

This is the peer reviewed version of the following article : Precise ablation zone segmentation on CT images after liver cancer ablation using semi-automatic CNN-based segmentation, Which has been published in final form at <https://doi.org/10.1002/mp.17373> This article may be used for non-commercial purposes in accordance with Wiley Terms and Conditions for Use of Self-Archived Versions. This article may not be enhanced, enriched or otherwise transformed into a derivative work, without express permission from Wiley or by statutory rights under applicable legislation. Copyright notices must not be removed, obscured or modified. The article must be linked to Wiley's version of record on Wiley Online Library and any embedding, framing or otherwise making available the article or pages thereof by third parties from platforms, services and websites other than Wiley Online Library must be prohibited.

1 **Precise Ablation Zone Segmentation on CT Images after**  
2 **Liver Cancer Ablation using Semi-automatic CNN-based**  
3 **Segmentation**

4 Quoc Anh Le<sup>1</sup>, Xuan Loc Pham<sup>2</sup>, Theo van Walsum<sup>3</sup>, Viet Hang Dao<sup>4,5</sup>,  
5 Tuan Linh Le<sup>6</sup>, Daniel Franklin<sup>7</sup>, Adriaan Moelker<sup>3†</sup>, Vu Ha Le<sup>1,2</sup>,  
6 Nguyen Linh Trung<sup>1</sup>, Manh Ha Luu<sup>1,2,3</sup>

7 <sup>1</sup> AVITECH, VNU University of Engineering and Technology, Vietnam

8 <sup>2</sup> Department of Radiology and Nuclear Medicine, Erasmus MC, the Netherlands

9 <sup>3</sup> FET, VNU University of Engineering and Technology, Vietnam

10 <sup>4</sup> Internal Medicine Faculty, Hanoi Medical University, Vietnam

11 <sup>5</sup> The Institute of Gastroenterology and Hepatology, Vietnam

12 <sup>6</sup> Diagnostic Imaging and Interventional Radiology Center, Hanoi Medical University Hospital, Vietnam

13 <sup>7</sup> School of Electrical and Data Engineering, University of Technology Sydney, Australia

14  
15  
16 Version typeset May 27, 2024

17  
18 Corresponding author: Manh Ha Luu. E-mail: halm@vnu.edu.vn

19 Corresponding address: 707 E3, 144 Xuan Thuy road, Cau Giay district, Hanoi, Vietnam

20 Postal code: 100000

## Abstract

22

**Background:** Ablation zone segmentation in contrast-enhanced computed tomography (CECT) images enables the quantitative assessment of treatment success in the ablation of liver lesions. However, fully-automatic liver ablation zone segmentation in CT images still remains challenging, such as low accuracy and time-consuming manual refinement of the incorrect regions.

23

24

25

26

27

**Purpose:** Therefore, in this study, we developed a semi-automatic technique to address the remaining drawbacks and improve the accuracy of the liver ablation zone segmentation in the CT images.

28

29

30

**Methods:** Our approach uses a combination of a CNN-based automatic segmentation method and an interactive CNN-based segmentation method. Firstly, automatic segmentation is applied for coarse ablation zone segmentation in the whole CT image. Human experts then visually validate the segmentation results. If there are errors in the coarse segmentation, local corrections can be performed on each slice via an interactive CNN-based segmentation method. The models were trained and the proposed method was evaluated using two internal datasets of post-interventional CECT images ( $n_1 = 22$ ,  $n_2 = 145$ ; 62 patients in total) and then further tested using an external benchmark dataset ( $n_3 = 12$ ; 10 patients).

31

32

33

34

35

36

37

38

39

**Results:** To evaluate the accuracy of the proposed approach, we used Dice Similarity Coefficient (*DSC*), average symmetric surface distance (*ASSD*), Hausdorff Distance (*HD*), and volume difference (*VD*). The quantitative evaluation results show that the proposed approach obtained mean *DSC*, *ASSD*, *HD*, and *VD* scores of 94.0%, 0.4 mm, 8.4 mm, 0.02, respectively, on the internal dataset, and 87.8%, 0.9 mm, 9.5 mm, and -0.03 respectively, on the benchmark dataset. We also compared the performance of the proposed approach to that of five well-known segmentation methods; the proposed semi-automatic method achieved state-of-the-art performance on ablation segmentation accuracy, and on average, 2 minutes are required to correct the segmentation. Furthermore, we found that the accuracy of the proposed method on the benchmark dataset is comparable to that of manual segmentation by human experts ( $p = 0.55$ , *t*-test).

40

41

42

43

44

45

46

47

48

49

50

**Conclusions:** The proposed semi-automatic CNN-based segmentation method can be used to effectively segment the ablation zones, increasing the value of CECT for assessment of treatment success. For reproducibility, the trained models, source code, and demonstration tool are publicly available at [https://github.com/lqanh11/Interactive\\_AblationZone\\_Segmentation](https://github.com/lqanh11/Interactive_AblationZone_Segmentation).

51

52

53

54

55

# 56 Contents

57	<b>I. Introduction</b>	<b>1</b>
58	<b>II. Methods</b>	<b>5</b>
59	II.A. Automatic segmentation . . . . .	5
60	II.B. Interactive CNN-based segmentation . . . . .	7
61	II.C. Combination scheme . . . . .	8
62	<b>III. Experiments and results</b>	<b>9</b>
63	III.A. Datasets . . . . .	9
64	III.A.1. The details of the datasets: . . . . .	9
65	III.A.2. Annotation and ground-truth . . . . .	11
66	III.A.3. Preprocessing . . . . .	11
67	III.B. Hardware and implementation in details . . . . .	12
68	III.B.1. Implementation of the CNNs . . . . .	12
69	III.B.2. Ablation zone segmentation demonstration tool . . . . .	12
70	III.C. Experiments setup and results . . . . .	13
71	III.D. Evaluation criteria . . . . .	13
72	III.D.1. Automatic ablation zone segmentation . . . . .	14
73	III.D.2. Define optimal model for click-based interactive segmentation . . . . .	15
74	III.D.3. Semi-automatic segmentation performance . . . . .	19
75	<b>IV. Discussion</b>	<b>22</b>
76	<b>V. Conclusions</b>	<b>25</b>
77	<b>Ethical statement</b>	<b>27</b>
78	<b>Conflicts of interest</b>	<b>27</b>
79	<b>Acknowledgments</b>	<b>27</b>
80	<b>References</b>	<b>27</b>

## 1. Introduction

Liver cancer has a high mortality rate, and its incidence increases yearly<sup>1</sup>. According to statistics from GLOBOCAN 2020, liver cancer ranked fifth in the number of new cases and third in the number of cancer deaths<sup>2</sup>. Early detection and treatment of liver cancer are crucial for improving treatment outcomes<sup>3</sup>. Thermal ablation such as radiofrequency ablation (RFA) and microwave ablation (MWA) are considered as curative treatment options for patients with early-stage liver cancer can not undergo an open surgical procedure<sup>4</sup> and can be performed with minimal discomfort for the patient, and with a short recovery time<sup>5</sup>. RFA and MWA also have a low risk of complications compared to surgery and a low risk of side effects compared to chemotherapy and radiation therapy<sup>6</sup>.

During the procedure, the interventionist uses ultrasound (US) or CT to guide the insertion of a thin needle through the patient's skin and into the target lesion<sup>7</sup>. Once the needle tip is placed in position, heat is generated at the needle tip to destroy the malignancy by creating a region of cell destruction, also known as the ablation zone. A major drawback of RFA and MWA is the high recurrence rate after treatment, especially in the local ablation site. It has been reported that ablation of tumors with a size ranging from 2 to 5 cm in diameter results in a recurrence rate of 26.4%<sup>8</sup>.

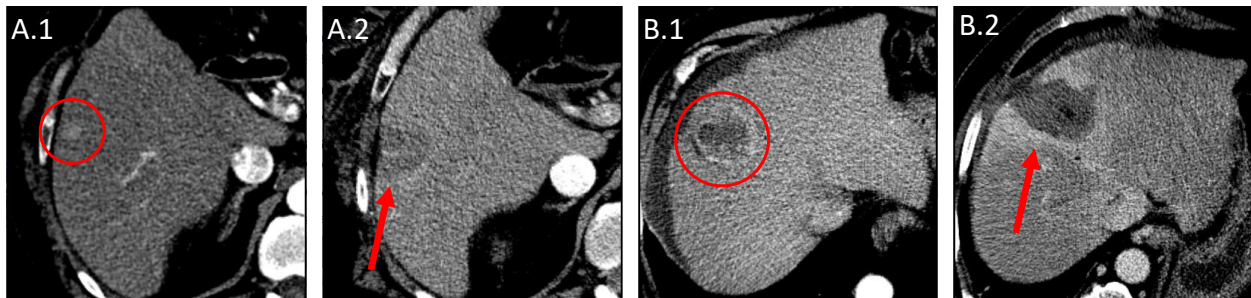


Figure 1: The liver tumor (red circle) in the pre-intervention CT image (A.1 and B.1) and its corresponding ablation zone (red arrow) in the post-intervention image (A.2 and B.2).

To evaluate the ablation, a CECT scan is usually performed at the end of the intervention to visualize the ablation zone (see Figure 1)<sup>9</sup>. By assessing the ablation zone in the CECT image, the physician can determine whether the procedure was completely successful, whether additional ablation needs to be performed, and what the safety margins of the ablation zone correspond to the lesion<sup>10</sup>. It is especially important to accurately assess the margins of the ablation zone. In current clinical practice, this assessment is performed visually by the interventional radiologist. Precise segmentation of the ablation zone in the CECT image enables quantitative assessment and may provide improved confidence in the outcome.

