |
| AVD | 0.62 | 0.56 | −0.42 | | −0.14 | 0.20 | X | 0.47 | 0.43 | −0.47 | −0.47 | ? |
| Dice | 0.47 | 0.46 | 0.76 | −0.14 | | 0.66 | 0.73 | −0.76 | −0.63 | 0.80 | 0.71 | ? |
| Precision | 0.62 | 0.31 | 0.38 | 0.20 | 0.66 | | 0.26 | −0.34 | −0.18 | 0.31 | 0.23 | ? |
| Recall | 0.22 | 0.55 | 0.56 | X | 0.73 | 0.26 | | −0.64 | −0.62 | 0.69 | 0.67 | ? |
| ASSD | X | X | −0.69 | 0.47 | −0.76 | −0.34 | −0.64 | | 0.91 | −0.93 | −0.93 | ? |
| HD 95 | X | X | −0.58 | 0.43 | −0.63 | −0.18 | −0.62 | 0.91 | | −0.81 | −0.94 | ? |
| SDT_small | X | X | 0.75 | −0.47 | 0.80 | 0.31 | 0.69 | −0.93 | −0.81 | | 0.91 | ? |
| SDT_large | X | 0.12 | 0.66 | −0.47 | 0.71 | 0.23 | 0.67 | −0.93 | −0.94 | 0.91 | | ? |
| Uncertainty | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | |

**Fig. 5(b) Uncertain, small and empty reference annotations (BRATS 2019 non-enhancing tumor): Like in (a), overlap metrics (Dice, Recall, and Precision) and AVD demonstrate a strong positive correlation with the reference volume (V), while distance metrics (ASSD, HD 95, SDT) exhibit insignificant correlation.**

**Fig. 5(c)** Certain and Large reference annotations (BRATS 2019 whole tumor)

| | V | V̂ | VS | AVD | Dice | Precision | Recall | ASSD | HD 95 | SDT_small | SDT_large | Uncertainty |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| V | | 0.98 | 0.21 | 0.38 | 0.42 | 0.35 | 0.17 | −0.12 | −0.13 | 0.12 | 0.08 | ? |
| V̂ | 0.98 | | 0.24 | 0.34 | 0.45 | 0.25 | 0.29 | −0.16 | −0.16 | 0.17 | 0.12 | ? |
| VS | 0.21 | 0.24 | | −0.79 | 0.68 | 0.13 | 0.37 | −0.65 | −0.63 | 0.67 | 0.52 | ? |
| AVD | 0.38 | 0.34 | −0.79 | | −0.40 | 0.05 | −0.24 | 0.55 | 0.52 | −0.56 | −0.46 | ? |
| Dice | 0.42 | 0.45 | 0.68 | −0.40 | | 0.48 | 0.64 | −0.86 | −0.83 | 0.86 | 0.70 | ? |
| Precision | 0.35 | 0.25 | 0.13 | 0.05 | 0.48 | | X | −0.33 | −0.30 | 0.27 | 0.21 | ? |
| Recall | 0.17 | 0.29 | 0.37 | −0.24 | 0.64 | X | | −0.60 | −0.62 | 0.64 | 0.57 | ? |
| ASSD | −0.12 | −0.16 | −0.65 | 0.55 | −0.86 | −0.33 | −0.60 | | 0.96 | −0.95 | −0.86 | ? |
| HD 95 | −0.13 | −0.16 | −0.63 | 0.52 | −0.83 | −0.30 | −0.62 | 0.96 | | −0.96 | −0.90 | ? |
| SDT_small | 0.12 | 0.17 | 0.67 | −0.56 | 0.86 | 0.27 | 0.64 | −0.95 | −0.96 | | 0.88 | ? |
| SDT_large | 0.08 | 0.12 | 0.52 | −0.46 | 0.70 | 0.21 | 0.57 | −0.86 | −0.90 | 0.88 | | ? |
| Uncertainty | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | |

**Fig. 5(c) Certain and Large reference annotations (BRATS 2019 whole tumor): Low correlation between the volumes to metric values.**

**Fig. 5(d)** Certain and small reference annotations (Spinal Cord gray matter)

| | V | V̂ | VS | AVD | Dice | Precision | Recall | ASSD | HD 95 | SDT_small | SDT_large | Uncertainty |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| V | | 0.94 | X | 0.48 | 0.40 | X | X | X | X | X | X | X |
| V̂ | 0.94 | | 0.51 | X | 0.45 | X | 0.40 | −0.24 | −0.44 | 0.37 | 0.37 | X |
| VS | X | 0.51 | | −0.65 | 0.58 | −0.59 | 0.71 | −0.44 | −0.43 | 0.57 | 0.57 | −0.32 |
| AVD | 0.48 | X | −0.65 | | X | 0.62 | −0.48 | X | X | −0.33 | −0.33 | X |
| Dice | 0.40 | 0.45 | 0.58 | X | | X | 0.76 | −0.92 | −0.58 | 0.87 | 0.87 | −0.72 |
| Precision | X | X | −0.59 | 0.62 | X | | −0.57 | X | X | X | X | X |
| Recall | X | 0.40 | 0.71 | −0.48 | 0.76 | −0.57 | | −0.80 | −0.58 | 0.73 | 0.73 | −0.46 |
| ASSD | X | −0.24 | −0.44 | X | −0.92 | X | −0.80 | | 0.60 | −0.85 | −0.85 | 0.70 |
| HD 95 | X | −0.44 | −0.43 | X | −0.58 | X | −0.58 | 0.60 | | −0.42 | −0.42 | 0.37 |
| SDT_small | X | 0.37 | 0.57 | −0.33 | 0.87 | X | 0.73 | −0.85 | −0.42 | | 1.00 | −0.79 |
| SDT_large | X | 0.37 | 0.57 | −0.33 | 0.87 | X | 0.73 | −0.85 | −0.42 | 1.00 | | −0.79 |
| Uncertainty | X | X | −0.32 | X | −0.72 | X | −0.46 | 0.70 | 0.37 | −0.79 | −0.79 | |

