003

004

005

008

009

010

011

012

013

014

015

016

017

018

019

020

021

023

024

025

026

027

028

029

0.31

032

033

034

035

036

037

038

039

040

041

043

044

045

046

057

060

061

065

075

079

080

087

088

089

090

093

# Towards clinical application of liver, vessel, and tumor segmentation using partially labeled data

Anonymized authors \*

## Abstract

Accurate delineation of liver parenchyma, intrahepatic vessels, and tumors (LVT) may aid earlier tumor detection, consistent response assessment, and surgical planning for patients with liver cancer. Deep learning (DL) may enable such automated delineation, but available CT datasets are fragmented and partially labeled, making them unsuited for end-to-end training. We investigate a single-head, 3D segmentation framework that learns from such fragmented data by: (i) loss masking per class or voxel to ignore missing annotations, (ii) using multihot targets and the anatomical hierarchy inherent to liver, vessels, and tumors, to handle overlapping structures without class competition. In controlled ablations that simulate partial-label training, this multi-label masked strategy reliably outperforms masked multi-class baselines, avoids precision collapse, and improves tumor overlap and lesion detection sensitivity. Scaling training to multiple partially labeled datasets, the model surpasses full-resolution nnU-Net on an external clinical cohort, with higher tumor and vessel segmentation performance. We conduct a qualitative retrospective case study to illustrate the clinical potential of the LVT application. We find that LVT models can enable earlier detection of metastasis by six months, longitudinal size tracking aligned with radiologist measurements, 3D tumor-vessel visualization for surgical planning, and stable inter-phase liver volumetry (2% deviation). These results show that multi-label masked learning enables robust, clinically relevant LVT segmentation from partially labeled datasets.

#### 1 Introduction

Effective management of liver cancer, including patient follow-up and surgical treatment, relies on patient-specific understanding of the liver parenchyma, intrahepatic vasculature, and tumor burden. Accurate delineation of these structures may enable a range of impactful clinical tasks: earlier and more reliable tumor detection during patient follow-up, objective and consistent longitudinal response assessment, preoperative virtual planning with 3D visualization, and automatic liver volumetry to estimate functional reserve before major resec-

tions [1-7].

Deep learning (DL) based segmentation models have the potential to automatically produce high-quality segmentations of the liver parenchyma, intrahepatic vessels, and hepatic tumors (LVT) [8, 9]. However, publicly available 3D annotation of liver, vessels, and tumors, to train such models, remains scarce and fragmented. This fragmentation has constrained clinically geared liver applications to single-task models, limiting generalizability and complicating clinical deployment.

In this paper, we address the scarcity of fully labeled CT liver, vessel, and tumor segmentation data and the fragmentation of labeled datasets. To achieve generalizability with little data, we leverage a recently proposed augmentation strategy for contrast-enhanced CT liver images called Random windowing [10, 11]. Furthermore, to exploit datasets with partial labels, we explore multiple segmentation strategies capable of learning from partial labels and potentially overlapping structures end-to-end. Ultimately, we try to answer the question: How to leverage partially labeled datasets with overlapping structures in LVT segmentation? We identify that multi-label binary segmentation with a masked loss and multi-hot encoded labels, to allow class overlaps, can balance the loss contribution of partial labels and better learn from overlapping classes (Figure 1).

We demonstrate the effectiveness of our approach with quantitative evaluation against the strong nnU-Net baseline[12]. To complement the quantitative evaluation and to demonstrate the clinical potential of automatic DL based segmentation of LVT structures, we qualitatively evaluate a clinical case study that highlights the potential of such models.

The case study illustrates how automatic LVT predictions could enable earlier detection of liver tumors, track tumor size over time comparable to radiologist measurements, and deliver 3D visualizations of tumor—vessel relationships. We also show that automated liver volumes remain stable across contrast phases, supporting volumetric assessment in clinical workflows.

Our contributions are twofold:

- 1. We analyze how to efficiently leverage partially labeled data with overlapping regions to segment the liver, vessels, and tumors in CT images.
- 2. We demonstrate real-world clinical potential of 094

<sup>\*</sup>Corresponding Author.

096

097

098

099

100

101

102

103

104

105

106

107

108

109

111

112

113

114

115

116

118

119

120

121

122

123

124

125

126

127

128

129

130

131

137

147

150

151

152

154

159

160

161

162

163

167

168

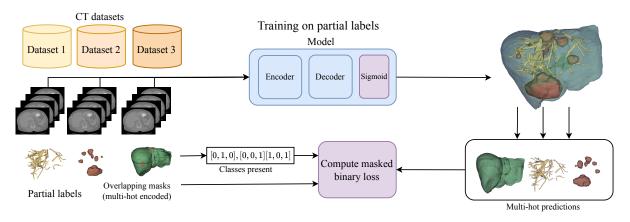


Figure 1. A schematic overview displaying our approach for learning robust and simultaneous liver, vessel, and tumor segmentation from multiple partially labeled datasets. Formulating the objective functions for binary segmentation allows us to define overlapping classes and mask the loss based on the missing labels.

LVT models through a combination of quantitative evaluation on challenging clinical data and qualitative retrospective cases that highlight earlier tumor detection, longitudinal monitoring, 3D surgical planning, and consistent volumetry.

We aim to build on these aforementioned advances to make a unified LVT segmentation model from the available public datasets and demonstrate its potential for the clinic.

#### 2 Related work

CT-based segmentation of LVT has often been addressed with task-specific models, as many of the impactful and available datasets provide only a subset of these labels. Methodological developments have been driven largely by medical segmentation challenges, and public datasets from these, such as LiTS/MSD Liver for liver and tumor, MSD HepaticVessel for hepatic vessels [8, 9]. Additionally, 3Dircadb provides the complete label set, but on a comparatively small cohort (20 cases) [13].

Architecturally, 3D encoder-decoder CNNs and the self-configuring nnU-Net remain strong baselines for medical 3D applications despite recent advances in vision transformers [12, 14–16]. Vision transformers have shown promise on large, diverse datasets, but often underperform CNN in small-to-moderate data regimes typical of clinical CT cohorts [16–18].

Reliable LVT segmentation in CT must be performed on contrast-enhanced images to effectively see the intrahepatic structures, such as vasculature and tumor. Training DL models on contrastenhanced CT images on limited datasets is challenging due to the high image variability across contrast phases and between patients for a given phase. Recently, Random windowing was proposed as a CT-specific augmentation scheme to expose the model to realistic phase variability due to contrastenhancement [10, 11]. It samples clinically plausible HU windows stochastically during training, and has been shown to improve robustness in CT segmentation of liver tumors.

#### Learning from partial labels 2.1

As public datasets rarely share a complete LVT label space, combining them during training must be done with care to avoid class conflicts. A typical challenge is handling missing classes, because treating them as background in a softmax multi-class setup can cause conflicting regions. Prior work has addressed missing annotations via masked or partial-label losses that ignore unlabeled classes during training, preventing spurious gradients from missing annotations [19–21]. To this end, binary formulations with sigmoid outputs are particularly suitable, as they avoid the competition inherent in softmax and allow perclass supervision wherever labels are present [20, 22]. To mitigate noisy gradients from a multi-class segmentation setting using softmax, the loss masking must happen on a per-voxel basis, ignoring signals from all non-foreground regions [19].

Alternative strategies include weakly or semi- 155 supervised learning using pseudo-labels [20, 23], and multi-task/multi-head designs [24, 25] where each dataset supervises a subset of heads while sharing an encoder. While effective in some settings, pseudo-labels can introduce confirmation bias from erroneous predictions, and multi-task heads can be difficult to calibrate across datasets.

Compared to approaches that stitch together separate binary models or rely on pseudo-labels to fill missing classes, we investigate segmentation models with a single segmentation head and an end-to-end training pipeline. In this setting, we investigate how to benefit from and train on partially labeled datasets using various loss masking strategies.

## 2.2 Overlapping classes

Within the context of partial label supervision, overlapping classes are often treated as separate classes, masking the loss contribution of any missing class during training [21, 26]. Although possible in certain settings, it does not address the potential label conflicts across datasets, and could lead to suboptimal performance (Section 4.2).