104 While manual segmentation of the ablation zone is tedious and time-consuming and thus is not  
105 feasible in clinical practice, computer-aided methods can be used to greatly reduce the required cognitive  
106 load. The main challenge of precise ablation zone segmentation is that the post-interventional CT  
107 images are noisy, with inhomogenous intensities inside the ablation zone. In addition, for clinical use,  
108 the segmentation processing time should be sufficiently fast (in order of minutes). Although there  
109 have been many studies developing methods for liver lesion segmentation<sup>11</sup>, only a few studies have  
110 focused on ablation zone segmentation. Egger et al. (2015)<sup>12</sup> presented a semi-automatic method,  
111 based on an interactive, graph-based contouring approach, for segmenting the ablation zone in twelve  
112 post-interventional CECT images. The obtained segmentation is compared to manual segmentations  
113 performed by two medical experts, achieving a mean Dice Similarity Coefficient (*DSC*) score of 77%.  
114 Wu et al. (2021)<sup>13</sup> applied a region-growing method to segment the ablation zone in the CT image,  
115 followed by fuzzy c-mean clustering and then refined by cyclic morphological processing, which achieved  
116 the mean *DSC* score of 75%. Recently, deep learning-based methods have been applied to segment the  
117 ablation zone. He et al. (2021)<sup>14</sup> used a multi-scale patch-based 3D Residual Attention U-Net (RA-  
118 UNet) to segment the ablation zone in multiphase CT images, achieving median *DSC* scores of 83% and  
119 89% for arterial and portal venous phase, respectively. It was reported that the hepatic enhancement in  
120 the arterial phase is not sufficient to discriminate the ablation zone's edge, resulting in the difficulty for  
121 the segmentation. Anderson et al. (2022)<sup>15</sup> investigated Hybrid-WNet to segment the ablation zone in  
122 CECT images, reporting a median *DSC* score of 79% and a median surface distance of 0.76 mm. From  
123 the reported state-of-the-art results, it is clear that precise segmentation of the ablation zone in CT  
124 images remains a challenging problem. In addition, to the best of our knowledge, none of the previous  
125 studies evaluated the methods on external datasets. Therefore, the purpose of this study is to propose  
126 and assess an effective method for precise ablation zone segmentation. In addition, we will evaluate the  
127 method using both internal and external datasets to verify the performance of the method.

128 In the last decade, deep learning-based methods, especially Convolutional Neural Networks (CNN)  
129 and Transformers, have demonstrated their indisputable effectiveness across a variety of fields including  
130 medical image analysis. Ronneberger et al. (2015)<sup>16</sup> introduced U-Net, which consists of a combi-  
131 nation of encoder-decoder structures for efficient automatic segmentation of medical images. Since  
132 then, numerous variants of U-Net have been proposed, such as Residual U-Net<sup>17</sup>, which integrates the  
133 residual path into the original structure, and H-DenseUNet<sup>18</sup>, which combines 2D and 3D Dense U-Net.  
134 Isensee et al (2021)<sup>19</sup> further demonstrated the effectiveness of U-Net in medical image segmentation  
135 via the release of nnU-Net, which has become a well-known platform for automatic training and organ

136 segmentation pipelines. Recently, the Vision Transformer technique (ViT) proposed by Dosovitskiy et  
137 al. (2020)<sup>20</sup> demonstrated its potential in computer vision, and is now being applied to the problem  
138 of medical image segmentation. For instance, UNETR<sup>21</sup> leverages the U-shaped encoder-decoder ar-  
139 chitecture but replaces the CNN-based encoding branch with ViT. Chen et al. (2021)<sup>22</sup> also realize  
140 the potential of the ViT-based encoder scheme and introduce TransUNet, but instead of a pure ViT  
141 encoding branch, the authors propose a hybrid architecture where a multi-resolution CNN is initially  
142 utilized to produce input feature maps for the ViT encoding block. CoTr<sup>23</sup> further improves the hybrid  
143 architecture in TransUNet by fully leveraging the multi-scale feature maps from the CNN instead of only  
144 using the lowest resolution. Subsequently, these multi-scale feature maps are flattened and encoded by  
145 Deformable ViT, which greatly reduces the complexity of ViT and therefore significantly boosts its  
146 performance.

147 Unlike automatic segmentation methods, *semi-automatic* segmentation methods require some  
148 level of human interactions to complete the segmentation procedure. The interaction can be point-  
149 clicking, scribbling, rectangle/circle initial drawing, and manual tuning of parameter values. Several  
150 semi-automatic segmentation methods have been developed for medical image segmentation<sup>24</sup>. Region-  
151 growing is one of the most popular semi-automatic segmentation methods<sup>25</sup>. Region-growing is initial-  
152 ized by a seed point with a predefined threshold interval and then expands within a connected region.  
153 The contrast between the object and the background acts as an important factor for a successful  
154 segmentation. Chan and Vese (2001)<sup>26</sup> introduced an interactive segmentation method based on a  
155 level-set algorithm, where the user provides an initial contour around the object. The method can suc-  
156 cessfully segment objects without clear boundaries. Grab-cut<sup>27</sup> is a well-known interactive segmentation  
157 method in which the user provides source and sink regions via manual interactions. Furthermore, sev-  
158 eral conventional methods such as the Robust Statistics Segmenter<sup>28</sup>, Otsu & Picking<sup>29</sup>, and Geodesic  
159 Segmenter<sup>30</sup> were investigated for interactive segmentation of objects in medical images.

160 Recently, interactive CNN-based segmentation methods have been investigated and shown to out-  
161 perform traditional interactive methods, achieving higher accuracy with fewer user interactions. Deep-  
162 Cut<sup>31</sup> and ScribbleSup<sup>32</sup> are among the first interactive CNN-based segmentation frameworks which  
163 embed user-provided bounding boxes or scribbles into CNN models. DeepGeoS<sup>33</sup> performed interactive  
164 segmentation by using geodesic distance transforms of scribbles as additional input channels to CNNs.  
165 Furthermore, Wang et al. (2018)<sup>34</sup> proposed an image-specific fine-tuning method to make a CNN  
166 model adaptive to a specific test image. Although promising results have been reported, the method  
167 still has some limitations. For example, it requires the user to provide a bounding box for the object and

---

168 scribbling on the background and foreground for the correction which may be inconvenient in practice.  
169 Luo et al. (2021)<sup>35</sup> introduced MIDeepSeg, a minimally interactive segmentation method based on a  
170 CNN and exponentiated geodesic distance, to segment both seen and unseen objects appearing in the  
171 training dataset. However, the framework requires the user to click several points at the edges of the  
172 object on each slice the object presents to create a ROI for the object, which may be inconvenient to  
173 use in applications with tight timing constraints. Sun et al. (2022)<sup>36</sup> proposed a graph convolutional  
174 neural network for segmentation tasks; however, the requirement of clicking points at the boundaries or  
175 dragging predicted points is also not well-suited to time-critical tasks. Sofiiuk et al. (2022)<sup>37</sup> proposed  
176 a click-based interactive segmentation, Reviving Iterative Training with Mask Guidance (RITM), that  
177 iteratively uses the differences of the previous prediction segmentation and the ground truth to provide  
178 additional prior information to train the model and improve segmentation prediction accuracy.

179 Generally, fully automatic segmentation methods are convenient and fast for global segmentation  
180 of the entire image. However, they frequently have small segmentation errors that need to be manually  
181 corrected. Theoretically, for CNN-based approaches, the more data involved in training the deep learning  
182 models, the more accurate the model can become. However, it is not clear how much data should be  
183 used for a specific medical image segmentation application. Furthermore, collecting large amounts  
184 of data in the medical image field with accurate segmentations is challenging. In contrast, a human  
185 expert can control an interactive method to segment a region correctly. Nevertheless, using interactive  
186 segmentation methods to fully segment a structure may not be practical because of the excessive  
187 amount of interactions required. In addition, we hypothesize that clicking points inside a region is more  
188 convenient than clicking points at the edges, or scribbling using the mouse. *Therefore, our key idea to*  
189 *solve the problem of ablation zone segmentation is to utilize an interactive CNN-based segmentation*  
190 *method in which the user clicks points at the incorrect segmentation regions to refine the segmentation*  
191 *provided by an automatic segmentation method.*

192 The overview of the remainder of the paper is as follows: Section II. describes the proposed  
193 interactive CNN-based segmentation method in detail. Subsequently, Section III. describes extensive  
194 experiments to assess the performance of the proposed solution. Next, the experimental results are  
195 discussed in Section IV.. Finally, Section V. summarizes the findings of this study.



## 196 II. Methods

197 In order to precisely segment the ablation zone in a CECT image, our strategy is to combine automatic  
198 segmentation with interactive segmentation. Firstly, the automatic segmentation method is applied to  
199 segment the ablation zone in the whole CT volume. Next, a human expert reviews the automatic seg-  
200 mentation and then uses RITM<sup>37</sup> as a CNN-based interactive segmentation method to fix the incorrect  
201 segmentations via clicking points at the local error locations. Finally, the segmentations are combined  
202 by mixing the probability maps at the local location of the two methods. The underlying assumption  
203 is that the more accurate the automatic segmentation is, the fewer human interactions are needed for  
204 error correction. In this study, we evaluate several automatic segmentation frameworks for the initial  
205 segmentation in Section III.D.1.. In addition, the interactive method enables click-based segmentation  
206 which is one of the simplest interaction types. The pipeline of the proposed approach is illustrated in  
207 Figure 2. The following sections describe each component in the pipeline in detail.

### 208 II.A. Automatic segmentation

209 In the first step of the proposed pipeline, an automatic segmentation network is utilized to segment the  
210 ablation zone from the CT image. In this study, we evaluate well-known CNN-based and Transformer-  
211 based networks to find a suitable one, aiming for fast inference time and high accuracy. Here are  
212 descriptions of four segmentation methods.

- 213 • 3D U-Net, introduced by Cicek et al. (2016)<sup>38</sup>, is an extension of U-Net architecture designed  
214 to segment objects in 3D data by processing them with corresponding 3D operations. 3D U-Net  
215 consists of an encoder-decoder structure with a skip connection to capture high-level and low-  
216 level features of the 3D image and produce full-resolution segmentation. The 3D U-Net has been  
217 extensively used in various medical image segmentation tasks.
  - 218 • UNETR, proposed by Hatamizadeh et al. (2022)<sup>21</sup>, leverages the potential of ViT in sequence  
219 representation learning, making it highly effective in segmenting objects in images. UNETR  
220 utilizes the strength of the U-Net architecture but replaces the CNN-based encoder with the  
221 Transformers-based encoder. The original image is split into 3D patches and a linear projec-  
222 tion of these patches is applied to produce the input for the ViT-based encoder. UNERT has  
223 demonstrated state-of-the-art performance in several medical image segmentation tasks.
-