**Fig. 5(d) Certain and small reference annotations (Spinal Cord gray matter): With the exception of Dice and AVD, there is a lack of correlation between the metrics and the reference volumes. Dice and SDT exhibit a high correlation with uncertainty.**

**Fig. 5: Correlation matrices of Spearman coefficient for data sets and metrics. X indicates insignificant correlations with $p > 0.05$. Overall correlation patterns among metrics (e.g. Dice and SDT) remain similar over the data sets. The correlation between Dice and uncertainty, as well as the reference volume, is reproducible in all datasets, albeit to varying degrees.**

**Table 4: Results of Segmentation Task Performance** [1]

| Categories | Metrics[2] | NCCT Ischemic core | | BRATS 2019 Non-enhancing tumor | | BRATS 2019 Whole tumor | | Spinal Cord Gray matter | |
|---|---|---|---|---|---|---|---|---|---|
| **Volume** | **VS** | 0.58 | ± 0.09 | 0.78 | ± 0.03 | 0.97 | ± 0 | 0.92 | ± 0.02 |
| | **AVD** | 4.48 | ± 1.17 | 4.41 | ± 0.95 | 4.95 | ± 0.85 | 0.08 | ± 0.04 |
| **Overlap** | **Dice** | 0.56 | ± 0.04 | 0.60 | ± 0.04 | 0.93 | ± 0.01 | 0.83 | ± 0.02 |
| | **Precision** | 0.69 | ± 0.11 | 0.66 | ± 0.06 | 0.94 | ± 0.01 | 0.89 | ± 0.03 |
| | **Recall** | 0.20 | ± 0.08 | 0.58 | ± 0.05 | 0.93 | ± 0.01 | 0.79 | ± 0.03 |
| **Distance** | **ASSD** | 2.34 | ± 0.41 | 2.50 | ± 0.15 | 0.94 | ± 0.07 | 0.12 | ± 0.02 |
| | **HD 95** | 8.30 | ± 1.51 | 8.00 | ± 0.55 | 2.87 | ± 0.3 | 0.50 | ± 0.05 |
| | **SDT small** [3] | 0.61 | ± 0.05 | 0.56 | ± 0.03 | 0.93 | ± 0.01 | 0.84 | ± 0.08 |
| | **SDT large** [4] | 0.86 | ± 0.03 | 0.85 | ± 0.02 | 0.99 | ± 0 | 0.84 | ± 0.08 |

[1] median ± 95% Confidence Interval (bootstrapped)
[2] VS = Volumetric Similarity, AVD = Absolute Volume Difference, ASSD = Average Surface Distance, HD 95 = Hausdorff Distance 95th percentile, SDT = Surface Dice at Tolerance
[3] Surface Dice at Tolerance with 2mm for NCCT and BRATS 2019 models and 0.05mm for the Spinal Cord model
[4] Surface Dice at Tolerance with 5mm for NCCT and BRATS 2019 models and 0.1mm for the Spinal Cord model

**Table 5: Results of Image-Classification Task**[1]

| Categories | Metrics[2] | NCCT Ischemic core | | BRATS 2019 Non-enhancing tumor | | BRATS 2019 Whole tumor | | Spinal Cord Gray matter[4] |
|---|---|---|---|---|---|---|---|---|
| **Class imbalance** | $p_i$ | 0.46 | | 0.89 | | 1.00 | | - |
| **Image-level** | **Sensitivity** | 0.67 | ± 0.04 | 0.93 | ± 0.01 | 1.00 | ± 0 | - |
| | **Specificity** | 0.98 | ± 0.01 | 0.53 | ± 0.09 | | | - |
| | **$F_1$-score** | 0.79 | ± 0.03 | 0.94 | ± 0.01 | 1.00 | ± 0 | - |
| | **ACC** | 0.84 | ± 0.02 | 0.89 | ± 0.02 | 1.00 | ± 0 | - |
| | **AUC** | 0.91 | ± 0.02 | 0.86 | ± 0.03 | | | - |

[1] median ± 95% Confidence Interval (bootstrapped)
[2] 1ml threshold,
[3] ACC=Accuracy, AUC=Area under the Curve
[4] healthy cohort, no threshold for pathology set

distance. VS is therefore suitable if volume agreement is the major clinical concern, as for some applications in neuroimaging, like stroke (Powers et al., 2018).

In theory, VS may be less appropriate for clinical datasets and segmentation tasks that heavily rely on spatial information, such as those involving multiple sclerosis (Filippi et al., 2019). However, our observations indicate a consistent moderate to strong correlation with overlap metrics (e.g., Dice coefficient ranging from 0.58 to 0.84) and distance metrics (e.g., SDT$_{small}$ ranging from 0.57 to 0.82), particularly in datasets where reference annotations are uncertain and small in size.

*Reward of Agreement of Emptiness:*

For cases with empty references and predicted masks, VS returns the optimal value of 1. Therefore, VS is suitable for data sets with expected empty reference masks. Nevertheless, we recommend setting a threshold for very small volumes (<1ml), because the frequency of empty reference or predicted masks could screw the distribution of values compared to other metrics.

**AVD:**

Our findings suggest, no advantage of AVD over VS for data sets with a small median of reference volumes and uncertainty.

In contrast to VS, AVD does not normalize to the sum of reference and predicted volumes. Larger reference volumes have potentially larger volume differences, resulting in a notably positive correlation between AVD and reference volumes across all data sets. In datasets with a wide spread of reference volumes (Figure 1(b)), it is unclear whether a reduction of AVD as a metric leads to slightly improved performance for large reference volumes or substantially for small reference volumes. This ambiguity can introduce bias when comparing model performance within and across datasets, as evidenced by inconsistent correlation patterns with overlap and distance agreement metrics in our correlation analysis.