Another established approach is cascaded models, which segment organ regions of interest (ROI) like the liver, before specialized models are trained for vessels and tumors within the organ ROI [9]. However, this approach leads to extra compute overhead during training and inference compared to an end-to-end pipeline in a complete label space.

Semantic segmentation on overlapping and partial labels has been addressed by probabilistic approaches that aggregate predictions [27] and multilabel approaches [28], but the approach within medical and LVT applications is largely underexplored.

## 2.3 Clinical potential for liver, vessel, and tumor segmentation

The clinical implications of LVT models are substantial. Below, we identify four routine scenarios where the LVT application has clinical potential: surgical planning, automatic liver volumetry, longitudinal patient follow-up, and tumor detection. In our final case study, we evaluate the current utility of LVT segmentation for each use case.

Surgical planning using accurate 3D spatial delineation of tumors to surrounding hepatic vasculature allows for more precise surgical planning. In an ongoing tangential study of the clinical impact of 3D liver models at Hospital1 (Anonymized), initial results suggest significant improvements in surgical planning when a 3D LVT model was used along with traditional radiological images (CT/MRI).

The integration of automatically generated 3D models into surgical procedures may significantly impact the management of complex hepatic resections [5]. 3D visualization of tumor-vessel relationships can aid surgeons when navigating in challenging anatomical landscapes [4, 7], and may reduce unintended vessel injury and improve resection margins.

Automatic liver volumetry further aids surgical decision-making by providing essential data needed for assessing hepatic functional reserve, which is particularly beneficial for major liver resections [2, 3].

**During patient follow-up**, the LVT model has potential to benefit the clinical follow-up of oncology patients undergoing chemotherapy. By automatically measuring tumor sizes and the liver volume, the model can provide a consistent and objective assessment of tumor response over time [6]. This may facilitate timely therapeutic decisions, allowing

clinicians to optimize the treatment based on tumor volume changes.

Tumor detection using a segmentation model enables precise identification of potential tumor regions in the CT scan. For patients at risk of liver metastasis, DL based tumor segmentation tools may help radiologists detect tumors early, improving diagnosis and treatment for patients.

In general, a performant liver, vessel, and tumor segmentation model has the potential to be a clinical tool and may improve diagnostic accuracy, enhance therapeutic planning, optimize patient follow-up, and increase surgical safety in hepatic tumor management.

## 3 Methodology

In this section, we describe how we adjust the default segmentation setup and training regime to accommodate partial labels, and how to learn from overlapping structures in the liver. Finally, we illustrate the clinical impact of our application by performing a qualitative retrospective analysis of the follow-up of typical liver metastasis patients.

### 3.1 Training on partially labeled data

In the partial label setting, computing the loss over regions with missing labels or unlabeled background regions requires special care. To simplify training on such datasets, we avoid evaluating the loss over unlabeled or potentially ambiguous regions in a binary and multi-class segmentation setting.

We achieve this by formulating the objective functions with a weight mask w or W to ignore the contribution of a given class or voxel, respectively, depending on the loss formulation.

Binary losses can be computed per-class on the sigmoid probabilities of the segmentation network. For each class k of a given sample with N voxels, we compute the mean voxel loss  $\ell_k$  and mask it with the per-class weight  $w_k$ , essentially removing the loss contribution for the classes with missing labels. The masked binary loss  $\mathcal{L}_B$  for partially labeled samples is thus computed for each voxel i and class as

$$\mathcal{L}_B = \frac{1}{\sum_{k=1}^K w_k} \sum_{k=1}^K w_k \frac{1}{N} \sum_{i=1}^N \ell_{ik}.$$
 (1) 267

In multi-class segmentation, the loss  $\ell_i$  is computed for each voxel over all classes, using one-hot encoded labels and softmax output probabilities. Masking out the loss contribution for ambiguous regions must therefore be performed on a per-voxel basis with  $W_i \in \{0,1\}$ , yielding the categorical loss for partially labeled sample  $\mathcal{L}_C$ 

$$\mathcal{L}_C = \frac{1}{\sum_{i=1}^{N} W_i} \sum_{i=1}^{N} \ell_i W_i. \tag{2}$$

278

279

280

281

282

283

284

285

286

287

289

290

291

292

293

295

296

297

298

299

300

302

303

304

305

306

307

309

310

311

312

313

314

317

318

319

321

322

323

324

325

326

327

335

344

345

346

348

349

351

352

353

354

357

358

360

361

365

366

368

369

370

375

376

380

This formulation removes the loss contribution of all non-foreground voxels, essentially giving us the mean foreground loss. This eliminates the problem of ambiguous regions, also when the background is modeled explicitly.

#### Multi-class and multi-label segmentation

Multi-class segmentation (MC) predicts exclusive classes using softmax activation function and the categorical loss  $\mathcal{L}_C$  with weight mask W from Equation 2. MC predicts K+1 classes and models the background explicitly in the initial channel. The explicit background channel will potentially interfere with all missing ROI. Therefore, W must mask out the loss contribution of all non-foreground regions for classes not present in the PS training set.

Multi-label segmentation (ML) uses the sigmoid activation independently for each K output channels. Training a ML model with one-hot labels for all datasets with full or partial supervision is comparable to the MC setup wrt. the training signal from exclusive classes. Loss masking is done per channel using w in Equation 1. We refer to this setup as Multi-label exclusive segmentation (MLx).

## 3.2 Segmenting overlapping structures

Given complete annotations, class exclusive segmentation setup has the benefit of yielding unambiguous regions and explicit information about boundaries between classes. However, for partially labeled datasets, class exclusivity is not guaranteed, and overlapping classes risk regions of conflicting supervision due to overlap.

To avoid this issue across partially labeled datasets, ML can be trained with overlapping classes if labels are represented as multi-hot vectors, with per-voxel labels  $y_{ik} \in \{0,1\}$  for each class k. This avoids competition between classes during training and allows supervision of whichever labels are available for that sample. In the context of partial labels, w can seamlessly be applied in Equation 1 to mask the loss contribution of unlabeled classes.

#### Anatomical liver hierarchy

In the context of LVT segmentation, vessels and tumors are anatomically contained within the liver. In our following experimental settings, we enforce this anatomical hierarchy by mapping vessel and tumor label positives into the liver channel:  $y_{i,L} \leftarrow y_{i,L} \lor y_{i,V} \lor y_{i,T}$ . This mapping is applied on-the-fly for datasets that provide one-hot labels, yielding a consistent multi-hot label space across datasets. When a class is not annotated for a sample, its loss weight  $w_k = 0$ , so it does not contribute to the objective.

## 4 Experiments

In this section, we investigate multiple strategies for training segmentation models on partially labeled datasets and how potentially overlapping classes impact the performance. We present the experiments and their results sequentially, and use the novel insight to inform our clinical case study.

Experimental setup. All experiments are performed under an identical medical image segmentation setup, where the objective is a segmentation map with mutually exclusive classes. We focus on end-to-end training pipelines using a U-Net-like architecture [29], with several modifications for robust training and results. The detailed experimental setup and evaluation settings can be found in the Appendix A.2.

### 4.1 Learning from partial labels

We construct a controlled experiment simulating a training setting with multiple datasets with partial/missing labels. Specifically, we create partially labeled training sets from different partitions of one fully annotated source dataset. This lets us test and evaluate various approaches without considering noise from distribution shifts from other data and label sources.

Simulating partial labeled training. Based on the 303 images from the HepaticVessel dataset [8], with vessel and tumor segmentation labels, and the auxiliary liver segmentation labels from Tian et al. [30], we randomly sample 5 datasets of similar size. Specifically, 20 % fully annotated (LVT masks) are reserved as hold-out test set (HV test), 20 % are used as fully supervision (FS) training with complete LVT annotations, 20 % have partial supervision (PS) with tumor mask only, 20 % with PS vessel mask, and 20 % with PS from liver mask only. Note that the liver mask comprises the complete liver organ, without "cutouts" for the vessel and tumor classes. In this regard, the liver overlaps with the vessel and tumor masks, similar to a real setting with partially labeled datasets. Further dataset details can be found in Table A.1.