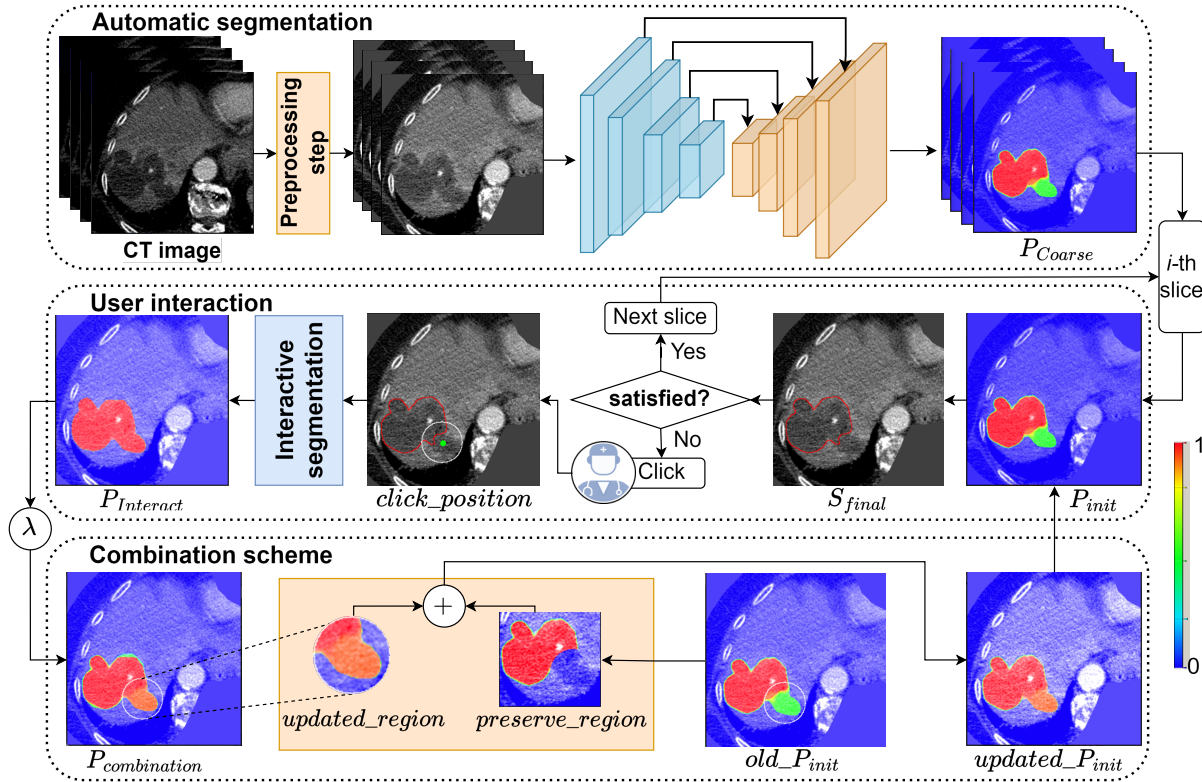


Figure 2: The pipeline of the proposed approach for the semi-automatic ablation zone segmentation. The 3D coarse segmentation is predicted using nnU-Net, which is then locally corrected in each slice by the user using the interactive segmentation via the combination scheme. The probability map, the output of segmentation models, overlapped with the CT image and visualized by a color map. The color bar indicates the probability prediction value of the segmentation. The white circle marks the region of interest.

224 • nn-UNet, developed by Isensee et al.(2021)<sup>19</sup>, focuses on optimizing and improving the perfor-  
 225 mance of the U-Net architecture by introducing various enhancements. nn-UNet employs novel  
 226 data augmentation techniques, training strategies, and model configurations to achieve better  
 227 segmentation results.

228 • CoTr, proposed by Xie et al. (2021)<sup>23</sup>, also consists of an encoder-decoder structure like other  
 229 segmentation networks. In the encoder part, CNNs and Transformers are used in the CoTr's  
 230 architecture. The CNN's role is to extract feature representations from the input image. Then,  
 231 the deformable Transformer (DeTran) is used to model long-range dependencies within the  
 232 extracted feature maps. Combining the strengths of CNNs and Transformers, CoTr addresses  
 233 complex tasks requiring local and global context understanding in image analysis.

## 234 II.B. Interactive CNN-based segmentation

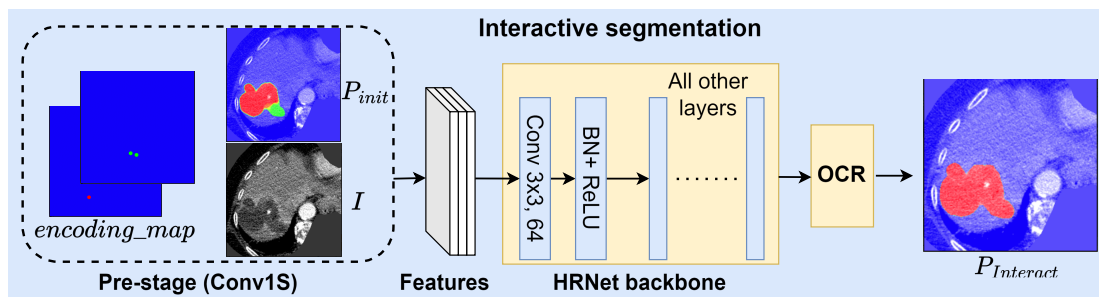


Figure 3: The architecture of RITM. User click points are encoded in binary disks. The positive and negative click points are shown in green and red in the *encoding\_map*, respectively.

235 We aim to use slice-based segmentation for the correction as in clinical routine: the medical expert  
 236 needs to check every single axial slice containing the ablation zone after the treatment<sup>12</sup>.

237 As post-interventional CT images are characterized by inhomogeneous intensities and noise in the  
 238 ablation zone (see Figure 1), traditional methods may not perform well in such conditions. It has been  
 239 demonstrated that CNN-based segmentation methods are able to deal with inhomogeneous regions and  
 240 noise<sup>39</sup>. Therefore, we use RITM<sup>37</sup>, a 2D interactive CNN-based segmentation method, to revise the  
 241 prediction of the automatic segmentation and obtain the final ablation zone segmentation. The RITM  
 242 consists of two parts: the click encoding block and the backbone, as shown in Figure 3. The idea  
 243 of the interactive segmentation network is to encode the user clicks and feed them into the network’s  
 244 backbone to generate a prediction. We encode the user click in a binary disk, which achieved effective  
 245 performance compared to other click encoding schemes<sup>40</sup>.

246 As a segmentation network, RITM also contains a semantic segmentation backbone. In this study,  
 247 we choose High-Resolution Net combined with Object-Contextual Representations (HRNet+OCR) as  
 248 the backbone of the interactive segmentation method. The HRNet+OCR is a promising architecture  
 249 specifically designed for producing high-resolution outputs<sup>34</sup>. The HRNet+OCR backbone model was  
 250 pre-trained using the ImageNet dataset, meaning the backbone’s input is a three-channel image. How-  
 251 ever, in the interactive segmentation, the input includes additional features such as a guided mask and  
 252 the click-encoding map. To adapt the pre-trained model, we employ a convolution block known as  
 253 Conv1S, introduced by Sofiiuket al. (2022)<sup>37</sup>. The architecture of Conv1S is designed to ensure that  
 254 the output feature’s channel matches the input of the backbone’s first convolutional layer (64 channels).

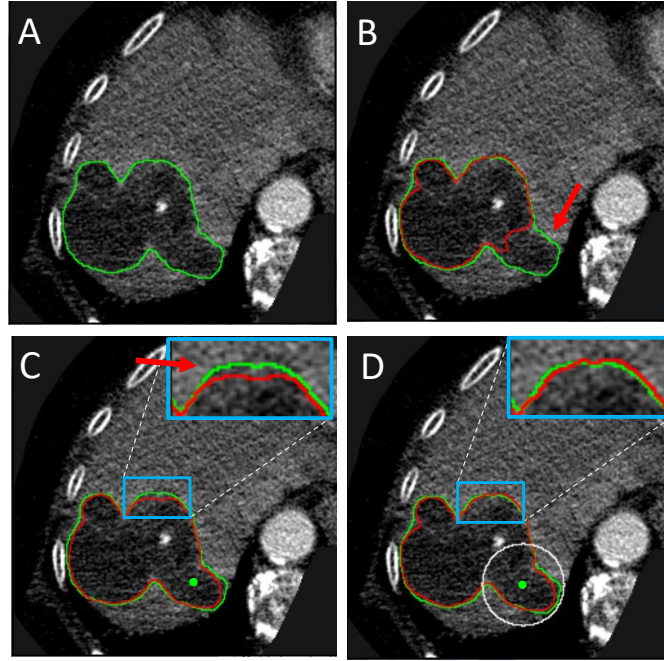


Figure 4: Illustration of interactive segmentation using the click-based semi-automatic method w and w/o the combination scheme: the original image with ground truth segmentation (green contour) of an ablation zone (A), the segmentation (red contour) obtained from the automatic model (B), the segmentation (red contour) obtained from the original interactive model (C), and the segmentation (red contour) obtained using the proposed combination scheme (D). The green dot represents the user’s click. The white circle is the ROI. The red arrows mark the mislabeled regions.

### 255 II.C. Combination scheme

256 Since the segmentation obtained from the automatic model might have mislabeled regions (see Figure  
 257 4.B), the interactive model is applied to revise the prediction of the automatic segmentation. Nev-  
 258 ertheless, the segmentation obtained from the interactive model may also have mislabeled regions far  
 259 from the clicked points (see Figure 4.C). Consequently, we combine the two CNN-based segmentations  
 260 using their probability predictions with a spatial constraint from the clicked points. The main idea is  
 261 to construct a weighted voter between the two probability predictions within the regions of interest  
 262 (ROIs), assuming that the user clicks inside the incorrect segmentation regions. The pseudo-code of  
 263 the combination scheme is shown in Algorithm 1.