### 4.3.2. Overlap Agreement

**Dice:**

*Robustness toward Uncertainty:*

We observed that the Dice correlates more with uncertainty compared to other metrics (-0.62 to -0.72). This indicates that the Dice value is influenced not only by the extent of overlap but also by the level of uncertainty.

In a theoretical context, let's consider two scenarios. In the best-case scenario, a model outperforms the experts (as determined by the majority vote of reference annotation) by correctly classifying voxels. In this ideal situation, all $FP$ are $TP$, and all $FN$ are $TN$. However, the denominator contains the sum of $|M^1| = TP + FN$ and $|\hat{M}^1| = TP + FP$ (Table 1) and would disproportionately increase and lead to a lower Dice value. As a result, the performance of the models is underestimated. In the worst-case scenario, a model is inferior to the experts in classifying voxels correctly; all $FP$ are truly $FP$ and all $FN$ are truly $FN$. The Dice value does not change. As a result, Dice is biased toward the worst-case scenario. Hence, the Dice over-penalizes overlap disagreement in the presence of uncertainty between the experts' masks with a lower value.

*Independence from Reference Volume:*

In our study, we consistently observed a positive correlation between the reference volume and the $IR$ across all datasets, ranging from 0.40 to 0.56, with the NCCT dataset exhibiting the highest correlation. We hypothesized that the size of the target object affects the Dice value. More specifically, we investigated how the $IR$ impacts the likelihood of a voxel being classified as $TP$ because the Dice primarily rewards accurate voxel assignment to $TP$ ($\frac{2TP}{2TP+FP+FN}$). To validate our hypothesis, we analyzed the Dice value on a random model using the parameter $p$ (Section 2.4.1). The value of $p$ can be directly calculated from the $IR$ and represents the probability of a voxel in the prediction mask being classified as belonging to the target object class (see Section 3.3.2). We plot the Dice curve of the random model (dark red line) and compared it to all data sets (Figure 4). If $p$ is very low at 0.01 (1% of the brain), then the expected Dice of the random model is 0.02. If $p$ is 0.5 (50% of the brain), the expected Dice is much higher at 0.5 (dashed line). The regression lines for the random model, NCCT, and BRATS non-enhancing tumor show a positive monotonic tendency of the Dice values with higher $p$. This behavior is also present in BRATS 2019 whole tumor and Spinal cord models with larger $p$, but less (shallower slope of gray and orange lines).

We infer that a high imbalance ratio is more likely to produce lower Dice values. Location and volume errors for small reference annotations may be more penalized than larger reference annotations, making the Dice a sub-optimal choice of metric for data sets with small reference annotations and a wide distribution of reference volumes.

*Reward of Volume and Location Agreement:*

The numerator of the Dice, which comprises $2TP$, represents the voxels assigned to both the reference and prediction masks. The maximization of this value occurs when there is a high agreement in terms of both location and volume between the masks. We empirically see that the Dice rewards of volume and location agreement with a consistent, moderate to strong corre-

lation to VS and distance metrics across all data sets.

*Reward of Agreement of Emptiness:*

The Dice does not reward the agreement of emptiness between the reference and predicted mask, but returns "NaN". We found a high number of cases in the NCCT data set with a Dice value of 0. Investigation showed, that the Dice is zero if target objects are right next to each other and also zero if they are far from each other, especially for small reference volumes. This may lead to a disproportionate count of cases with Dice equal to zero. Depending on the clinical context, very small reference volumes (i.e. <1 ml) may be excluded from the evaluation of segmentation metrics. This is done to avoid introducing bias to the overall performance without obtaining meaningful information.

Instead, we suggest image-classification metrics to evaluate very small reference volumes or empty reference annotations masks. For example, a case with $V < 1$ml may be better evaluated by image classification metrics than by a segmentation metric. We implemented this idea with the USE-Evaluator, where a lower volume threshold can be set that will exclude studies with $V < threshold$ and automatically initializes an image-classification evaluation.

**Recall and Precision:**

Overall, Recall and Precision show similar behavior compared to the Dice, but only capture certain aspects of overlap agreement, and should be evaluated with other segmentation metrics and in the context of the clinical question.

They differ in their consideration of $FP$ and $FN$ in the denominator. Precision rewards $TP$ relative to the predicted volume, $TP + FP = |\hat{M}^1|$, and Recall $TP$ relative to the reference volume, $TP + FN = |M^1|$ (Table 1).

Especially Recall showed a correlation to uncertainty in the NCCT and Spinal cord data set (-0.60 and -0.43). One can argue that in the setting of high-class imbalance, the models learn to classify voxels with high entropy less frequently to $|\hat{M}^1|$, because the chance of being correct if classified to $|\hat{M}^0|$ is higher, increasing $FN$ (Leevy et al., 2018). We then get a higher denominator for Recall, thus an underestimation of uncertain reference volumes.

Similarly to the Dice, Recall, and Precision do not reward the agreement of emptiness. We, therefore, recommend setting a threshold for very small reference volumes and evaluating such cases with image-classification metrics.

### 4.3.3. Distance Agreement

Overall, we found that distance metrics, especially SDT, show favorable behavior in the context of small and uncertain reference annotations, while still exhibiting a consistent correlation to metrics that measure volume and overlap agreement.

**SDT:**

*Robustness toward Uncertainty and Independence from Reference Volume:*

SDT assigns cardinalities to surface voxels based on their proximity to the nearest surface voxel in either the reference or predicted mask. This approach emulates the behavior of the Dice while serving as a distance metric. However, contrary to Dice, if the reference and predicted volumes are right

next to each other and within the border region $\hat{B}^t$, SDT still measures this agreement. This becomes particularly advantageous when a lower signal caused by pathophysiological factors and modality-related effects introduces more uncertainty in the outer regions of the target object compared to its inner regions, i.e. like a stroke on NCCT. Compared to the Dice, we found weaker correlations to both the U-score and reference volume for the NCCT data set (-0.37, respectively).

In the Spinal cord data set, there is a correlation between SDT to the U-score. Image analysis of a few distinct cases with high uncertainty and low SDT value revealed deteriorating image quality in the cranial and caudal slices of the spinal cord, which is suggested to be the primary source of this relationship.