The question we want to answer is "How can we leverage auxiliary datasets with partial labels to improve segmentation performance over only using the fully labeled training set?". To this end, we compare the MC and MLx end-to-end segmentation setups with their different loss masking strategies with 25 % FS with and without auxiliary PS datasets. For comparison, we also provide the 100 % FS training alternative.

We report the mean segmentation performance on the HV test and the external Ircad [13] test set

384

385

386

387

388

389

391

392

393

394

395

397

398

399

400

401

402

403

405

406

407

408

409

410

413

414

415

416

417

420

421

422

423

425

426

428

430

431

432

433

434

435

440

441

442

443

451

452

454

455

456

466

467

475

476

481

482

483

489

490

after 5-fold cross-validation training on the combined full and partially labeled data splits. We measure the Dice similarity coefficient (DSC) on the liver, vessels, and tumors of the respective test sets. Our full evaluation strategy can be found in Appendix A.2.3.

Binary segmentation benefits from partial supervision. Based on the results, presented in Table 1, we make the following observations:

- (1) In the fully supervised settings, the segmentation DSC are comparable for the liver class of MLx and MC, and higher or on par for MLx on the vessel and tumor classes across both datasets.
- (2) With partial supervision and multi-class segmentation, the DSC performance collapse compared to full supervision for almost all classes. However, the exception is segmentation performance on Ircad vessels, which exceeds all other settings. Upon closer inspection, the liver and vessel recalls of MC∩PS are actually the highest across all datasets, while the precision is lowest. This can explain the extremes in DSC, because it suggests that the model over-segments with many false positives. In the HepaticVessel dataset, it is a clear disadvantage as the vessel labels are minimal and to some degree lacking, but an advantage in the Ircad dataset, which has more dilated and detailed vessel structures. We suspect the cause of over-segmentation is the individual and unbalanced supervision each class receives in the masked loss of the categorical loss formulation. While an increase and drop in recall and precision, respectively, are observed also in the MLx∩PS setup, the DSC does not suffer as severely. We attribute this to the binary loss formulations, which natively balance foreground/background better, also in the partially labeled settings.
- (3) Contrary to the MC setup, MLx benefits from the auxiliary partially labeled data in all settings. The results suggest that the masked binary loss formulation in the multi-label setup can learn from the available data, without interfering destructively with the unlabeled classes.

#### 4.2 Training on overlapping classes

Although our desired output space is exclusive, with each voxel in the liver belonging to either the liver, vessel, or tumor class, it might be suboptimal and unnecessary during training. As binary outputs in the segmentation head allow multi-label training with overlapping classes, we investigate how the potentially conflicting regions across partially labeled datasets contribute to downstream performance.

In the same controlled environment as our partial label experiment, we ablate the effect of training on ambiguous regions. Specifically, we compare MLx and ML trained with one-hot and multi-hot labels

(Section 3.2), respectively, in the FS and PS settings. 437

Liver tumor segmentation is sensitive to label conflicts. We report the segmentation DSC and the tumor detection sensitivity computed on the connected components of the predictions. Based on the results presented in Table 2, we make the following observations:

(1) Vessel segmentation is largely unaffected by the conflicting labels in the exclusive training setup. We suspect it to be a consequence of segmenting the small vessel structure in the comparatively large surrounding liver. As the vessel structures are small, the MLx model can learn to produce multi-label, rather than exclusive, class outputs without being punished significantly in the loss, as only the FS training set has complete labels with vessel "cutouts" that punish such behaviour. (2) For the tumor DSC in the non-overlapping baseline, the performance is significantly worse compared to the overlapping version. Contrary to the vessel class, the tumor class is more massive, which leads to a larger loss impact when the model predicts the liver without cutouts for the FS set. (3) The tumor detection sensitivity drops as a consequence of partial supervision on ambiguous liver and tumor labels for both test sets. The impact of this result is pivotal, as it is not a matter of slightly worse or better segmentation overlap, but more liver tumors that are being detected. Lesion precision remains similar for both methods on Ircad, and elevated for ML on the HepaticVessel test set.

## 4.3 Learning from public datasets with partial labels

For the clinical case study, we aim to build on public CT liver datasets with partial labels to scale up the training data. Following the lessons from our previous experiments, we train the LVT model under the multi-label, class-masked regime described in Section 3.2. For a solid quantitative baseline, we evaluate against the strong, but specialized nnU-Net.

We scale up partial labeled training using the complete HepaticVessel dataset (vessel + tumor) with additional liver labels from [30], the LiTS dataset [9] (liver + tumor), and Ircad (liver, vessel, tumor). The datasets are further described in Appendix A.1. We use Random windowing for CT augmentation and to mitigate cross-dataset shift. For all other training configurations, we follow the nnU-Net setup.

For quantitative evaluation, we compare against full-resolution nnU-Net baselines trained on MSD Liver (liver+tumor) and MSD HepaticVessel (vessel). We report DSC and normalized surface dice (NSD) on an external test set of contrast-enhanced CT images from ExDS (anonymized). Following [8], we use 7 mm tolerance for liver and 3 mm for vessels

**Table 1.** Segmentation DSC reported on the liver, vessel, and tumor classes of the HepaticVessel test set and the external Ircad test set. We report the segmentation performance along with the proportion of full supervision (FS) in the training set, whether partial labeled datasets (PS) were used as auxiliary training signal, and the segmentation head used, multi-class (MC) vs. multi-label exclusive (MLx).

				HepaticVessel			Ircad	
FS	PS	Head	Liver	Vessel	Tumor	Liver	Vessel	Tumor
25 %	×	MC MLx	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	$\begin{array}{c} 0.579 \pm 0.010 \\ 0.600 \pm 0.004 \end{array}$	$\begin{array}{c} 0.517 \pm 0.024 \\ 0.536 \pm 0.009 \end{array}$	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	$\begin{array}{c} 0.372 \pm 0.033 \\ 0.439 \pm 0.012 \end{array}$	$0.484 \pm 0.013 \\ 0.501 \pm 0.003$
25 %	<b>√</b> ✓	MC MLx	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	$\begin{array}{c} 0.428 \pm 0.006 \\ 0.629 \pm 0.003 \end{array}$	$\begin{array}{c} 0.221 \pm 0.026 \\ 0.561 \pm 0.018 \end{array}$	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	$0.489 \pm 0.026$ $0.466 \pm 0.011$	$\begin{array}{c} 0.341 \pm 0.045 \\ 0.525 \pm 0.037 \end{array}$
100 %	×	MC MLx	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	$0.638 \pm 0.008$ $0.698 \pm 0.004$	$\begin{array}{c} 0.615 \pm 0.023 \\ 0.790 \pm 0.010 \end{array}$	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	$0.397 \pm 0.008$ $0.468 \pm 0.010$	$\begin{array}{c} 0.462 \pm 0.017 \\ 0.573 \pm 0.025 \end{array}$

**Table 2.** We ablate the effect of ambiguous regions, due to partial labels, during training. By allowing overlapping classes through multi-hot encoded labels, the binary segmentation head avoids mixed signals from the partially labeled liver dataset (lacking vessel and tumor). Allowing overlapping classes in the partial supervision setting leads to improved tumor segmentation.

		LVT		HepaticVessel			Ircad	
FS	PS	overlap	Vessel	Tumor	Sensitivity	Vessel	Tumor	Sensitivity
25 %	×	×	$ \begin{vmatrix} 0.600 \pm 0.004 \\ 0.601 \pm 0.009 \end{vmatrix} $	$\begin{array}{c} 0.536 \pm 0.009 \\ 0.535 \pm 0.016 \end{array}$	$\begin{array}{c} 0.734 \pm 0.034 \\ 0.740 \pm 0.033 \end{array}$	$ \begin{vmatrix} 0.439 \pm 0.012 \\ 0.437 \pm 0.015 \end{vmatrix}$	$\begin{array}{c} 0.501 \pm 0.003 \\ 0.472 \pm 0.031 \end{array}$	$\begin{array}{c} 0.639 \pm 0.031 \\ 0.663 \pm 0.037 \end{array}$
25 %	<b>√</b>	×	$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$		$0.779 \pm 0.027$ $0.818 \pm 0.037$			

**Table 3.** Evaluation of our multi-label segmentation network and the full-res nnU-Net trained on MSD Liver (liver + tumor) and MSD HepaticVessel (vessel). The models are evaluated on contrast-enhanced CT images from the ExDS external dataset.