264 Firstly, the automatic segmentation network coarsely predicts the ablation zone from the CT slice  
 265  $I \in \mathbb{R}^{512 \times 512}$  to obtain the probability prediction  $P_{coarse} \in \Omega^{512 \times 512}, \Omega \triangleq [0, 1]$  which is then defined  
 266 as the initial prediction  $P_{init}$  for the combination scheme. Next, the initial prediction is thresholded to  
 267 obtain the final segmentation  $S_{final} \in \Theta^{512 \times 512}, \Theta \triangleq \{0, 1\}$ . If the segmentation is not satisfying, the

**Algorithm 1** Combination scheme

---

```

1: Input: CT slice:  $I$ , Coarse prediction:  $P_{coarse}$ , Kernel size:  $K$ , Weighted value:  $\lambda$ ,
   Threshold value:  $thrsh$ 
2: Output: Final segmentation:  $S_{final}$ 
3:  $P_{init} \leftarrow P_{coarse}$ 
4:  $S_{final} \leftarrow P_{init} > thrsh$ 
5: while  $S_{final}$  is unsatisfied do
   /* User define a click in the mislabeled region */
6:    $encoding\_map, click\_positon \leftarrow \text{USER\_CLICK}()$ 
7:    $P_{interact} \leftarrow \text{PREDICTOR}(I, P_{init}, encoding\_map)$ 
8:    $circle\_mask\_1, circle\_mask\_0 \leftarrow \text{CIRCLE}(click\_position, K)$ 
   /* Update  $P_{init}$  */
9:    $P_{combination} \leftarrow \lambda \times P_{interact} + (1-\lambda) \times P_{init}$ 
10:   $updated\_region \leftarrow P_{combination} \times circle\_mask\_1$ 
11:   $preserve\_region \leftarrow P_{init} \times circle\_mask\_0$ 
12:   $P_{init} \leftarrow updated\_region + preserve\_region$ 
   /* Thresholding to get  $S_{final}$  */
13:   $S_{final} \leftarrow P_{init} > thrsh$ 
14: end while

```

---

268 user corrects the mis-segmented regions using a negative point click to correct a false positive region  
269 or a positive point click to correct a false negative region. Based on the clicked points, the interactive  
270 segmentation network generates a probability prediction (PREDICTOR), which refers to the interactive  
271 prediction  $P_{interact} \in \Omega^{512 \times 512}$ . In addition, we define ROIs from the positions of user clicks by dilating  
272 the clicking positions with a kernel size  $K$  (CIRCLE). The initial probability prediction is then updated  
273 by combining the previous initial probability prediction  $P_{init}$  and the interactive probability prediction  
274  $P_{interact}$  with a weighted parameter  $\lambda$  within the local correcting ROIs. Subsequently, the ablation zone  
275 segmentation is corrected only in the ROIs. To this end, the process is repeated until the final ablation  
276 zone segmentation  $S_{final}$  is satisfied. The effect of weighted parameter  $\lambda$  and kernel size  $K$  on the  
277 segmentation accuracy will be assessed in Section III.D.2..

## 278 III. Experiments and results

### 279 III.A. Datasets

#### 280 III.A.1. The details of the datasets:

281 This study involved three datasets, comprising a total of 179 contrast-enhanced CT scans, from two  
282 medical centers.

---

283 The first dataset, denoted as  $EMC_A$  dataset, includes CT scan images (arterial phase) from 22  
 284 patients who underwent ablation treatment of liver lesions. The CT scans were acquired at Erasmus  
 285 MC using Siemens CT scanners and reconstructed with an axial matrix size of  $512 \times 512$ , with a pixel  
 286 spacing of 0.8 mm, a slice thickness of 3 mm, and the number of slices ranged from 40 to 70 slices.

287 The second dataset, referred to as  $EMC_B$  dataset, was reused from our prior study<sup>41</sup> and consists  
 288 of 145 multiphase CT scans from 40 patients who underwent ablation treatment of liver cancer at  
 289 Erasmus MC. The CT scan was acquired while patients were in the intra-intervention room and patients  
 290 went to the medical center for the follow-up procedure. The CT scans were reconstructed with an axial  
 291 matrix size of  $512 \times 512$ , with the slice thickness ranging from 0.4 to 5 mm, the pixel spacing ranging  
 292 from 0.59 to 0.98 mm, and the number of slices ranging from 19 to 672 slices.

293 The third dataset, the Benchmark dataset, was obtained from the Medical University of Leipzig,  
 294 Saxony, Germany and included 12 CT scans from 10 patients<sup>12</sup>. The CT scans are acquired in the  
 295 portal venous phase, reconstructed with an axial matrix size of  $512 \times 512$ , pixel spacing ranging from  
 296 0.68 to 0.78 mm, slice thickness ranging from 1 to 2 mm, and the number of slices ranging from 52 to  
 297 232.

298 Each CT scan contains one to three ablation zones. The mean diameter of the ablation zones in  
 299 the training set was  $55 \pm 18.6$  mm. In 12 CT scans in the training set, the needle is visible. In the  
 300 validation set, the mean diameter of the ablation zones is  $54.6 \pm 23.8$  mm. Three CT scans in the  
 301 validation set contain the needle. In the EMC testing set, the mean diameter of the ablation zones was  
 302  $54.5 \pm 18.4$  mm. The needle was visible in seven CT scans in this set. For the Benchmark dataset, the  
 303 mean diameter of the ablation zones was  $62.2 \pm 16.3$  mm. The needle is visible in six CT scans in this  
 304 dataset.

305 The specifics of data division can be found in Table 1.

Table 1: Number of CT volumes and slices for training, validation, and testing used in this study. The numbers in parentheses are the number of 2D slices.

	Dataset	# Patients	Arterial	Portal venous	Total
<b>Training</b>	$EMC_A$	18	18	-	18
	$EMC_B$	29	36 (5638)	35 (5379)	71 (11017)
<b>Validation</b>	$EMC_A$	4	4	-	4
	$EMC_B$	18	5 (862)	13 (1353)	18 (2215)
<b>Testing</b>	$EMC_B$	11	31 (3006)	25 (3141)	56 (6147)
	Benchmark	10	-	12 (1525)	12 (1525)

### 306 III.A.2. Annotation and ground-truth

307 For the ablation zone of  $EMC_A$  and  $EMC_B$  datasets, a technician manually created segmentations using  
308 Mevislab version 3.4.3. Next, the segmentations were corrected/verified by a physician, serving as the  
309 ground truth. In the Benchmark dataset, the ground truth of the ablation zones was obtained from two  
310 medical experts who manually drew segmentations using Mevislab<sup>12</sup>.

### 311 III.A.3. Preprocessing

312 In the preprocessing step of the CT scan, we limited the prediction range to the liver region. Firstly,  
313 liver segmentation is automatically acquired using a nnU-Net segmentation network, which has been  
314 trained on the LiTS dataset<sup>11</sup>. Then, the liver segmentation is dilated with a kernel size of (30, 30, 1).  
315 The dilated liver segmentation is used as a mask for the ablation zone segmentation.

316 In the interactive segmentation, we truncated the CT image using a HU range, which is defined in  
317 Section III.D.2. for an optimal value. Then, we converted the clipped CT images into a three-channel  
318 format using the *OpenCV* library for input into the interactive segmentation model.

319 The preprocessing step of the automatic segmentation methods is performed using the default  
320 settings described in the original articles. For nnU-Net, nnU-Net (fine-tuning), and CoTr, these models  
321 are implemented within the nnU-Net framework. Consequently, they share similar preprocessing steps,  
322 which are automatically determined based on the characteristics of the training dataset. Firstly, all data  
323 is cropped to retain only the region containing nonzero values. Subsequently, the data is resampled  
324 to the median voxel spacing of the entire training dataset, where third-order spline interpolation and  
325 nearest interpolation are applied for the image data and segmentation label, respectively. Following  
326 resampling, the data is normalized by clipping intensity values to the [0.5, 99.5] percentiles of the entire  
327 training dataset's intensity range. This is based on z-score normalization, which is computed based on  
328 the mean and standard deviation of all intensity values collected from the dataset.<sup>42</sup> For 3D U-Net and  
329 UNETR, the models are implemented using the MONAI library. The preprocessing steps include various  
330 data transformations, which are resampled to a fixed spacing, clipped, and scaled the intensity into the  
331 range of 0 to 1<sup>43</sup>. Since the training samples are not extensive, data augmentation is also applied to  
332 the data to prevent the overfitting problem, which are random crop, random rotate, random flip, and  
333 random elastic deformations<sup>42</sup>.

---

## 334 III.B. Hardware and implementation in details

### 335 III.B.1. Implementation of the CNNs

336 This study is conducted using a Ubuntu 20.04 workstation, with an Intel® Core™ i9-10900K CPU, 64GB  
337 RAM, RTX 8000 GPU. The source code is implemented in Python 3.6 with Pytorch 1.10 integrated  
338 with CUDA 11.3. The implementations of automatic segmentation models and RITM are based on  
339 their authors' GitHub repositories <sup>a b c d</sup>.

340 We implemented well-known CNN-based and Transformer-based networks, including 3D U-Net<sup>38</sup>,  
341 nnUNet<sup>19</sup>, UNETR<sup>21</sup>, and CoTr<sup>23</sup> for automatic segmentation performance comparison. We trained  
342 CNN-based and Transformer-based models with the training EMC<sub>B</sub> dataset mentioned in Section III.A..  
343 For nnU-Net and CoTr, two models were constructed using the same self-configuration framework for  
344 training and testing. For 3D U-Net and UNETR, we trained the models using the tutorials from MONAI  
345 with default parameters. To assess the performance of the automatic segmentation model when more  
346 data is involved in the training model, we conducted an experiment with nnU-Net. We trained the nnU-  
347 Net model with the EMC<sub>A</sub> dataset and then employed the fine-tuning technique to train the model with  
348 the EMC<sub>B</sub> dataset; we refer to this experiment as nnU-Net (fine-tuning). All automatic segmentation  
349 models were trained for 1000 epochs, and the training time for nnU-Net and CoTr was approximately  
350 30 hours, while the training time for U-Net and UNETR was about 25 hours.

351 To train the RITM model, an interactive sampling procedure is required. We reused the procedure  
352 described in the original RITM paper, in which the sampled point is obtained by applying a morphological  
353 erosion operation of the mislabeled region<sup>37</sup>. In addition, we used the *DiceCE* loss function, which was  
354 used in the automatic segmentation network nnU-Net. We reuse the RITM model, which was trained on  
355 the COCO+LVIS dataset<sup>44,45</sup> as the pre-trained model for the ablation zone segmentation task. Note  
356 that before the training stage, the backbone HRnet was already pre-trained with the ImageNet dataset.  
357 We then employed the fine-tuning technique to train the RITM model with the EMC<sub>B</sub> training dataset.  
358 We trained RITM for 500 epochs with the default parameters from the original article.

### 359 III.B.2. Ablation zone segmentation demonstration tool

360 For an easy demonstration, we adopt the demonstration tool, which was created by Sofiiuk et al.  
361 (2022)<sup>37</sup>, using the *Tkinter* library. Nevertheless, the original demonstration tool was initially designed  
362 for 2D images, which necessitates modifications for working with CT scans. The most significant

---



363 modification involves embedding the proposed method, incorporating 3D CT scans and enabling users  
 364 to adjust the display of slices in the 3D CT scan using the mouse. The demonstration tool and video  
 365 are publicly available at Github <sup>e</sup>.

### 366 III.C. Experiments setup and results

### 367 III.D. Evaluation criteria

368 In this study, we use the Dice Similarity Coefficient (*DSC*), Average Symmetric Surface Distance  
 369 (*ASSD*), Hausdorff Distance (*HD*), and Volume difference (*VD*) as metrics for evaluation of the pro-  
 370 posed methods.