*Reward of Volume and Location Agreement:*
SDT shares similarities with overlap measures, due to its reliance on the spatial relationships among surface voxels and the direct influence of the object size on $|B^t|$. Consequently, SDT captures both the agreement in location and volume, which is further supported by its strong correlation with volume and distance metrics across all data sets (0.52-0.87).

*Reward of Agreement of Emptiness:*
In the presence of empty reference masks, all distance metrics return "inf". Similarly, to overlap metrics, distance metrics may need a lower bound volume threshold to evaluate empty and small volume reference masks with image-classification metrics.

**HD 95 and ASSD:**

Given that HD 95 and ASSD are metrics based on distance (as discussed in Section 2.1), which implies that if the model predicts a slightly different volume, HD 95 and ASSD should still yield values close to an optimal result, primarily capturing location accuracy and allowing volume error.

Since HD 95 and ASSD are distance-based metrics (as explained in Section 2.1), they primarily assess location accuracy while accommodating for volume errors. Values are close to the optimal even if the model's predicted volume slightly deviates from the reference.

Consistent with this, we found that HD 95 and ASSD exhibited mostly no correlations with reference volumes and uncertainty. However, similar to SDT, we observed a correlation between ASSD and the U-score in the Spinal Cord data set, likely attributed to low image quality in the cranial and caudal slices in distinct cases.

Overall, HD 95 and ASSD show robustness to uncertainty and reference volume, however, mostly measure distance agreement. Even though they empirically show strong correlations with volume and overlap metrics across all datasets, SDT should be preferred as a metric, if volume and location agreement is crucial.

*4.4. Evaluation of Image-level classification Metrics*

We propose a simultaneous evaluation of image-level classification metrics besides segmentation metrics to ensure an unbias evaluation of model performance when trained on data sets that include cases with uncertain, small, or empty reference annotations. Negligible reference volumes below a certain *threshold* may only be evaluated with image classification metrics

For example, Liu et al. only included positive cases and proposed image classification metric, LDR (Liu et al., 2021). However, the agreement in image-level classification is not assessed in the case of empty reference masks or small-volume cases. Clinical tests are unlikely to exclusively be performed on patients with a present pathology.

Data sets with negative cases or cases with negligible reference volumes would be more representative of the distribution of patients in clinical practice. This can have major implications for the idea of mostly positive or ambiguous cases being read by a radiologist and negative cases confidently evaluated by an algorithm (Wang et al., 2021).

In this section, we briefly highlight how a class imbalance between positive and negative/small-volume cases also introduces evaluation biases for inter-models and inter-data set comparison.

**Sensitivity, Specificity, and F$_1$-score:** Whether to use Specificity or Sensitivity as the primary image-classification metric depends on the clinical context. For this study, we found higher Sensitivity, Specificity, and F$_1$-score associated with higher $p_i$ (Supplemental material Figure 7). However, further studies are needed for more general statements.

**ACC:** The ACC evaluates the agreement in image-level classification in the case of $TN_i$ and $TP_i$ (Maier-Hein et al., 2022). As for Sensitivity, Specificity, and F$_1$-score, we found that a higher ACC value is associated with higher $p_i$ (Supplemental material, Figure 7).

**AUC:** AUC is a standard multi-threshold classification metric to evaluate a predictor and is not defined for populations where only one class is present (therefore only NCCT and BRATS non-enhancing tumor) (Maier-Hein et al., 2022). As the true discrete class in this setting is defined by the volume *threshold*, the AUC reveals information on how well the models classify volumes. We found that AUC does not change with $p_i$, suggesting AUC is a more robust metric for unbalanced data sets.

## 5. Limitations

The first limitation is that we evaluate metrics behavior and reference annotation uncertainty in only three medical neuroimaging data sets and examine four different target objects. We introduce methodologies aimed to be applied to a broader range of medical imaging data sets, allowing for a comprehensive examination of our findings. The second limitation is that the choice of baseline models might influence the correlations between metrics. In order to mitigate this, we choose nnUNet, a model that is generalizable to many medical segmentation tasks (Isensee et al., 2021). Furthermore, correlation does not prove causation. For example, the correlation between reference volumes and uncertainty to the value of metrics does not imply that a higher reference volume value causes a higher metric value. The correlation of Dice and reference volumes have been found in previous works (Taha and Hanbury, 2015; Liu et al.,

2021; Commowick et al., 2018; Maier-Hein et al., 2022) however analysis of data sets properties, in-depth analysis, quantification of the uncertainty were missing.

## 6. Conclusion

We notice a mismatch between dataset properties in challenge-winning segmentation models and cases encountered in clinical practice. Some commonly used metrics (i.e. Dice score) might not capture whether models' performance generalize well to the distribution of images encountered in clinical practice. In particular, (i) the presence of uncertainty in reference annotations causes misleading values, (ii) small reference volumes lead to unreasonable low metric values, (iii) empty reference annotations cause a return of "NaN", "inf" or zero. For a data set with uncertain, small, and empty reference annotations, we suggest that model performance generalizes better to clinical practice when evaluated by the Surface Dice at Tolerance. We further proposed to set a lower volume threshold for very small volumes or empty reference masks and use image-level classification metrics such as AUC ( USE-Evaluator).

It is crucial to evaluate the performance of the model using multiple metrics that effectively encompass the specific objectives of the clinical segmentation task. These objectives can vary significantly across different areas of clinical practice. To facilitate the selection of appropriate metrics, we recommend referring to Table 6.

We highlight the difficulty of comparing models trained to address different clinical problems. While uncertain, small, and empty reference annotations require a rethinking of evaluation, it also increases the value an algorithmic tool provides because the underlying task is hard for human experts.