Task	Metric	nnU-Net	LVT (ours)
Tumor	DSC NSD	$0.723 \pm 0.145$ $0.706 \pm 0.214$	$0.778 \pm 0.106 \\ 0.771 \pm 0.136$
Liver	DSC NSD	$0.912 \pm 0.07$ $0.951 \pm 0.058$	$0.898 \pm 0.066$ $0.959 \pm 0.051$
Vessels	DSC NSD	$\begin{array}{c} 0.545 \pm 0.051 \\ 0.788 \pm 0.061 \end{array}$	$0.575 \pm 0.048 \\ 0.808 \pm 0.053$

and tumors (Table 3).

Our model outperforms nnU-Net on tumors and vessels across DSC and NSD, and achieves higher liver NSD with slightly lower liver DSC, consistent with minor over-segmentation addressed by the surface-tolerant NSD metric. These results indicate that learning from additional partially labeled datasets with a multi-label, masked loss and Random windowing improves robustness and clinical relevance. We next present qualitative retrospective analyses in Section 5.

## 5 Case study

Up until this point, we have validated our methods from a quantitative perspective. In this section, we shift our focus to the clinical practice and highlight the clinical usefulness of the LVT application.

#### 5.1 Tumor detection

Retrospectively analyzing the longitudinal CT scans of a patient and comparing the predictions with the radiology reports from the follow-up allows us to identify if the model could have assisted in the early detection of tumors. Such retrospective analysis helps identify when what might seem like a false positive tumor prediction by the model actually was a missed tumor by the radiologist.

For a given patient surgically treated for colorectal cancer at Hospital1 (anonymized), with a high risk of developing liver metastasis, we obtained predictions for the contrast-enhanced liver CT scans from the follow-up studies in the patient pathway from both our LVT model and the nnU-Net baseline. After the patient's initial treatment, they had no metastasized liver cancer for the following 1,5 years, but in March 2009, a 4 cm tumor was discovered in the left liver lobe. In the preceding CT scan, 6 months before, the radiologist stated that there were no suspicious lesions in the liver. However, retrospective analysis with our LVT model marked a small tumor region in the left liver lobe in the same scan, 6 months prior to the radiologist identifying a liver tumor in that

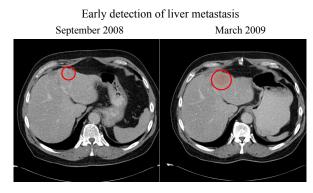


Figure 2. Comparison of CT images of the same patient 6 months apart. In the precedent scan from September 2008, the radiologist identified no suspicious lesions in the liver. 6 months later, the radiologist found a tumor measuring 4 cm in diameter. Our LVT model marked a corresponding lesion displayed in the image in the former image, 6 months before the radiologist.

same region (Figure 2). The prediction from the nnU-Net did not detect this tumor.

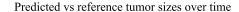
This example illustrates the potential of DL assisted image analysis and how it may lead to significantly earlier detection of liver metastasis.

### 5.2 Follow-up and tumor monitoring

A key consideration when treating a patient with, e.g., liver metastasis, is the size of the metastatic region over time. Tracking the lesion's size helps assess how the patient responds to the treatment they receive. Decreasing tumor size suggests that the patient responds well to the treatment, while growth indicates tumor resistance to the treatment.

We retrospectively analyze a patient's CT liver scans during the follow-up period and assess how the LVT model performs automatic size measurements of the tumor. We report the largest dimension of the tumor in the x-y plane and compare it against the radiologist's measurements at the time of the study.

We present the results in Figure 3 and find the extracted measurements to correlate well with the radiologist's measurements. During the follow-up period, the patient experienced an initial period of tumor growth after a lesion of liver metastasis was located in liver segment 4. The radiologist's and the model's predictions align well during this period. After the initial tumor growth, the patient was considered inoperable and began chemotherapy with a good response, leading to tumor regression. However, in two scans during the regression, the model yielded false positives in another liver segment, compared with the radiologists' findings. Six months post-treatment, disease progression was again observed, and despite further management, the malig-



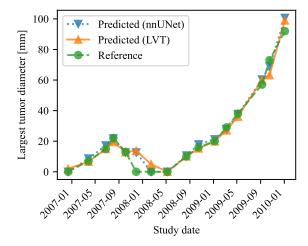


Figure 3. Comparison of tumor measurements from the LVT model's predictions and the radiologist's reference measurements. The patient develops liver metastasis in segment 4 of the liver and responds well to the treatment initially. After a period of tumor regression, the malignancy continued to advance.

nancy continued to progress. During this critical period, the LVT model matches the radiologist in tumor detection and size prediction.

### 5.3 Surgery planning

Surgical resection of tumors is, in most cases, considered the only cure for liver metastasis [1]. As the liver is a complex organ with eight independent anatomical segments with its own blood supplies from the hepatic arteries and portal veins, the tumor's precise location and relation to the vessels are crucial. Manual segmentation of these structures is too costly and rare in clinical practice. Automatic segmentation tools provide totally novel insight for the multidisciplinary team and surgeons treating the patient.

Our LVT segmentation model is able to precisely delineate the tumor and blood vessels in high-quality contrast-enhanced CT images of the liver. For a patient at Hospital1 (Anonymized), we retrospectively obtain LVT predictions from their CT images to illustrate the output when visualized in 3D software. The results are shown in Figure 4 and display the liver and delineation of a liver tumor in segment 7 with its surrounding vessels. The 3D view of the LVT predictions makes the evaluation of proximity to the structures surrounding the tumor. The visualizations are produced with 3D Slicer image computing platform [31], where the user can edit the predictions if needed, and subsequently make precise measurements before final assessment.

600

601

602

603

604

605

607

608

609

610

611

612

614

615

616

617

618

619

620

621

622

623

635

636

637

638

647

649

650

651

653

654

655

663

665

666

667

669

673

678

**Figure 4.** Illustration of the potential of automatic 3D segmentation of liver, vessels, and tumors in CT images. For a patient who is considered for surgical resection of liver metastasis, the 3D visualization of the liver parenchyma, hepatic vessels, and liver tumor can be a valuable support in surgical planning.

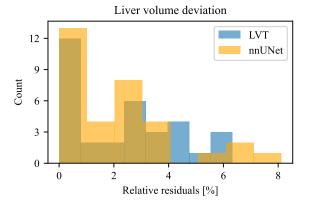


Figure 5. Relative deviation of estimated liver volume between images of different contrast phases.

#### 5.4 Automatic liver volumetry

Prior to major liver resections, CT volumetry can be used to measure the liver volume and estimate the hepatic functional reserve for the patient.

We aim to illustrate the efficacy of the LVT model for this purpose in a clinical setting by retrospectively analyzing unlabeled images from the clinic. We assume that a patient's liver on images from the same day is similar in size, and we aim to compare the model's predictions across contrast phases. To this end, we retrieved 33 images from 14 patient studies, where there are 2 or 3 contrast-enhanced images in each study, with images in the arterial, venous, or late phase. We use the LVT model to obtain liver masks for each patient and compute the liver volumes for each image. As the liver sizes vary from 1.3 L to 2.5 L across patients, we report the relative residuals in percentages of the mean liver volume of each study and plot the results in Figure 5.

We find the liver volume deviation from the reference to be within  $\approx 2\%$  for all cases, and show that the liver measurements are consistent across images of the same patient. This is within the margin of what is expected, as the measured volume usually varies slightly between images of different contrast phases [3].