- 371 • **Dice similarity coefficient (DSC):** Suppose  $A$  and  $B$  represent the ground truth and predicted  
 372 segmentation of the ablation zone of a 3D CT image, respectively. The *DSC* measures how good  
 373 the overlap between  $A$  and  $B$  is. The *DSC* value of 0 means no overlap, and 100 means perfect  
 374 overlap. The more overlap between  $A$  and  $B$ , the closer the *DSC* score to 100%.

$$375 \quad DSC(A, B) = \frac{2|A \cap B|}{|A| + |B|} \times 100\% \quad . \quad (1)$$

- 376 • **Average symmetric surface distance (ASSD):** Suppose  $S(A)$ ,  $S(B)$  represent all surface voxels  
 377 on the ground truth ( $A$ ) and the predicted ablation zone segmentation ( $B$ ). Voxel  $v_A$  and  $v_B$  are  
 378 arbitrary voxels belonging to  $A$  and  $B$ , respectively. We define the shortest path from  $v_A$  to  $S(B)$   
 379 or  $v_B$  to  $S(A)$  as follows:

$$380 \quad d(v_{A'}, S(B)) = \min_{v_{BA} \in S(B)} \|v_A - v_{BA}\| \quad , \quad (2)$$

$$381 \quad d(v_{B'}, S(A)) = \min_{v_{AB} \in S(A)} \|v_B - v_{AB}\| \quad , \quad (3)$$

383 where  $v_{BA}$  means the point in  $S(B)$  that draws the shortest distance from point  $v_A$  and similar  
 384 for  $v_{AB}$ . The *ASSD* metric measures the average gap between the boundary of  $A$  and  $B$ . The  
 385 formula for *ASSD* is written as follows:

$$386 \quad ASSD(A, B) = \frac{\sum_{v_A \in S(A)} d(v_{A'}, S(B)) + \sum_{v_B \in S(B)} d(v_{B'}, S(A))}{S(A) + S(B)} \quad . \quad (4)$$

- 387 • **Hausdorff distance (HD):** The *HD* shows the maximum distance between the boundary of  $A$   
 388 and  $B$ . The formula for *HD* is as follows:

$$389 \quad HD(A, B) = \max \left( \max_{v_A \in S(A)} d(v_{A'}, S(B)), \max_{v_B \in S(B)} d(v_{B'}, S(A)) \right) \quad . \quad (5)$$

- 390 • **Volume Difference (VD)**: Suppose  $V_A$  and  $V_B$  represent the volumes of the ground truth and  
 391 the predicted segmentation of the ablation zone in a 3D CT image, respectively. The  $VD$  metric  
 392 measures the volume difference between these two volumes without considering their overlap. A  
 393  $VD$  value close to 0 indicates that the size of the prediction closely matches the size of the ground  
 394 truth. The formula for  $VD$  is as follows:

$$395 \quad VD(A, B) = \frac{2(V_B - V_A)}{V_B + V_A} . \quad (6)$$

- 396 • **Precision**: The formula for Precision is written as:

$$397 \quad Precision = \frac{TP}{TP + FP} , \quad (7)$$

398 where  $TP$  is the number of correctly identified ablation voxels, and  $FP$  is the number of over  
 399 segmentation voxels of the predicted ablation zone segmentation.

- 400 • **Recall**: The formula for Recall is written as:

$$401 \quad Recall = \frac{TP}{TP + FN} , \quad (8)$$

402 where  $FN$  is the number of incorrectly identified ablation voxels.

403 In addition, we also use the Area under a curve metric ( $AUC$ ) to evaluate the performance of the  
 404 segmentation methods.

### 405 III.D.1. Automatic ablation zone segmentation

406 In this experiment, we investigate the performance of four state-of-the-art automatic segmentation net-  
 407 works on ablation zone segmentation: 3D U-Net<sup>38</sup>, UNETR<sup>21</sup>, nnU-Net, nnU-Net (fine-tuning)<sup>19</sup> and  
 408 CoTr<sup>23</sup>. The experimental results regarding the comparison of the automatic ablation zone segmen-  
 409 tation of four well-known methods are summarized in Table 2. The evaluation is based on the three  
 410 metrics:  $DSC$ ,  $HD$  and  $ASSD$ , and on the two test sets from the  $EMC_B$  testing dataset, with the final  
 411 segmentation achieved using a threshold of 0.5. We also list the number of failed cases, when there is  
 412 no overlap between the ground truth and the predicted segmentation. The highest mean  $DSC$  values  
 413 are 81.2% and 88.4% for CT images acquired in the arterial and portal venous phases, respectively.

414 Regarding the processing time, UNETR has the lowest processing time (4-5 seconds for a CT volume  
 415 on average), while nnU-Net (fine-tuning) requires a slightly longer processing time of approximately 7

Table 2: Performance comparison of ablation zone segmentations among the well-known automatic segmentation methods. The bold numbers are the highest mean scores.

Dataset	Method	DSC	HD (mm)	ASSD (mm)	VD	# failure	Processing time (s)
EMC (Arterial) n = 31	3D U-Net	61.7 ± 18.7	66.7 ± 78.1	8.2 ± 8.8	0.61 ± 0.39	13	4.4 ± 3
	UNETR	48.8 ± 20.8	203.1 ± 64.8	16.7 ± 11.2	0.56 ± 0.46	9	3.8 ± 2
	nnU-Net	<b>81.2 ± 12.8</b>	34.1 ± 29.7	<b>2.9 ± 4.5</b>	<b>0.26 ± 0.29</b>	1	31.9 ± 11.7
	CoTr	76.6 ± 20.6	39.4 ± 29.7	3.6 ± 5.8	0.34 ± 0.39	1	25.3 ± 9.1
	nnU-Net (fine-tuning)	80 ± 15.6	<b>29.9 ± 30.3</b>	3.4 ± 5.4	0.3 ± 0.33	2	7.1 ± 3.5
EMC (Portal venous) n = 25	3D U-Net	58.4 ± 26.9	125.3 ± 94.9	12.5 ± 16.1	0.62 ± 0.54	1	5.3 ± 3.8
	UNETR	52.9 ± 20.9	240 ± 84.2	23.4 ± 21.4	0.34 ± 0.58	0	4.9 ± 4
	nnU-Net	86.1 ± 13.7	44.3 ± 37.3	2 ± 2.1	0.15 ± 0.3	0	33 ± 12.8
	CoTr	87 ± 14.2	42.4 ± 36.3	1.9 ± 2.3	0.16 ± 0.31	1	28.2 ± 9.7
	nnU-Net (fine-tuning)	<b>88.4 ± 11.2</b>	<b>25.6 ± 20.2</b>	<b>1.4 ± 1.5</b>	<b>0.13 ± 0.26</b>	0	6.6 ± 4.7

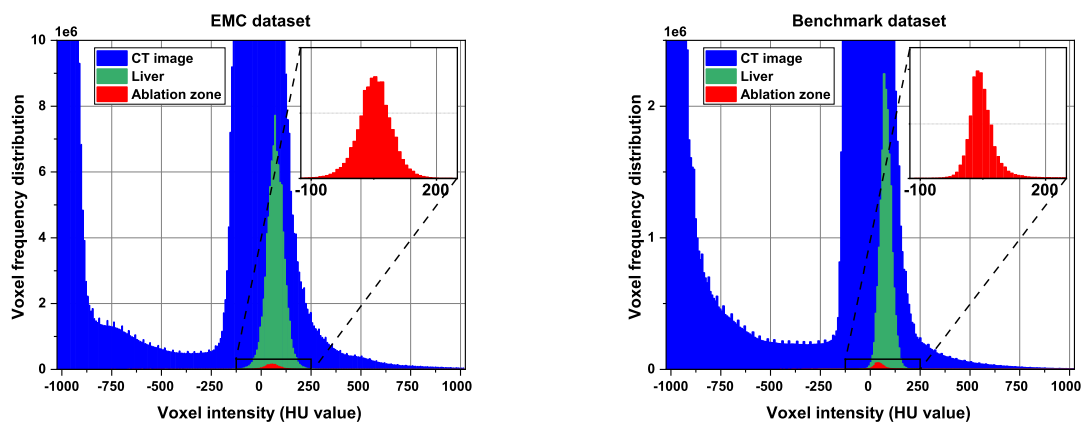


Figure 5: The histogram of voxel intensity in the CT scans in the EMC dataset (left) and the Benchmark dataset (right).

416 seconds for a CT volume, on average. Furthermore, nnU-Net trained by EMC dataset only has the  
 417 longest processing time— an average of approximately 30 seconds for a CT volume. The reason is that the  
 418 nnU-Net framework is designed for patch-based segmentation, which means that the framework needs  
 419 to define the patch size and the patch separation strategy based on the training dataset. The nnU-Net  
 420 (fine-tuning) performs a preprocessing stage based on the EMC<sub>A</sub> dataset, therefore the preprocessing  
 421 strategy is different from that of nnU-Net and CoTr, resulting in a lower number of split patches. Thus,  
 422 the processing time of nnU-Net (fine-tuning) is significantly reduced compared to nnU-Net and CoTr.  
 423 As a result, we choose the nnU-Net (fine-tuning) model as the automatic method in the proposed  
 424 approach since it has shown high accuracy and sufficiently fast processing time.

### 425 III.D.2. Define optimal model for click-based interactive segmentation

426 a. *HU* truncation range assessment:

Table 3: Performance assessment of interactive segmentation model with several the HU truncation range on CT images. NoC@85% and NoC@90% are the average number of required clicks to achieve mean *DSC* scores of 85% and 90%, respectively.

Model	HU truncation ranges	NoC@85	NoC@90
<b>RITM + DiceCE</b> <b>(Baseline)</b>	-100 to 200	<b>3.62</b>	<b>6.01</b>
	-160 to 240	3.64	6.24
	-100 to 400	3.82	6.53
	-1024 to 1024	4.26	6.46

427 In this experiment, we demonstrate the value of the HU truncation range on ablation zone seg-  
 428 mentation. First, we plot the histogram of the voxel intensity of two datasets to show the distribution  
 429 of ablation zone intensity in the CT image in Figure 5. It can be seen that the range of -100 to 200 HU  
 430 contains the most ablation zone voxel intensity (larger than 99%). Furthermore, we evaluate the impact  
 431 of HU truncation on the performance of the interactive segmentation model. Four HU truncation ranges  
 432 are employed for this purpose. Firstly, the HU range of -100 to 200 is utilized by He et al. for automatic  
 433 ablation zone segmentation<sup>14</sup>. Secondly, the HU range of -160 to 240 is widely used in various methods  
 434 participating in the Liver Tumor Segmentation Benchmark (LiTS)<sup>11</sup>. Thirdly, the HU range of -100  
 435 to 400 is often used for liver segmentation<sup>11</sup>. Finally, the HU range of -1024 to 1024 represents the  
 436 entire HU range of a CT image. We use the RITM (baseline) model with DiceCE loss to perform this  
 437 experiment. The models are trained using truncated CT image datasets. We evaluate using 100 2D  
 438 images randomly selected from the validation set. Table 3 indicates that the range of -100 to 200 HU  
 439 achieved the minimum number of clicks required compared to other HU truncation ranges. Hence, we  
 440 used the HU range of -100 to 200 for truncating the CT image in the interactive segmentation model.