## 7. Acknowledgements

**References**

Akeret, K., Bas van Niftrik, C.H., Sebök, M., Muscas, G., Visser, T., Staartjes, V.E., Marinoni, F., Serra, C., Regli, L., Krayenbühl, N., Piccirelli, M., Fierstra, J., 2021. Topographic volume-standardization atlas of the human brain. medRxiv , 2021.02.26.21251901URL: `https://www.medrxiv.org/content/medrxiv/early/2021/03/01/2021.02.26.21251901.full.pdf`, doi:10.1101/2021.02.26.21251901.

Albers, G.W., Marks, M.P., Kemp, S., Christensen, S., Tsai, J.P., Ortega-Gutierrez, S., McTaggart, R.A., Torbey, M.T., Kim-Tenser, M., Leslie-Mazwi, T., Sarraj, A., Kasner, S.E., Ansari, S.A., Yeatts, S.D., Hamilton, S., Mlynash, M., Heit, J.J., Zaharchuk, G., Kim, S., Carrozzella, J., Palesch, Y.Y., Demchuk, A.M., Bammer, R., Lavori, P.W., Broderick, J.P., Lansberg, M.G., 2018. Thrombectomy for stroke at 6 to 16 hours with selection by perfusion imaging. New England Journal of Medicine 378, 708–718. URL: `https://www.nejm.org/doi/full/10.1056/NEJMoa1713973https://www.nejm.org/doi/pdf/10.1056/NEJMoa1713973?articleTools=true`, doi:10.1056/NEJMoa1713973.

**Table 6: Suggestions for Choosing Meaningful Metrics for Data Sets with Uncertain, Small and Empty Reference Annotation**

| Category | Metric[1] | Robustness toward Uncertainty in Reference Annotation | Independence from Volume of Reference Annotation | Reward of Volume and Location Agreement | Reward of Agreement of Emptiness |
|---|---|---|---|---|---|
| **Volume** | **VS** | ✓ | ✓ | - | ✓ |
| | **AVD** | ✓ | - | - | ✓ |
| **Overlap** | **Dice** | - | - | ✓ | - set $threshold^2$ |
| | **Recall** | - | - | ✓ | - set $threshold^2$ |
| | **Precision** | - | - | ✓ | - set $threshold^2$ |
| **Distance** | **HD 95** | ✓ | ✓ | - | - set $threshold^2$ |
| | **ASSD** | (✓) | - | ✓ | - set $threshold^2$ |
| | **SDT small** | ✓ | ✓ | ✓ | - set $threshold^2$ |
| | **SDT large** | ✓ | ✓ | ✓ | - set $threshold^2$ |

[1] VS = Volumetric Similarity, AVD = Absolute Volume Difference,ASSD = Average Surface Distance, HD 95 = Hausdorff Distance 95th percentile, SDT = Surface Dice at Tolerance, [2] set *threshold* volume = below this volume threshold images are considered to have no lesion

Amukotuwa, S., Straka, M., Aksoy, D., Fischbein, N., Desmond, P., Albers, G., Bammer, R., 2019. Cerebral blood flow predicts the infarct core. Stroke 50, 2783–2789. URL: https://www.ahajournals.org/doi/abs/10.1161/STROKEAHA.119.026640, doi:doi:10.1161/STROKEAHA.119.026640.

Bakas, S., Akbari, H., Sotiras, A., Bilello, M., Rozycki, M., Kirby, J.S., Freymann, J.B., Farahani, K., Davatzikos, C., 2017. Advancing the cancer genome atlas glioma mri collections with expert segmentation labels and radiomic features. Sci Data 4, 170117. doi:10.1038/sdata.2017.117. 2052-4463 Bakas, Spyridon Orcid: 0000-0001-8734-6482 Akbari, Hamed Orcid: 0000-0001-9786-3707 Sotiras, Aristeidis Bilello, Michel Rozycki, Martin Kirby, Justin S Freymann, John B Farahani, Keyvan Davatzikos, Christos U24 CA189523/CA/NCI NIH HHS/United States Dataset Journal Article Research Support, N.I.H., Extramural 2017/09/06 Sci Data. 2017 Sep 5;4:170117. doi: 10.1038/sdata.2017.117.

Bakas, S., Reyes, M., Jakab, A., Bauer, S., Rempfler, M., Crimi, A., Shinohara, R.T., Berger, C., Ha, S.M., Rozycki, M., 2018. Identifying the best machine learning algorithms for brain tumor segmentation, progression assessment, and overall survival prediction in the brats challenge. arXiv preprint arXiv:1811.02629 .

Becker, A.S., Chaitanya, K., Schawkat, K., Muehlematter, U.J., Hötker, A.M., Konukoglu, E., Donati, O.F., 2019. Variability of manual segmentation of the prostate in axial t2-weighted mri: A multi-reader study. European Journal of Radiology 121, 108716. URL: https://www.sciencedirect.com/science/article/pii/S0720048X19303663https://www.sciencedirect.com/science/article/pii/S0720048X19303663?via%3Dihub, doi:https://doi.org/10.1016/j.ejrad.2019.108716.

Bertels, J., Eelbode, T., Berman, M., Vandermeulen, D., Maes, F., Bisschops, R., Blaschko, M.B., 2019. Optimizing the dice score and jaccard index for medical image segmentation: Theory and practice, in: Medical Image Computing and Computer Assisted Intervention–MICCAI 2019: 22nd International Conference, Shenzhen, China, October 13–17, 2019, Proceedings, Part II 22, Springer. pp. 92–100.

Brosch, T., Peters, J., Groth, A., Stehle, T., Weese, J., 2018. Deep learning-based boundary detection for model-based segmentation with application to mr prostate segmentation, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer. pp. 515–522.

Caradu, C., Spampinato, B., Vrancianu, A.M., Bérard, X., Ducasse, E., 2021. Fully automatic volume segmentation of infrarenal abdominal aortic aneurysm computed tomography images with deep learning approaches versus physician controlled manual segmentation. Journal of Vascular Surgery 74, 246–256.e6. URL: https://www.sciencedirect.com/science/article/pii/S0741521420325106, doi:https://doi.org/10.1016/j.jvs.2020.11.036.