## 6 Limitations and future work 625

Despite recent efforts, consistent and reliable segmentation of liver tumors and vessels remains a difficult task. Table 3 reports improvements over the baseline, but DSC scores remain below 0.8 for tumors and vessels, which suggests that false positives and missed structures are to be expected. False positives may arise when cysts, necrotic tissue, or low attenuation liver regions are misclassified as tumors, or when blurred vessel boundaries lead to over-segmentation. Conversely, missed structures may occur for small and early-stage tumors, and poorly contrasted or thin vessels.

In the partial label experiments in subsection 4.1 and subsection 4.2, we observed certain inconsistencies in vessel segmentation performance across the HV Test set and Ircad dataset. These inconsistencies are symptoms of different label characteristics of the vessels in the two datasets, which have different quality and level of detail. We therefore recommend careful evaluation when comparing these datasets.

While our retrospective analyses demonstrate the potential for clinical utility of the LVT model, integrating it into real-time clinical workflows remains an open challenge. To further identify the strengths and limitations of DL-based liver, vessel, and tumor applications, we recommend thorough clinical validation with expert supervision to validate the model's impact on patient outcomes and clinicians' workflows.

## 7 Conclusion

This study explores multi-label and multi-class approaches in the context of CT liver, vessel, and tumor segmentation, to effectively handle overlapping and potentially ambiguous regions from partially labeled datasets. We find a binary multi-label segmentation setup with class-wise loss masking to work well for this setting. Allowing overlapping regions in the label space enables the use of public datasets with partial labels during training to learn simultaneous liver, vessel, and tumors labels in CT images. Our results show that our approach is particularly beneficial for tumors and vessels, allowing us to benefit from datasets with partial and ambiguous labels.

We evaluate the LVT model on clinical data to illustrate the potential for real-world utility in the clinic. In retrospective analysis of real patients, we demonstrated that the model has the potential to detect tumors earlier than the radiologist, accurately track tumor progression, provide 3D visualization of complex liver structures, and reliably perform liver volumetry for real patients. These results underscore the potential for AI-driven tools for diagnostic accuracy, optimizing treatment planning, and improving patient outcomes.

707

708

709

710

711

713

714

715

717

718

719

720

721

722

723

724

725

728

729

730

731

732

735

737

738

740

741

742

749

750

752

754

758

767

777

784

## References

- L.-X. Liu, W.-H. Zhang, and H.-C. Jiang. 681 "Current treatment for liver metastases from 682 colorectal cancer". In: World Journal of Gas-683  $troenterology: WJG~9.2~(2003), \, \mathrm{pp.}~193–200.$ 684 DOI: 10.3748/wjg.v9.i2.193. 685
- A. Guglielmi, A. Ruzzenente, S. Conci, A. 686 Valdegamberi, and C. Iacono. "How Much 687 Remnant Is Enough in Liver Resection?" In: 688 Digestive Surgery 29.1 (2012), pp. 6–17. DOI: 689 10.1159/000335713. 690
- M. C. Lim, C. H. Tan, J. Cai, J. Zheng, and 691 A. W. C. Kow. "CT volumetry of the liver: 692 Where does it stand in clinical practice?" In: 693 Clinical Radiology 69.9 (2014), pp. 887–895. 694 DOI: 10.1016/j.crad.2013.12.021.
- A. P. Brady. "Error and discrepancy in radi-696 ology: inevitable or avoidable?" In: Insights 697 into Imaging 8.1 (2017), pp. 171-182. DOI: 698 10.1007/s13244-016-0534-1.
- A. Miyata, A. Junichi, Y. Kawaguchi, K. 700 Hasegawa, and N. Kokudo. "Simulation and 701 navigation liver surgery: an update after 2,000 virtual hepatectomies". In: Global Health & 703 Medicine 2.5 (2020), pp. 298-305. DOI: https: 704 //doi.org/10.35772/ghm.2020.01045. 705
  - S. K. Jeon, D. H. Lee, J. Park, K.-S. Suh, N.-J. Yi, S. K. Hong, and J. K. Han. "Tumor volume measured using MR volumetry as a predictor of prognosis after surgical resection of single hepatocellular carcinoma". In: European Journal of Radiology 144 (2021), p. 109962. DOI: 10.1016/j.ejrad.2021.109962.
  - P. Gavriilidis, B. Edwin, E. Pelanis, E. Hidalgo, N. de'Angelis, R. Memeo, L. Aldrighetti, and R. P. Sutcliffe. "Navigated liver surgery: State of the art and future perspectives". In: Hepatobiliary & Pancreatic Diseases International 21.3 (2022), pp. 226-233. DOI: 10.1016/j. hbpd.2021.09.002.
  - M. Antonelli, A. Reinke, S. Bakas, K. Farahani, A. Kopp-Schneider, B. A. Landman, G. Litjens, B. Menze, O. Ronneberger, R. M. Summers, B. Van Ginneken, M. Bilello, P. Bilic, P. F. Christ, R. K. G. Do, M. J. Gollub, S. H. Heckers, H. Huisman, W. R. Jarnagin, M. K. McHugo, S. Napel, J. S. G. Pernicka, K. Rhode, C. Tobon-Gomez, E. Vorontsov, J. A. Meakin, S. Ourselin, M. Wiesenfarth, P. Arbeláez, B. Bae, S. Chen, L. Daza, J. Feng, B. He, F. Isensee, Y. Ji, F. Jia, I. Kim, K. Maier-Hein, D. Merhof, A. Pai, B. Park, M. Perslev, R. Rezaiifar, O. Rippel, I. Sarasua, W. Shen, J. Son, C. Wachinger, L. Wang, Y. Wang, Y. Xia, D. Xu, Z. Xu, Y. Zheng, A. L. Simpson, L. Maier-Hein, and M. J. Cardoso. "The

- Medical Segmentation Decathlon". In: Nature Communications 13.1 (2022), p. 4128. DOI: 10.1038/s41467-022-30695-9.
- P. Bilic et al. "The Liver Tumor Segmenta-739 tion Benchmark (LiTS)". In: Medical Image Analysis 84 (2023), p. 102680. DOI: 10.1016/ j.media.2022.102680.
- [10]E. A. Østmo, K. K. Wickstrøm, K. Radiya, 743 M. C. Kampffmeyer, and R. Jenssen. "View it Like a Radiologist: Shifted Windows for Deep Learning Augmentation Of CT Im- 746 ages". In: 2023 IEEE 33rd International Workshop on Machine Learning for Signal Processing (MLSP). 2023, pp. 1–6. DOI: 10.1109/ MLSP55844.2023.10285978.
- [11]Anonymized. RandomWindowing. 751 Anonymized reference. 2025.
- F. Isensee, P. F. Jaeger, S. A. A. Kohl, J. 753 |12|Petersen, and K. H. Maier-Hein. "nnU-Net: a self-configuring method for deep learningbased biomedical image segmentation". In: Nature Methods 18.2 (2021), pp. 203–211. DOI: 757 10.1038/s41592-020-01008-z.
- L. Soler, A. Hostettler, A. Charnoz, J. Fasquel, 759 J. Moreau, A. Osswald, M. Bouhadjar, and J. 760 Marescaux. 3D image reconstruction for com- 761 parison of algorithm database: A patient specific anatomical and medical image database. Tech. rep. Strasbourg, France: IRCAD, 2010. 764 URL: https://www.ircad.fr/research/ data - sets / liver - segmentation - 3d - 766 ircadb-01/ (visited on 02/01/2024).
- O. Çiçek, A. Abdulkadir, S. S. Lienkamp, [14]T. Brox, and O. Ronneberger. "3D U-Net: 769 Learning Dense Volumetric Segmentation from Sparse Annotation". In: Medical Image Com- 771 puting and Computer-Assisted Intervention MICCAI 2016. Ed. by S. Ourselin, L. 773 Joskowicz, M. R. Sabuncu, G. Unal, and W. Wells. Lecture Notes in Computer Science. Cham: Springer International Publishing, 2016, pp. 424-432. DOI: 10.1007/978-3-319-46723-8\_49.
- F. Milletari, N. Navab, and S.-A. Ahmadi. 779 "V-Net: Fully Convolutional Neural Networks for Volumetric Medical Image Segmentation". In: 2016 Fourth International Conference on 3D Vision (3DV). 2016, pp. 565–571. DOI: 10. 783 1109/3DV.2016.79.
- F. Isensee, T. Wald, C. Ulrich, M. Baumgart- 785 ner, S. Roy, K. Maier-Hein, and P. F. Jäger. 786 "nnU-Net Revisited: A Call for Rigorous Val- 787 idation in 3D Medical Image Segmentation". 788 In: Medical Image Computing and Computer 789 Assisted Intervention – MICCAI 2024. Ed. by 790