441 *b. Weight & kernel size optimization:*

442 In this section, we examined the impact of weight  $\lambda$  and kernel size  $K$  on the mean number of clicks  
 443 required to achieve *DSC* scores of 85% and 90% (referred to as NoC@85% and NoC@90%). To achieve  
 444 this,  $\lambda$  is varied from 0.1 to 0.9 and  $K$  from 10 to 190 pixels, and the results were evaluated using 100  
 445 2D images randomly selected from the validation set. We experimented with the values of  $\lambda$  and  $K$  in  
 446 two strategies: keeping the value fixed and changing adaptively based on the number of clicks. In the  
 447 fixed strategy, the values of  $\lambda$  and  $K$  are fixed for the segmentation revising process. In the adaptive  
 448 strategy, each click point provided by the user increases the  $\lambda$  by 10% and decreases  $K$  by 10%. A 10%  
 449 change per click is substantial enough to alter the parameters meaningfully and avoid instability in the  
 450 segmentation refinement process. The findings in Figure 6 indicate that when  $\lambda$  is small (e.g., 0.1 and  
 451 0.3), the interactive network requires more clicks. When  $\lambda$  exceeds 0.5, there is no difference in the

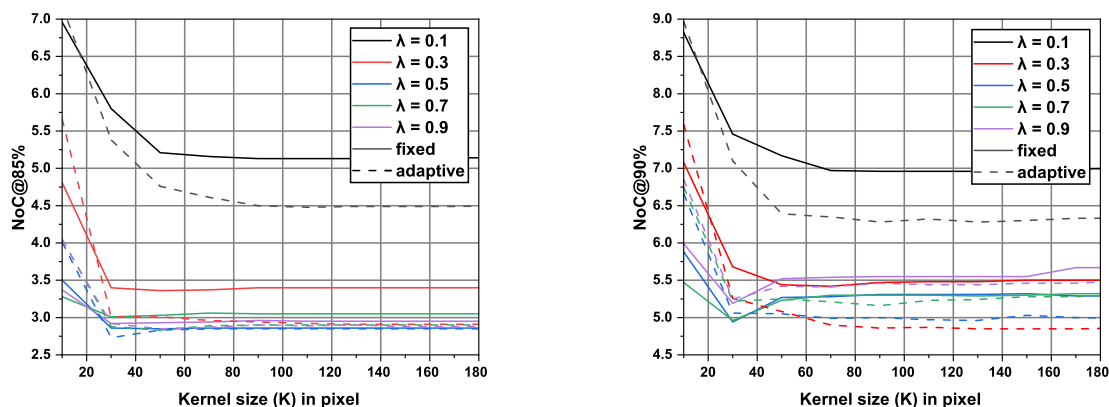


Figure 6: The effect of the mean number of clicks to achieve  $DSC$  of 85% (left) and 90% (right) w.r.t the weighted value ( $\lambda$ ) and kernel size ( $K$ ) in fixed strategy (solid line) and adaptive (dash line) strategy.

452 results. The reason is that when  $\lambda$  is less than 0.5, the weight of interactive segmentation is smaller  
 453 than that of automatic segmentation. Thus, it requires more clicks in the ablation zone regions in which  
 454 the automatic model predicts with low confidence scores. Regarding the kernel size  $K$ , the fewest clicks  
 455 are required when  $K$  is 30 pixels. For  $K$  is less than 30, the ROI may not cover all the large mislabeled  
 456 regions (e.g., Figure 4.B). When  $K$  exceeds 30, it produces an ROI with a large coverage area; when  
 457 the interactive network has mislabeled regions, it subsequently affects the final segmentation. For the  
 458 strategies of value adjustment, we observed that there were no significant differences in the minimum  
 459 value of  $NoC@85\%$  (with a difference of 0.13) and  $NoC@90\%$  (with a difference of 0.08) between the  
 460 fixed and adaptive strategies. Furthermore, utilizing fixed values can offer considerable advantages in  
 461 simplicity, flexibility, and user proactivity during the process of revising segmentation. Based on these  
 462 considerations, we selected fixed values of  $\lambda = 0.5$  and  $K = 30$  for the proposed method.

### 463 c. Interactive methods comparison:

464 This section presents the results of an experiment conducted on the EMC testing set to evaluate  
 465 the impact of the loss function and combination scheme on the mean number of clicks required to  
 466 achieve a mean  $DSC$  score of 85% and 90% (referred to as  $NoC@85\%$  and  $NoC@90\%$ ). Based on the  
 467 study of Sofiiuk et al. (2022), the RITM architecture outperformed several interactive segmentation  
 468 methods<sup>37</sup>. Therefore, we use the RITM network structure to perform this experiment. The experiment  
 469 involved three models: the baseline model, which is an interactive segmentation model without a guided  
 470 mask from automatic segmentation; the automatic initial model, which is an interactive segmentation  
 471 model with a guided mask from automatic segmentation; and the combination scheme model, which is

Table 4: Performance comparison of the loss functions and combination schemes on 2D CT images of the ablation zone. NoC@85% and NoC@90% are the average number of required clicks to achieve mean *DSC* scores of 85% and 90%, respectively. SPC is second per click.

	Model	NoC@85%	NoC@90%	SPC (s)
Baseline	RITM + NFL	5.13	8.03	0.061
	RITM + DiceCE	4.47	7.23	0.055
Automatic Initial	RITM + NFL	4.9	7.86	0.054
	RITM + DiceCE	4.07	6.85	0.051
Combination scheme	RITM + NFL	4.79	6.98	0.057
	RITM + DiceCE (selected)	<b>3.74</b>	<b>6.22</b>	0.054

472 an interactive segmentation model that utilizes the combination scheme described in section II.C.. To  
 473 compare the loss function, we trained the RITM model using Normalized Focal Loss (*NFL*), which was  
 474 used in the original work by Sofiiuk et al. (2022)<sup>37</sup>, and using *DiceCE* loss, a loss function that has  
 475 been used in various medical image segmentation studies. The results, as shown in Table 4, indicate  
 476 that using *DiceCE* loss performs better than using the *NFL* loss in terms of the mean number of clicks  
 477 required. Additionally, the selected combination scheme achieved the highest results. As a result, we  
 478 selected the combination scheme RITM model (RITM + *DiceCE*) for further evaluation.

479 In the next experiment, we compare the performance of the proposed method with the baseline  
 480 interactive segmentation method (RITM) and conventional approach (manual segmentation) on a pilot  
 481 dataset. The pilot dataset contains 10 CT volumes, which are randomly selected from the testing set,  
 482 and contains from one to three ablation zones per volume. Two medical image analysis technicians (3  
 483 years and 1 year of experience), referred to as User 1 and User 2, respectively, utilized the developed  
 484 tool with the two interactive segmentation methods to segment the ablation zone in the pilot dataset.  
 485 Additionally, two users manually annotated the ablation zone slice-by-slice using the Mevislab software.  
 486 Two users perform the ablation zone annotation until the satisfaction is met. The ablation zone appears  
 487 as a non-enhancing area of low attenuation in the CT image. The users segment the ablation zone slice-  
 488 by-slice by delineating the attenuation as a typical procedure<sup>13</sup>. To assess the impact of the interaction  
 489 on the segmentation accuracy, the *DSC* score of the whole 3D CT volume is recalculated when a new  
 490 click is provided by a technician. We also plot the mean *DSC* score of manual segmentation from two  
 491 technicians. The results are shown in Figure 7 (left). It can be seen that, for both of the technicians,  
 492 using the baseline interactive method requires an average of more than 250 clicks to achieve a mean  
 493 *DSC* score of 88%. In contrast, for the proposed method, both of the technicians require averages of 53  
 494 and 94 clicks to achieve saturated mean *DSC* scores of approximately 91.1% and 92.4%, respectively,

495 indicating that the proposed method outperformed RITM (baseline) in terms of  $DSC$  score with the  
 496 same amount of clicks. Furthermore, mean  $DSC$  scores of 92.4% and 90.9% are achieved by User 1 and  
 497 User 2 in manual segmentation, respectively. We also achieved the mean  $DSC$  of 92.4% between the  
 498 ablation zone manually annotated by the two technicians. The experiment demonstrated a high level of  
 499 inter-observation agreement between the manual segmentation annotations made by two technicians.  
 500 The annotation time for each case was recorded. The average annotation time is shown in Figure 7  
 501 (right). The results experimental show that using the proposed method, the average annotation time  
 502 is reduced by approximately 40% and 60% compared with the baseline interactive method and manual  
 503 segmentation, respectively.

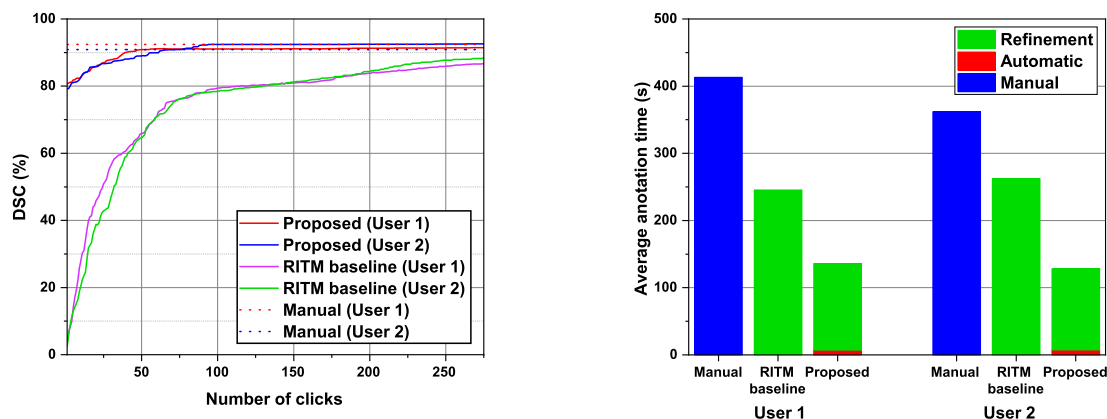


Figure 7: Experiment on the pilot dataset. Mean  $DSC$  scores w.r.t the number of clicks by two users for segmenting the ablation zones in a CT volume (left). Average annotation time by two users in three manners: Manual, RITM baseline assistance, and proposed method assistance (right). RITM baseline is the interactive segmentation model trained using NFL loss.