Cheng, B., Girshick, R., Dollár, P., Berg, A.C., Kirillov, A., 2021. Boundary iou: Improving object-centric image segmentation evaluation, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 15334–15342.

Commowick, O., Istace, A., Kain, M., Laurent, B., Leray, F., Simon, M., Pop, S.C., Girard, P., Améli, R., Ferré, J.C., Kerbrat, A., Tourdias, T., Cervenansky, F., Glatard, T., Beaumont, J., Doyle, S., Forbes, F., Knight, J., Khademi, A., Mahbod, A., Wang, C., McKinley, R., Wagner, F., Muschelli, J., Sweeney, E., Roura, E., Lladó, X., Santos, M.M., Santos, W.P., Silva-Filho, A.G., Tomas-Fernandez, X., Urien, H., Bloch, I., Valverde, S., Cabezas, M., Vera-Olmos, F.J., Malpica, N., Guttmann, C., Vukusic, S., Edan, G., Dojat, M., Styner, M., Warfield, S.K., Cotton, F., Barillot, C., 2018. Objective evaluation of multiple sclerosis lesion segmentation using a data management and processing infrastructure. Scientific Reports 8, 13650. URL: https://doi.org/10.1038/s41598-018-31911-7https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6135867/pdf/41598_2018_Article_31911.pdf, doi:10.1038/s41598-018-31911-7.

Dewey, B.E., Zhao, C., Reinhold, J.C., Carass, A., Fitzgerald, K.C., Sotirchos, E.S., Saidha, S., Oh, J., Pham, D.L., Calabresi, P.A., van Zijl, P.C.M., Prince, J.L., 2019. Deepharmony: A deep learning approach to contrast harmonization across scanner changes. Magnetic Resonance Imaging 64, 160–170. URL: https://www.sciencedirect.com/science/article/pii/S0730725X18306490, doi:https://doi.org/10.1016/j.mri.2019.05.041.

Elguindi, S., Zelefsky, M.J., Jiang, J., Veeraraghavan, H., Deasy, J.O., Hunt, M.A., Tyagi, N., 2019. Deep learning-based auto-segmentation of targets and organs-at-risk for magnetic resonance imaging only

planning of prostate radiotherapy. Physics and Imaging in Radiation Oncology 12, 80–86. URL: https://www.sciencedirect.com/science/article/pii/S2405631619300569https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7192345/pdf/main.pdf, doi:https://doi.org/10.1016/j.phro.2019.11.006.

Filippi, M., Preziosa, P., Banwell, B.L., Barkhof, F., Ciccarelli, O., De Stefano, N., Geurts, J.J., Paul, F., Reich, D.S., Toosy, A.T., et al., 2019. Assessment of lesions on magnetic resonance imaging in multiple sclerosis: practical guidelines. Brain 142, 1858–1875.

Gautam, A., Raman, B., 2021. Towards effective classification of brain hemorrhagic and ischemic stroke using cnn. Biomedical Signal Processing and Control 63, 102178.

Heimann, T., Van Ginneken, B., Styner, M.A., Arzhaeva, Y., Aurich, V., Bauer, C., Beck, A., Becker, C., Beichel, R., Bekes, G., et al., 2009. Comparison and evaluation of methods for liver segmentation from ct datasets. IEEE transactions on medical imaging 28, 1251–1265.

Huttenlocher, D.P., Klanderman, G.A., Rucklidge, W.J., 1993. Comparing images using the hausdorff distance. IEEE Transactions on pattern analysis and machine intelligence 15, 850–863.

Isensee, F., Jaeger, P.F., Kohl, S.A.A., Petersen, J., Maier-Hein, K.H., 2021. nnu-net: a self-configuring method for deep learning-based biomedical image segmentation. Nat Methods 18, 203–211. URL: https://www.ncbi.nlm.nih.gov/pubmed/33288961https://www.nature.com/articles/s41592-020-01008-z.pdf, doi:10.1038/s41592-020-01008-z. isensee, Fabian Jaeger, Paul F Kohl, Simon A A Petersen, Jens Maier-Hein, Klaus H eng Research Support, Non-U.S. Gov't 2020/12/09 Nat Methods. 2021 Feb;18(2):203-211. doi: 10.1038/s41592-020-01008-z. Epub 2020 Dec 7.

Isensee, F., Schell, M., Pflueger, I., Brugnara, G., Bonekamp, D., Neuberger, U., Wick, A., Schlemmer, H.P., Heiland, S., Wick, W., et al., 2019. Automated brain extraction of multisequence mri using artificial neural networks. Human brain mapping 40, 4952–4964.

Janssens, R., Zeng, G., Zheng, G., 2018. Fully automatic segmentation of lumbar vertebrae from ct images using cascaded 3d fully convolutional networks, in: 2018 IEEE 15th international symposium on biomedical imaging (ISBI 2018), IEEE. pp. 893–897.

Jungo, A., Meier, R., Ermis, E., Blatti-Moreno, M., Herrmann, E., Wiest, R., Reyes, M., 2018. On the effect of inter-observer variability for a reliable estimation of uncertainty of medical image segmentation, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer. pp. 682–690.

Karimi, D., Dou, H., Warfield, S.K., Gholipour, A., 2020. Deep learning with noisy labels: Exploring techniques and remedies in medical image analysis. Medical image analysis 65, 101759.

Kuijf, H.J., Biesbroek, J.M., De Bresser, J., Heinen, R., Andermatt, S., Bento, M., Berseth, M., Belyaev, M., Cardoso, M.J., Casamitjana, A., et al., 2019. Standardized assessment of automatic segmentation of white matter hyperintensities and results of the wmh segmentation challenge. IEEE transactions on medical imaging 38, 2556–2568.

Leevy, J.L., Khoshgoftaar, T.M., Bauder, R.A., Seliya, N., 2018. A survey on addressing high-class imbalance in big data. Journal of Big Data 5, 42. URL: https://doi.org/10.1186/s40537-018-0151-6, doi:10.1186/s40537-018-0151-6.