831

832

833

834

835

836

839

840

841

849

850

851

852

853

855

861

863

867

876

885

886

887

895

896

897

902

- M. G. Linguraru, Q. Dou, A. Feragen, S. Gi-791 annarou, B. Glocker, K. Lekadir, and J. A. 792 Schnabel. Cham: Springer Nature Switzerland, 793 2024, pp. 488-498. DOI: 10.1007/978-3-031-794 72114-4\_47. 795
- A. Hatamizadeh, Y. Tang, V. Nath, D. Yang, |17|796 A. Myronenko, B. Landman, H. R. Roth, and 797 D. Xu. "UNETR: Transformers for 3D Medi-798 cal Image Segmentation". In: WACV. 2022, pp. 574-584. URL: https://openaccess thecvf . com / content / WACV2022 / html / 801 Hatamizadeh \_ UNETR \_ Transformers \_ for \_ 802 3D\_Medical\_Image\_Segmentation\_WACV\_ 803 2022\_paper.html (visited on 04/27/2023). 804
- Y. Tang, D. Yang, W. Li, H. R. Roth, B. Land-|18|805 man, D. Xu, V. Nath, and A. Hatamizadeh. 806 "Self-Supervised Pre-Training of Swin Trans-807 formers for 3D Medical Image Analysis". 808 In: 2022, pp. 20730-20740. URL: https: 809 //openaccess . thecvf . com / content / CVPR2022 / html / Tang \_ Self - Supervised \_ 811 Pre - Training of Swin Transformers \_ 812 for\_3D\_Medical\_Image\_Analysis\_CVPR\_ 813 2022\_paper.html (visited on 04/14/2023). 814
- G. González, G. R. Washko, and R. San 815 |19|José Estépar. "Multi-structure Segmentation 816 from Partially Labeled Datasets. Application 817 to Body Composition Measurements on CT 818 Scans". In: Image Analysis for Moving Organ, Breast, and Thoracic Images. Ed. by D. Stoyanov, Z. Taylor, B. Kainz, G. Maicas, R. R. 821 Beichel, A. Martel, L. Maier-Hein, K. Bha-822 tia, T. Vercauteren, O. Oktay, G. Carneiro, 823 A. P. Bradley, J. Nascimento, H. Min, M. S. 824 Brown, C. Jacobs, B. Lassen-Schmidt, K. Mori, 825 J. Petersen, R. San José Estépar, A. Schmidt-826 Richberg, and C. Veiga. Cham: Springer International Publishing, 2018, pp. 215–224. DOI: 10.1007/978-3-030-00946-5\_22. 829
  - O. Petit, N. Thome, A. Charnoz, A. Hostettler, |20|and L. Soler. "Handling Missing Annotations for Semantic Segmentation with Deep ConvNets". In: Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support: 4th International Workshop, DLMIA 2018, and 8th International Workshop, ML-CDS 2018, Held in Conjunction with MICCAI 2018, Granada, Spain, September 20, 2018, Proceedings. Berlin, Heidelberg: Springer-Verlag, 2018, pp. 20–28. DOI: 10. 1007/978-3-030-00889-5\_3.
- [21]J. Zhang, Y. Xie, Y. Xia, and C. Shen. 842 "DoDNet: Learning To Segment Multi-Organ 843 and Tumors From Multiple Partially La-844 beled Datasets". In: 2021, pp. 1195–1204. 845 URL: https://openaccess.thecvf.com/ 846 content/CVPR2021/html/Zhang\_DoDNet\_ 847

- Learning\_To\_Segment\_Multi-Organ\_and\_ Tumors\_From\_Multiple\_Partially\_CVPR\_ 2021\_paper.html (visited on 09/04/2025).
- L. Jiang, L. Y. Ma, T. Y. Zeng, and S. H. Ying. "UFPS: A unified framework for partially annotated federated segmentation in heterogeneous data distribution". In: Patterns 5.2 (2024). 854 DOI: 10.1016/j.patter.2024.100917.
- O. Petit, N. Thome, and L. Soler. "Iterative [23]confidence relabeling with deep ConvNets for organ segmentation with partial labels". In: 858 Computerized Medical Imaging and Graphics 91 (2021), p. 101938. DOI: 10.1016/j. 860 compmedimag.2021.101938.
- X. Fang and P. Yan. "Multi-Organ Segmentation Over Partially Labeled Datasets With Multi-Scale Feature Abstraction". In: IEEE Transactions on Medical Imaging 39.11 (2020), 865 pp. 3619–3629. DOI: 10.1109/TMI.2020. 866 3001036.
- G. Shi, L. Xiao, Y. Chen, and S. K. Zhou. 868 "Marginal loss and exclusion loss for partially supervised multi-organ segmentation". In: 870 Medical Image Analysis 70 (2021), p. 101979. 871 DOI: 10.1016/j.media.2021.101979.
- Y. Xie, J. Zhang, Y. Xia, and C. Shen. "Learn- 873 ing From Partially Labeled Data for Multi-Organ and Tumor Segmentation". In: IEEE Transactions on Pattern Analysis and Machine Intelligence 45.12 (2023), pp. 14905–14919. 877 DOI: 10.1109/TPAMI.2023.3312587.
- P. Bevandić, M. Oršić, I. Grubišić, J. Šarić, and S. Segvić. "Multi-Domain Semantic Segmentation With Overlapping Labels". In: 2022, 881 pp. 2615-2624. URL: https://openaccess. 882 thecvf . com / content / WACV2022 / html / Bevandic \_ Multi - Domain \_ Semantic \_ Segmentation\_With\_Overlapping\_Labels\_ WACV \_ 2022 \_ paper . html (visited on 09/12/2025).
- [28] S. Hwang, S. Lee, H. Kim, M. Oh, J. Ok, 888 and S. Kwak. "Active Learning for Semantic Segmentation with Multi-class Label Query". 890 In: Advances in Neural Information Pro- 891 cessing Systems 36 (2023), pp. 27020–27039. 892 URL: https://proceedings.neurips. cc / paper \_ files / paper / 2023 / hash / 559a0998fab1d19b80e7e43a5852401c Abstract - Conference . html (visited on 09/14/2025).
- O. Ronneberger, P. Fischer, and T. Brox. "U-898 Net: Convolutional Networks for Biomedical Image Segmentation". In: MICCAI. Vol. 9351. LNCS. Cham: Springer, 2015, pp. 234–241. 901 DOI: 10.1007/978-3-319-24574-4\_28.