### 504 III.D.3. Semi-automatic segmentation performance

505 In this section, we investigate the performance of the proposed method on 3D CT images with human  
 506 interaction. A technician uses the interactive segmentation tool, which integrates the proposed method,  
 507 to correct the ablation zone region in the 3D CT images. The evaluation on the  $EMC_B$  dataset, arterial  
 508 phase CT subset shows that the proposed segmentation method obtained a mean  $DSC$ ,  $HD$ ,  $ASSD$ ,  
 509 and  $VD$  of 92.3%, 6.5 mm, 0.5 mm, and 0.05, respectively. These metrics for the portal-venous subset  
 510 are 94%, 8.4 mm, 0.4 mm, and 0.02, respectively (see Table 5). The paired  $t$ -tests to those of the  
 511 nnU-Net (fine-tuning) obtained  $p$ - values which are less than 0.01, suggesting that the proposed method  
 512 statistically significantly improves the segmentation accuracy of the automatic method. In addition, the

Table 5: Performance comparison of ablation zone segmentations between the automatic/semi-automatic segmentation and proposed methods on the EMC<sub>B</sub> dataset. The bold numbers are the highest mean values.

Dataset	Method	DSC	HD (mm)	ASSD (mm)	VD	# fails	Processing time (s)
EMC <sub>B</sub> (Arterial)	nnU-Net (fine-tuning)	80 ± 15.6	29.9 ± 30.3	3.4 ± 5.4	0.3 ± 0.33	2	7.1 ± 3.5
	MONAI Label (Deepedit)	50.4 ± 23	99.9 ± 82.6	16.6 ± 20.1	0.64 ± 0.41	7	2.1 ± 1.1
	Proposed	<b>92.3 ± 3.6</b>	<b>6.5 ± 3.2</b>	<b>0.5 ± 0.3</b>	<b>0.05 ± 0.13</b>	0	121.5 ± 103.1
EMC <sub>B</sub> (Portal venous)	nnU-Net (fine-tuning)	88.4 ± 11.2	25.6 ± 20.2	1.4 ± 1.5	0.13 ± 0.26	0	6.6 ± 4.7
	MONAI Label (Deepedit)	60.3 ± 18.2	159.7 ± 137.2	35.1 ± 55.1	0.14 ± 0.42	2	2.4 ± 1.2
	Proposed	<b>94.0 ± 2.2</b>	<b>8.4 ± 5.9</b>	<b>0.4 ± 0.2</b>	<b>0.02 ± 0.06</b>	0	126.7 ± 105.8

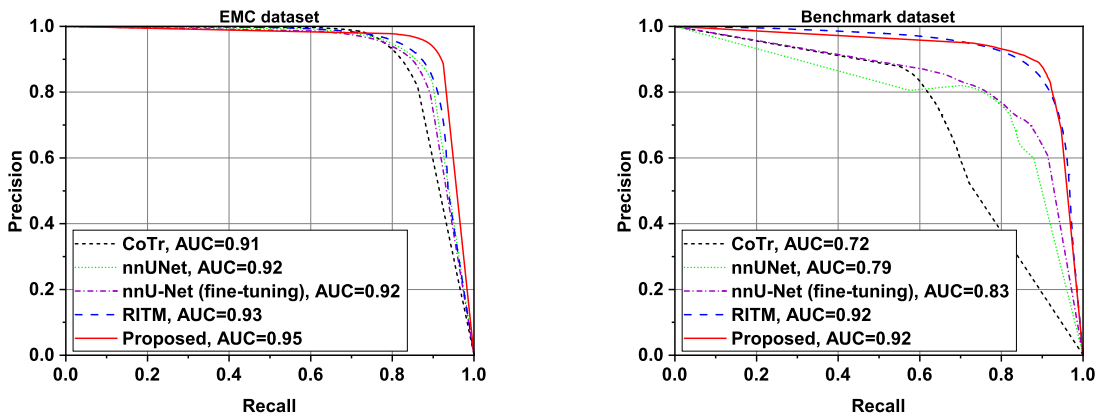


Figure 8: The Precision-Recall curve of ablation zone segmentation on the EMC dataset (left) and the Benchmark dataset (right).

513 proposed method successfully segmented all of the lesions in the EMC<sub>B</sub> dataset while the nnU-Net  
 514 fine-tuning model failed 2 cases in the EMC<sub>B</sub> arterial subset. Moreover, the Precision-Recall curve  
 515 of the proposed method shows highly precise ablation zone segmentation compared to the automatic  
 516 segmentation methods: CoTr, nnU-Net, nnU-Net fine-tuning (as depicted in Figure 8). Specifically, the  
 517 proposed method’s AUC scores are 0.92 and 0.95 for the Benchmark dataset and the EMC<sub>B</sub> dataset,  
 518 respectively, which are greater than those of the other automatic methods. Examples of ablation zone  
 519 segmentation by the proposed method and the other methods on EMC and Benchmark dataset are in  
 520 Figure 10.

521 To further assess the segmentation accuracy of the proposed method on the Benchmark dataset,  
 522 we compared the proposed method with inter-observer manual segmentation and the other well-known  
 523 segmentation methods, including CoTr, nnU-Net, nnU-Net (fine-tuning) and Graph-based contouring<sup>12</sup>.  
 524 Two experts labeled the Benchmark dataset. We use the labels created by the first expert as the



Table 6: Performance comparison of ablation zone segmentations on the Benchmark dataset. The bold numbers are the best scores. The statistics of Graph-based contouring method is listed from the original paper by Egger et al. (2015)<sup>12</sup>.

Method	DSC	HD (mm)	ASSD (mm)	VD	# failure	Processing time (s)
3D U-Net	62.9 ± 21.5	87.9 ± 78.8	7.8 ± 6.9	0.33 ± 0.64	1	4.7 ± 1.5
UNETR	60.8 ± 13.8	260.3 ± 75.3	25.2 ± 7.5	-0.02 ± 0.43	2	4.5 ± 2.2
nnU-Net	76.1 ± 18.4	40.5 ± 41	6.4 ± 8.8	-0.06 ± 0.46	0	29.8 ± 7.9
CoTr	78 ± 17	47 ± 45	7.3 ± 9.5	-0.19 ± 0.28	3	23.1 ± 7.6
nnU-Net (fine-tuning)	77.2 ± 16.6	33.9 ± 40.3	5.8 ± 8.2	-0.07 ± 0.39	0	5.5 ± 4.5
MONAI Label (Deepedit)	52.3 ± 20.8	105.8 ± 100	14.2 ± 13.8	0.42 ± 0.75	2	-
Proposed	<b>87.8 ± 6.8</b>	<b>9.5 ± 6.9</b>	<b>0.9 ± 0.5</b>	<b>-0.03 ± 0.07</b>	<b>0</b>	<b>134.3 ± 82.8</b>
Manual (inter-observer)	88.8 ± 3.3	8.6 ± 3.4	0.8 ± 0.2	-0.02 ± 0.07	0	-

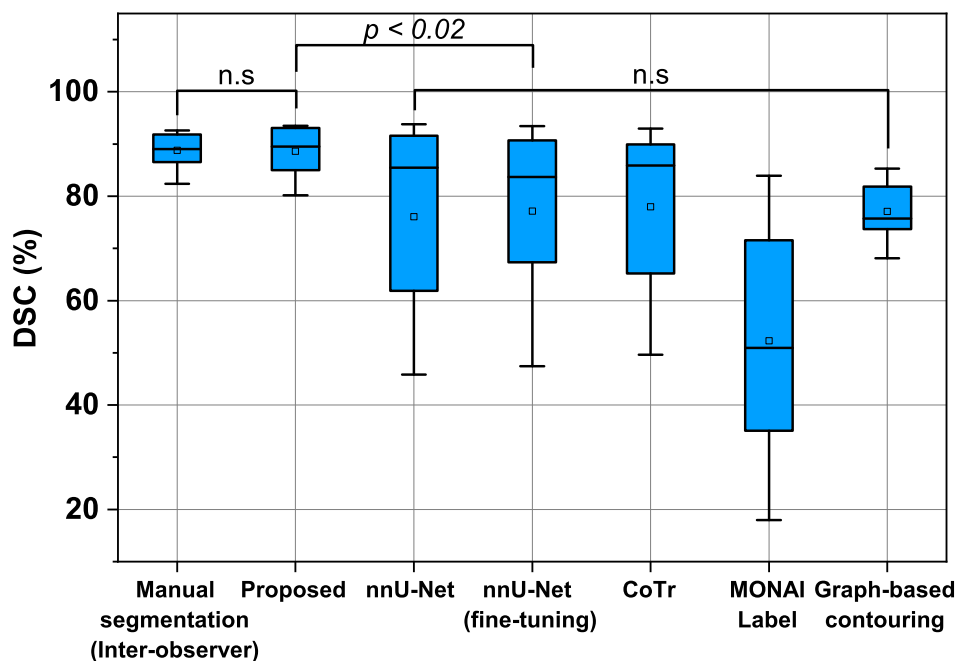


Figure 9: The boxplot of DSC scores among the manual segmentation, the proposed method, the automatic methods, and the classical interactive method<sup>12</sup> for ablation zone segmentation on the Benchmark dataset.