Litjens, G., Toth, R., Van De Ven, W., Hoeks, C., Kerkstra, S., van Ginneken, B., Vincent, G., Guillard, G., Birbeck, N., Zhang, J., et al., 2014. Evaluation of prostate segmentation algorithms for mri: the promise12 challenge. Medical image analysis 18, 359–373.

Liu, C.F., Hsu, J., Xu, X., Ramachandran, S., Wang, V., Miller, M.I., Hillis, A.E., Faria, A.V., Wintermark, M., Warach, S.J., Albers, G.W., Davis, S.M., Grotta, J.C., Hacke, W., Kang, D.W., Kidwell, C., Koroshetz, W.J., Lees, K.R., Lev, M.H., Liebeskind, D.S., Sorensen, A.G., Thijs, V.N., Thomalla, G., Wardlaw, J.M., Luby, M., The, S., investigators, V.I., 2021. Deep learning-based detection and segmentation of diffusion abnormalities in acute ischemic stroke. Communications Medicine 1, 61. URL: https://doi.org/10.1038/s43856-021-00062-8https://www.nature.com/articles/s43856-021-00062-8.pdf, doi:10.1038/s43856-021-00062-8.

Maier-Hein, L., Reinke, A., Christodoulou, E., Glocker, B., Godau, P., Isensee, F., Kleesiek, J., Kozubek, M., Reyes, M., Riegler, M.A., 2022. Metrics reloaded: Pitfalls and recommendations for image analysis validation. arXiv preprint arXiv:2206.01653 .

Mehta, R., Filos, A., Baid, U., Sako, C., McKinley, R., Rebsamen, M.,

15

Dätwyler, K., Meier, R., Radojewski, P., Murugesan, G.K., et al., 2022. Qubrats: Miccai brats 2020 challenge on quantifying uncertainty in brain tumor segmentation-analysis of ranking scores and benchmarking results. Journal of Machine Learning for Biomedical Imaging 1.

Menze, B.H., Jakab, A., Bauer, S., Kalpathy-Cramer, J., Farahani, K., Kirby, J., Burren, Y., Porz, N., Slotboom, J., Wiest, R., Lanczi, L., Gerstner, E., Weber, M.A., Arbel, T., Avants, B.B., Ayache, N., Buendia, P., Collins, D.L., Cordier, N., Corso, J.J., Criminisi, A., Das, T., Delingette, H., Ç, D., Durst, C.R., Dojat, M., Doyle, S., Festa, J., Forbes, F., Geremia, E., Glocker, B., Golland, P., Guo, X., Hamamci, A., Iftekharuddin, K.M., Jena, R., John, N.M., Konukoglu, E., Lashkari, D., Mariz, J.A., Meier, R., Pereira, S., Precup, D., Price, S.J., Raviv, T.R., Reza, S.M.S., Ryan, M., Sarikaya, D., Schwartz, L., Shin, H.C., Shotton, J., Silva, C.A., Sousa, N., Subbanna, N.K., Szekely, G., Taylor, T.J., Thomas, O.M., Tustison, N.J., Unal, G., Vasseur, F., Wintermark, M., Ye, D.H., Zhao, L., Zhao, B., Zikic, D., Prastawa, M., Reyes, M., Leemput, K.V., 2015. The multimodal brain tumor image segmentation benchmark (brats). IEEE Transactions on Medical Imaging 34, 1993–2024. doi:10.1109/TMI.2014.2377694.

Nikolov, S., Blackwell, S., Zverovitch, A., Mendes, R., Livne, M., De Fauw, J., Patel, Y., Meyer, C., Askham, H., Romera-Paredes, B., 2018. Deep learning to achieve clinically applicable segmentation of head and neck anatomy for radiotherapy. arXiv preprint arXiv:1809.04430 .

Powers, W.J., Rabinstein, A.A., Ackerson, T., Adeoye, O.M., Bambakidis, N.C., Becker, K., Biller, J., Brown, M., Demaerschalk, B.M., Hoh, B., 2018. 2018 guidelines for the early management of patients with acute ischemic stroke: a guideline for healthcare professionals from the american heart association/american stroke association. stroke 49, e46–e99.

Prados, F., Ashburner, J., Blaiotta, C., Brosch, T., Carballido-Gamio, J., Cardoso, M.J., Conrad, B.N., Datta, E., Dávid, G., De Leener, B., et al., 2017. Spinal cord grey matter segmentation challenge. Neuroimage 152, 312–329.

Schell, M., Tursunova, I., Fabian, I., Bonekamp, D., Neuberger, U., Wick, W., Bendszus, M., Maier-Hein, K., Kickingereder, P., 2019. Automated brain extraction of multi-sequence mri using artificial neural networks, European Congress of Radiology-ECR 2019.

Shusharina, N., Söderberg, J., Edmunds, D., Löfman, F., Shih, H., Bortfeld, T., 2020. Automated delineation of the clinical target volume using anatomically constrained 3d expansion of the gross tumor volume. Radiotherapy and Oncology 146, 37–43. URL: https://www.sciencedirect.com/science/article/pii/S0167814020300475, doi:https://doi.org/10.1016/j.radonc.2020.01.028.

Styner, M., Lee, J., Chin, B., Chin, M., Commowick, O., Tran, H., Markovic-Plese, S., Jewells, V., Warfield, S., 2008. 3d segmentation in the clinic: A grand challenge ii: Ms lesion segmentation. MIDAS journal 2008, 1–6.

Taha, A.A., Hanbury, A., 2015. Metrics for evaluating 3d medical image segmentation: analysis, selection, and tool. BMC medical imaging 15, 29–29. URL: https://pubmed.ncbi.nlm.nih.gov/26263899https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4533825/https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4533825/pdf/12880_2015_Article_68.pdf, doi:10.1186/s12880-015-0068-x. 26263899[pmid] PMC4533825[pmcid] 10.1186/s12880-015-0068-x[PII].