963

964

968

969

991

993

994

995

1002

- J. Tian, L. Liu, Z. Shi, and F. Xu. "Automatic Couinaud Segmentation from CT Vol-904 umes on Liver Using GLC-UNet". In: Machine 905 Learning in Medical Imaging. Ed. by H.-I. Suk, 906 M. Liu, P. Yan, and C. Lian. Lecture Notes 907 in Computer Science. Cham: Springer Inter-908 national Publishing, 2019, pp. 274–282. DOI: 909 10.1007/978-3-030-32692-0\_32.
- S. Pieper, M. Halle, and R. Kikinis. "3D [31]911 Slicer". In: 2004 2nd IEEE International 912 Symposium on Biomedical Imaging: Nano to Macro (IEEE Cat No. 04EX821). 2004, 632-914 635 Vol. 1. DOI: 10 . 1109 / ISBI . 2004 . 915 1398617. 916
- [32]M. Antonelli, A. Reinke, S. Bakas, K. Fara-917 hani, AnnetteKopp-Schneider, B. A. Landman, 918 G. Litjens, B. Menze, O. Ronneberger, R. M. 919 Summers, B. van Ginneken, M. Bilello, P. Bilic, 920 P. F. Christ, R. K. G. Do, M. J. Gollub, 921 S. H. Heckers, H. Huisman, W. R. Jarnagin, 922 M. K. McHugo, S. Napel, J. S. G. Pernicka, 923 K. Rhode, C. Tobon-Gomez, E. Vorontsov, H. Huisman, J. A. Meakin, S. Ourselin, M. 925 Wiesenfarth, P. Arbelaez, B. Bae, S. Chen, L. 926 Daza, J. Feng, B. He, F. Isensee, Y. Ji, F. Jia, 927 N. Kim, I. Kim, D. Merhof, A. Pai, B. Park, 928 M. Perslev, R. Rezaiifar, O. Rippel, I. Sara-929 sua, W. Shen, J. Son, C. Wachinger, L. Wang, 930 Y. Wang, Y. Xia, D. Xu, Z. Xu, Y. Zheng, 931 A. L. Simpson, L. Maier-Hein, and M. J. Car-932 doso. The Medical Segmentation Decathlon. 933 Tech. rep. arXiv:2106.05735. arXiv, 2021. URL: 934 http://arxiv.org/abs/2106.05735 (visited 935 on 05/24/2022). 936
- C.-Y. Lee, S. Xie, P. Gallagher, Z. Zhang, and [33]937 Z. Tu. "Deeply-Supervised Nets". In: Proceed-938 ings of the Eighteenth International Confer-939 ence on Artificial Intelligence and Statistics. 940 PMLR, 2015, pp. 562-570. URL: https:// 941 proceedings.mlr.press/v38/lee15a.html (visited on 01/12/2024). 943
- B. Xu, N. Wang, T. Chen, and M. Li. "Em-[34]944 pirical Evaluation of Rectified Activations in 945 Convolutional Network". In: Deep Learning 946 Workshop, ICML 2015 (2015). DOI: 10.48550/ 947 arXiv.1505.00853. 948
- [35]K. He, X. Zhang, S. Ren, and J. Sun. 949 "Deep Residual Learning for Image Recogni-950 tion". In: 2016, pp. 770-778. URL: https: 951 952 //openaccess . thecvf . com / content \_ 953 cvpr \_ 2016 / html / He \_ Deep \_ Residual Learning\_CVPR\_2016\_paper.html (visited 954 on 05/24/2022). 955
- [36]S. Ioffe and C. Szegedy. "Batch Normalization: 956 Accelerating Deep Network Training by Re-957 ducing Internal Covariate Shift". In: Proceed-958 ings of the 32nd International Conference on

- Machine Learning. PMLR, 2015, pp. 448–456. URL: https://proceedings.mlr.press/ v37/ioffe15.html (visited on 06/02/2022).
- I. Loshchilov and F. Hutter. "Decoupled Weight Decay Regularization". In: ICLR 2019 (2019). DOI: 10.48550/arXiv.1711.05101.
- D. Ulyanov, A. Vedaldi, and V. Lempitsky. 966 |38|Instance Normalization: The Missing Ingredient for Fast Stylization. 2016. URL: https: //arxiv.org/abs/1607.08022v3 (visited on 01/12/2024).
- F. C. Monteiro and A. C. Campilho. "Per- 971 formance Evaluation of Image Segmentation". 972 In: Image Analysis and Recognition. Ed. by D. 973 Hutchison, T. Kanade, J. Kittler, J. M. Kleinberg, F. Mattern, J. C. Mitchell, M. Naor, O. Nierstrasz, C. Pandu Rangan, B. Steffen, M. Sudan, D. Terzopoulos, D. Tygar, M. Y. Vardi, G. Weikum, A. Campilho, and M. S. Kamel. 978 Vol. 4141. Berlin, Heidelberg: Springer Berlin Heidelberg, 2006, pp. 248–259.
- S. Nikolov, S. Blackwell, A. Zverovitch, R. 981 Mendes, M. Livne, J. D. Fauw, Y. Patel, C. Meyer, H. Askham, B. Romera-Paredes, C. Kelly, A. Karthikesalingam, C. Chu, D. Carnell, C. Boon, D. D'Souza, S. A. Moinuddin, B. 985 Garie, Y. McQuinlan, S. Ireland, K. Hampton, K. Fuller, H. Montgomery, G. Rees, M. Suleyman, T. Back, C. O. Hughes, J. R. Ledsam, and O. Ronneberger. "Clinically Applicable Segmentation of Head and Neck Anatomy for Radiotherapy: Deep Learning Algorithm Development and Validation Study". In: Journal of Medical Internet Research 23.7 (2021), e26151. DOI: 10.2196/26151.

#### Appendix $\mathbf{A}$

#### A.1Datasets

In our experiments, we leverage multiple datasets that are publicly available, in addition to external test data from Hospital1 (Anonymized). The following sections describe the datasets used in our 1000 experiments. 1001

#### A.1.1Hepatic vessels dataset

The HepaticVessel dataset is from the Medical Seg- 1003 mentation Decathlon challenge [32] and consists of 1004 303 portal-venous phase CT scans from the US. The 1005 dataset has an out-of-plane voxel spacing ranging 1006 from 0.8 to 8.0 mm. The images contain the liver 1007 with segmented liver tumors and vessel structures. 1008

1011

1012

1013

1014

1016

1017

1018

1019

1020

1021

1022

1023

1024

1026

1027

1028

1029

1030

1031

1032

1033

1034

1036

1037

1038

1039

1040

1041

1043

1044

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

1063

1064

1069

#### A.1.2HepaticVessel Liver dataset

Building on the HepaticVessel dataset, a supplementary label set<sup>1</sup> of the liver (HV Liver) and Couinaud segments of the liver was released by Tian et al. [30]. The dataset contains all the same images as the HepaticVessel dataset, but the additional liver and Couinaud segmentation masks are created independently. We leverage the additional liver masks from HV Liver together with the HepaticVessel dataset in our experiments.

#### 3D-ircadb-01 dataset A.1.3

The 3D-ircadb-01-dataset<sup>2</sup> [13] (Ircad), contains 20 CT scans from France that are labeled with various organs, including liver, hepatic vessels, and any liver tumors. The scans from the IRCAD dataset are a subset of the LiTS dataset; however, with only the liver and tumor masks are present, and no vessel masks [9].

#### Liver tumor segmentation (LiTS) A.1.4 dataset

The LiTS dataset [9] contains 131 segmented CT volumes from different patients. The CT scans are from 7 different institutions in Canada, the Netherlands, Germany, France, and Israel. The CT images are contrast-enhanced and captured in the portalvenous phase and have an out-of-plane voxel spacing ranging from 0.7 to 5.0 mm. All images contain a rough segmentation mask of the liver in addition to a radiologist's segmentation of any liver tumors. The liver tumors are both primary and metastatic from colorectal, breast, and lung primary cancers.

The LiTS dataset has 20 volumes (volumes 28-47) overlapping with the 3D-ircadb-01 dataset [13]. This subset contains the same segmented liver and tumor masks as LiTS, in addition to vessel masks, which are used in the LVT application.

#### A.1.5External dataset

The External Dataset (ExDS) is under development at Hospital1 (Anonymized) and Institution1 (Anonymized) and is used for evaluation in our final experiments. The dataset is from a large database of CT images from the follow-up period of 376 patients that was treated for colorectal cancer from 2006 to 2011 at Hospital1. From this database, we have created two labeled subsets: ExDS LT, which contains liver and tumor masks, and ExDS V, which contains liver and vessel masks. The former is used for external validation and testing of the liver and tumor segmentation performance of the model and consists of 18 contrast-enhanced CT volumes with

Table A.1. An overview of the different datasets with partial labels used in this paper and the corresponding class weights  $w_k$  for classes liver, vessel and tumor in Equation 1. HV test and datasets are only used for testing purposes.

Dataset	Liver	Vessel	Tumor	Images	Class weights, $w$
FS LVT	<b>√</b>	<b>√</b>	<b>√</b>	61	[1, 1, 1]
PS L	✓	×	×	61	[1, 0, 0]
PS V	×	$\checkmark$	×	61	[0, 1, 0]
PS T	×	×	✓	60	[0, 0, 1]
HV test	$\checkmark$	$\checkmark$	$\checkmark$	60	_
LiTS	✓	×	✓	131	[1, 0, 1]
HepaticVessel	×	✓	✓	303	[0, 1, 1]
IRCAD	✓	$\checkmark$	$\checkmark$	20	[1, 1, 1]
HV Liver	$\checkmark$	×	×	303	[1, 0, 0]
ExDS LT	<b>√</b>	×	<b>√</b>	18	_
ExDS V	×	$\checkmark$	×	10	_

segmented liver and liver tumor masks. ExDS V 1059 contains 10 contrast-enhanced CT volumes of the 1060 liver with segmented liver vessels and is used to eval- 1061 uate the vessel segmentation performance in Section 1062 4.3.

#### A.1.6 Partial labels in datasets

As we train on datasets with partial labels for certain 1065 experiments, we present an overview of the present 1066 label classes for each dataset and the class weights 1067 w needed for Equation 1 is presented in Table A.1. 1068

#### A.2Experimental setup

All models in this paper are trained on 3D patches of 1070  $128 \times 128 \times 96$  voxels, sampled from training images 1071 resampled to isotropic voxel spacing of  $1 \times 1 \times 1$  mm 1072 using trilinear interpolation. The U-Net-like archi- 1073 tecture uses deep supervision [33] with two auxiliary 1074 heads at intermediate resolutions and LeakyReLU 1075 activations [34]. During training, patches are over- 1076 sampled from a foreground region with p = 0.333, 1077 and we apply the following augmentations in se- 1078 quence: random crop resizing applied with probabil- 1079 ity p = 0.2 and a scale factor  $\alpha \sim U(0.7, 1.4)$ , ran- 1080 dom rotation with p=0.2 and angle  $\beta \sim U(-30,30)$ , 1081 and random flip with p = 0.5 along all axes. We 1082 leverage Random windowing [11] for preprocess- 1083 ing and CT intensity augmentation, applying win- 1084 dow shifting and scaling independently with a total 1085 probability p = 0.3, sampling the Hounfield unit 1086 window parameters from W  $\sim U(11.5, 152.9)$  and 1087  $L \sim [141.2, 325.9]$  [10]. Training is done with the 1088 combined CE and Dice loss (Equation 3 and Equa- 1089) tion 4). All models are trained with a batch size of 1090 112 images across 8 GPU compute dies on 4 AMD 1091 MI250x GPUs on the LUMI supercomputer.

The PS experiments in Section 4.1 and 4.2 are 1093 performed with a residual encoder [35] and batch 1094 normalization [36]. Training is done with AdamW 1095 [37] optimizer with learning rate 0.001 and cosine 1096

<sup>&</sup>lt;sup>1</sup>Available at: https://github.com/GLCUnet/dataset <sup>2</sup>Available at: https://www.ircad.fr/research/ data-sets/liver-segmentation-3d-ircadb-01/

decay with warmup [37]. The models are trained 40 epochs with 100 steps each. The masking weights used for each dataset are listed in Table A.1.

The LVT model trained in Section 4.3 deviates from this setup to match the one used by the nnU-Net baseline. Specifically, we no residual connections in the encoder, instance normalization [38], stochastic gradient descent with weight decay optimizer [37] and polynomial learning rate decay. The model is trained on 448 000 training samples over 1000 epochs, which is comparable to the baseline, which sees 500 000 training samples.

#### A.2.1 Training loss

For each sample with K classes, the cross-entropy loss is defined per voxel as

$$\ell_i^{CE} = -\sum_{k=1}^K y_{ik} \log p_{ik}, \tag{3}$$

where  $y_{ik}$  and  $p_{ik}$  are the target and prediction from the one-hot encoded mask and softmax probabilities, respectively.

For the same network with sigmoid outputs, we obtain the case for binary cross-entropy, where the per-voxel per-class loss is defined as

$$\ell_{ik} = -[y_i \log p_i + (1 - y_i) \log(1 - p_i)]. \tag{4}$$

The dice loss is computed independently for each voxel and class, given the output probabilities, and is given by

$$\ell_{ik}^{Dice} = 1 - \frac{2 \cdot y_{ik} p_{ik}}{y_{ik}^2 + p_{ik}^2}.$$
 (5)

In the multi-class segmentation setup, it is typically reduced over the class dimension to match the pervoxel loss formulation of the  $\ell_i^{CE}$ .

#### A.2.2 Inference settings

As our models are trained on crops smaller than a typical CT image, we follow the sliding window inference pipeline of Isensee et al. [12] to obtain predictions. Specifically, each test volume is cropped into patches of  $128 \times 128 \times 96$  voxels, with 50 % overlap. The model predictions on each patch are aggregated to a complete output with a gaussian weighing, as the predictions are usually more stable towards the center. The final semantic output is obtained through the argmax across channels. For the binary segmentation outputs, we use a sigmoid threshold of p=0.5, and obtain mutually exclusive outputs by giving positives of the overlapping classes priority based on the heuristic hierarchy: tumor, vessel, liver, background.

For comparison with the nnU-Net in Section 5, we use the 5-fold cross-validation models to obtain an

ensemble prediction of each pred. During inference, 1145 we use test-time augmentation by flipping each crop 1146 along all axes. We also limit the final prediction to 1147 the largest connected component. These inference 1148 settings are also employed by the baseline. 1149

#### A.2.3 Evaluation

Precision and recall are common metrics for 1151 evaluating classification performance using the true 1152 positives (TP), false positives (FP), and false nega-1153 tive (FN) predictions. Precision measures the pro-1154 portion of predicted positives that are correct: 1155

$$Precision = \frac{TP}{TP + FP}, (6) 1156$$

while recall (sensitivity) measures the proportion of actual positives that are correctly identified: 1158

$$Recall = \frac{TP}{TP + FN}. (7) {1159}$$

Although pixel-wise precision and recall are not commonly used to evaluate segmentation masks, they 1161 can assist in diagnosing under and over-segmentation 1162 in models. Specifically, low recall tends to correspond to under-segmentation, and low precision to 1164 over-segmentationMonteiro and Campilho [39]. 1165

Dice similarity coefficient (DSC). To evalu- 1166 ate segmentation predictions against ground truth 1167 masks more reliably, we rely on the DSC, which 1168 measures volume overlap between predicted and 1169 true masks, X and Y, as the harmonic mean of the 1170 precision and recall: 1171

$$DSC = \frac{2|X \cap Y|}{|X| + |Y|}.$$
 (8) 1172

Normalized surface dice (NSD). While widely 1173 used, DSC treats all pixel errors equally, which may 1174 obscure clinically important mistakes (e.g., missing 1175 an entire tumor vs. scattered noise). Therefore, we 1176 additionally leverage NSD [40] in our clinical evaluation. NSD addresses this by comparing surfaces 1178 within a tolerance, defined per class in millimeters. 1179 Errors inside the tolerance do not reduce the score, 1180 making NSD more clinically meaningful. 1181

Lesion sensitivity. Based on a connected component analysis of the ground-truth and predicted tumor segmentation, we classify a tumor in the tumor segmentation, we classify a tumor in the tumor ground truth as detected if they have a corresponding prediction with ¿ 10 % overlap. Based on this tumor sitivity using Equation 7.

Reporting In most evaluations, we report the mean result for each model in the 5-fold cross-1190 validation evaluation. To showcase the variation 1191 between multiple runs of comparable methods, we 1192 use the standard deviation of performance between 1193 runs. The result in Table 3 deviates from this pro-1194 tocol, as the whole ensemble is used to obtain each 1195 prediction. We therefore report the per-case mean 1196 and standard deviation for this result.