Tiulpin, A., Finnilä, M., Lehenkari, P., Nieminen, H.J., Saarakkala, S., 2020. Deep-learning for tidemark segmentation in human osteochondral tissues imaged with micro-computed tomography, in: International Conference on Advanced Concepts for Intelligent Vision Systems, Springer. pp. 131–138.

Vania, M., Mureja, D., Lee, D., 2019. Automatic spine segmentation from ct images using convolutional neural network via redundant generation of class labels. Journal of Computational Design and Engineering 6, 224–232. URL: https://doi.org/10.1016/j.jcde.2018.05.002, doi:10.1016/j.jcde.2018.05.002.

de Vos, V., Timmins, K.M., van der Schaaf, I.C., Ruigrok, Y., Velthuis, B.K., Kuijf, H.J., 2021. Automatic cerebral vessel extraction in tof-mra using deep learning. arXiv preprint arXiv:2101.09253 .

Wang, B., Jin, S., Yan, Q., Xu, H., Luo, C., Wei, L., Zhao, W., Hou, X., Ma, W., Xu, Z., et al., 2021. Ai-assisted ct imaging analysis for covid-19 screening: Building and deploying a medical ai system. Applied Soft Computing 98, 106897.

Yang, F., Zamzmi, G., Angara, S., Rajaraman, S., Aquilina, A., Xue, Z., Jaeger, S., Papagiannakis, E., Antani, S.K., 2023. Assessing inter-annotator agreement for medical image segmentation. IEEE Access 11, 21300–21312. doi:10.1109/ACCESS.2023.3249759.

Zhu, R., Guo, Y., Xue, J.H., 2020. Adjusting the imbalance ratio by the dimensionality of imbalanced data. Pattern Recognition Letters 133, 217–223.

## 8. Supplemental Material

### 8.1. Dice Score of the Random Model

We aim to show how the Dice score depends on the volume of the target object, independently of the model performance. We do this by showing that the trend of the Dice score, with respect to volume, is present in the Dice score for a simple, random model. Furthermore, we show that this trend is replicated in multiple settings.

We define the random model for a parameter $IR$, as one where each voxel is chosen to be positive in the predicted mask with a probability $IR$ and there are exactly $IR \times |M^{ROI}|$ voxels in the target object. Note that the expected number of predicted positive voxels under this model is exactly $IR \times |M^{ROI}|$.

Under this model, we can compute the expected Dice score, $E_D$, across multiple draws. We use the standard combinations notation, $\binom{a}{b}$ to denote the number of orderings where we flip $b$ heads from $a$ coin flips. Note that by definition, $TP + FN = IR|M^{ROI}|$, the size of the target object.

$$E_D(p) = 2^{-|M^{ROI}|} \sum D(TP,TN,FP,FN)Pr[TP,TN,FP,FN] \tag{8}$$

where

$$D(\ldots) = \left( \frac{2 \times TP}{2 \times TP + FN + FP} \right) \tag{9}$$

and

$$Pr[\ldots] = \binom{TP+FN}{TP}\binom{TN+FP}{TN} \times p^{TP+FP}(1-p)^{FN+TP} \tag{10}$$

**Table 7: Reference Annotation Uncertainty: Inter-expert and Majority-expert Agreement**

| Categories | Metric[2] | NCCT Inter-expert[1] | | NCCT Majority-expert[1] | | Spinal Inter-expert[1] | | Spinal Majority-expert[1] | |
|---|---|---|---|---|---|---|---|---|---|
| **Uncertainty** | **U-score** | 0.87 | ± 0.05 | | | 0.39 | ± 0.02 | | |
| **Volume** | **VS** | 0.50 | ± 0.02 | 0.75 | ± 0.03 | 0.93 | ± 0.01 | 0.95 | ± 0.01 |
| | **AVD [ml]** | 10.50 | ± 2.11 | 4.25 | ± 0.98 | 0.13 | ± 0.02 | 0.10 | ± 0.02 |
| **Overlap** | **Dice** | 0.39 | ± 0.05 | 0.67 | ± 0.04 | 0.84 | ± 0.01 | 0.91 | ± 0.01 |
| | **Precision** | 0.39 | ± 0.05 | 0.64 | ± 0.06 | 0.83 | ± 0.02 | 0.95 | ± 0.01 |
| | **Recall** | 0.41 | ± 0.04 | 0.91 | ± 0.02 | 0.85 | ± 0.02 | 0.88 | ± 0.01 |
| **Distance** | **ASSD** | 4.75 | ± 0.54 | 2.03 | ± 0.23 | 0.11 | ± 0.01 | 0.07 | ± 0.01 |
| | **HD 95 [mm]** | 16.73 | ± 2.12 | 9.49 | ± 1.15 | 0.50 | ± 0.14 | 0.48 | ± 0.13 |
| | **SDT small**[3] | 0.40 | ± 0.03 | 0.67 | ± 0.02 | 0.85 | ± 0.07 | 0.93 | ± 0.04 |
| | **SDT large**[4] | 0.59 | ± 0.05 | 0.83 | ± 0.03 | 0.85 | ± 0.07 | 0.93 | ± 0.04 |

[1] per case and data set median ± 95% Confidence Interval (bootstrapped)
[2] VS = Volumetric Similarity, AVD = Absolute Volume Difference, ASSD = Average Surface Distance, HD 95 = Hausdorff Distance 95th percentile, SDT = Surface Dice at Tolerances
[3] Surface Dice at Tolerance with 2mm for NCCT and BRATS 2019 models and 0.05mm for the Spinal Cord model
[4] Surface Dice at Tolerance with 5mm for NCCT and BRATS 2019 models and 0.1mm for the Spinal Cord



Fig. 6: Data Sampling and Partition of 5-fold-Cross-Validation



Fig. 7: Line plot of image classification metrics value over $p_i$ for the NCCT and BRATS 2019 non-enhancing tumor and whole tumor data set, respectively