525 ground truth, while the labels created by the second expert are used to represent the inter-observer  
 526 manual segmentation. From Table 6, the mean *DSC*, *HD*, *ASSD*, and *VD* scores achieved by the  
 527 proposed method were 87.8%, 9.5mm, 0.9 mm, and -0.03, respectively. The inter-observer manual  
 528 segmentation achieved mean *DSC*, *HD*, *ASSD*, and *VD* scores of 88.8%, 8.64 mm, 0.8 mm, and -0.02,  
 529 respectively. In addition, we applied a *t*-test on the *DSC* scores of the methods. As shown in Figure 9,  
 530 there is no statistically significant difference between the proposed method and inter-observer manual

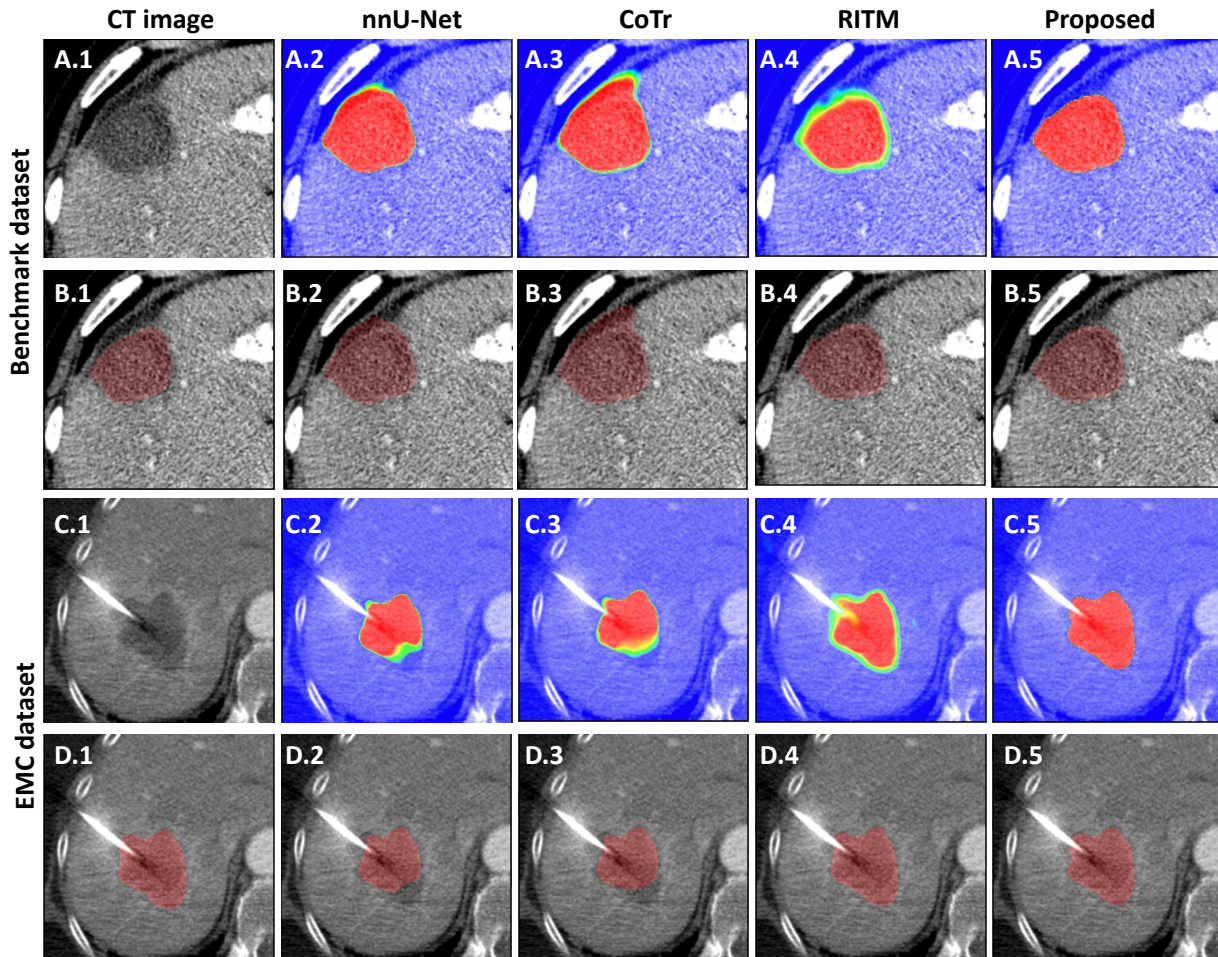


Figure 10: Examples of ablation zone segmentation on EMC and Benchmark dataset of the methods with the segmentation ground truths (B.1 and D.1). The original images (A.1 and C.1) are overlaid by the probability predictions (A.2-5; C2-5) and the thresholded segmentations of each corresponding method (B.2-5; D2-5).

531 segmentation ( $p$ -value = 0.55). While the performances of CoTr, nnU-Net, nnU-Net (fine-tuning),  
 532 and Graph-based contouring are not statistically significantly different, the proposed method obtained  
 533 a statistically significantly better performance compared to those of the methods ( $p$ -values  $\leq 0.02$ ).  
 534 Furthermore, the mean processing time to correct the ablation zone is 134 seconds.

## 535 IV. Discussion

536 In this study, we have proposed and evaluated a semi-automatic method for accurate segmentation of the  
 537 ablation zone in the post-interventional liver tumor ablation CT images. The ablation zone segmentation  
 538 accuracy was compared to five state-of-the-art segmentation methods using both internal and publicly

539 available external datasets. Extensive experiments were carried out to assess the performance of the  
540 proposed method. In addition, we also developed a tool for demonstrating the effectiveness of the  
541 method. The demonstration tool and the source code were made publicly available for research purposes.

542 Table 2 displays the ablation zone segmentation accuracy of four automatic segmentation meth-  
543 ods. The results indicate that nnU-Net performs better than the other baseline methods in automatic  
544 segmentation, showing the effectiveness of the self-configuration framework for automatic ablation zone  
545 segmentation. Furthermore, using the  $t$ -test, we found no statistically significant differences between  
546 the results of nnU-Net and CoTr (with a  $p$ -value ranging from 0.2 to 0.5). The accuracy of both  
547 methods is comparable to those of the state-of-the-art ablation zone segmentation method reported by  
548 He et al. (2021)<sup>14</sup> and Anderson et al. (2022)<sup>15</sup>. We also found that the accuracy of the methods  
549 is reduced when using arterial phase CT images compared to the portal venous phase, which is consis-  
550 tent with the study by He et al. (2021)<sup>14</sup>. However, the performance of the automatic segmentation  
551 methods suggests that they are still unreliable and insufficient for clinical use. In addition, Table 6  
552 shows that the accuracy of the method decreases when performed on the external dataset, indicating  
553 that ablation zone segmentation is still a challenge for fully automatic methods. On the other hand, the  
554 proposed semi-automatic segmentation method achieves state-of-the-art performance on ablation zone  
555 segmentation and outperforms the other methods, yielding mean  $DSC$  scores of 92.3%, 94.0%, and  
556 87.8% on the Arterial  $EMC_B$ , Portal venous  $EMC_B$  and Benchmark datasets, respectively (see Table 5,  
557 6, and 8), which are remarkably better than the mean  $DSC$  scores reported by Anderson et al. (2022)<sup>15</sup>  
558 (79%). The means of  $ASSD$  cores of the proposed method on both the  $EMC$  dataset and Benchmark  
559 dataset are less than 1 mm, which is smaller than the ideal ablation zone safety margin of 10 mm<sup>46</sup>  
560 and equivalent to the median surface distance reported by Anderson et al. (2022)<sup>15</sup> (0.76 mm).

561 From the experiment with MONAI Label, we see that MONAI Label, a state-of-the-art semi-  
562 automatic segmentation method for medical images<sup>47</sup>, yielded low accuracy with the embedded click-  
563 based method. This is because the default preprocessing setting of the MONAI Label is investigated for  
564 multiple organ segmentation (BTCV Challenge), which may not be optimal for a single class (ablation  
565 zone segmentation only). In the resampling step, a fixed spacing is applied for the entire data, and a  
566 large spacing (spacing of [1.5, 1.5]) for the axial plane resampling makes lost information. Additionally,  
567 the patch-based segmentation with a small patch size (the size of [96, 96, 32]) results in a class  
568 imbalance during the training phase. Since the ablation zone region is small compared to the whole  
569 CT volume, the number of patches that contain the background only is larger than the number of  
570 patches that contain the ablation zone. This is the evidence to explain the performance of models

---

571 implemented using the MONAI library (3D U-Net, UNETR, and MONAI Label) achieved low accuracy.  
572 In contrast, the proposed method has the advantage of the self-configuration framework (nnU-Net),  
573 which can automatically adapt the preprocessing step based on the training data. Furthermore, using  
574 a 2D interactive segmentation model (RITM) for the refinement of the ablation zone in a slice-based  
575 manner supports the technician in the refinement of the ablation zone as a typical procedure.

576 Figure 9 also indicates that the proposed method achieved the best performance compared to  
577 other automatic methods. Our observation is that the  $EMC_B$  dataset is less noisy than the Benchmark  
578 dataset. Thus, the performance of the methods on the  $EMC_B$  dataset is higher than that on the  
579 Benchmark dataset. In addition, the manual segmentations of the Benchmark dataset were created  
580 by two experts. By quantitatively evaluating the inter-observer variability of manual segmentation, we  
581 found that the proposed method achieved a segmentation accuracy comparable with the inter-observer  
582 variation in terms of  $DSC$ . The  $p$ -value of 0.55 ( $t$ -test) suggests that there is no statistically significant  
583 difference between the proposed method and the inter-observer manual technique (see Figure 9).

584 From Figure 7, we can see that the proposed method needs a lower number of required clicks  
585 compared to the original RITM. This is because the proposed method takes advantage of the automatic  
586 segmentation to reduce the workload for the user. This indicates that the performance of the automatic  
587 method is also an essential factor in reducing the number of clicks. The larger the number of clicks is,  
588 the more time and inconvenient it is for the operator to obtain a good segmentation. From Table 5  
589 and 6, we can see that the average processing time of the proposed method for a volume is around 2  
590 minutes, which is small compared to the average operation time of a MWA/RFA session of 112 -149  
591 minutes<sup>48</sup>. Moreover, it can be seen from Figure 7 that the mean saturated  $DSC$  scores of both users  
592 are slightly different. Since the proposed method requires human interaction, we suppose the difference  
593 is caused by inter-observer variation.

594 Our study has some limitations. Firstly, although we conducted the study with both internal and  
595 external datasets, and achieved state-of-the-art performance, the number of CT images in the external  
596 dataset is only 12 CT volumes, thus conclusions on generality should be drawn with care. By sharing  
597 our source code and demonstration tool, we expect other researchers can easily reproduce our obtained  
598 results and perform testing on larger external datasets. Secondly, the developed demonstration tool  
599 was derived from the work by Sofiiuk et al.(2022)<sup>37</sup>, which was not originally designed for medical  
600 application purposes. As a result, the number of interactions might not be optimal yet. In this study,  
601 we consider the tool for the demonstration purpose only. Further studies may require a better design for

602 the user interface and user experience. Another solution which could be considered is to integrate the  
603 proposed method with existing medical image tools such as MONAI<sup>43</sup>, ITK-SNAP<sup>49</sup> and 3D Slicer<sup>50</sup>.  
604 In addition, although the accuracy of the semi-automatic method is comparable with the accuracy of the  
605 manual segmentation by experts, it still contains errors in the final segmentation with a mean *HD* score  
606 of 9.5 mm (compared that of the 8.6 mm inter-observer score,  $p$ -value = 0.67). These errors seem to  
607 be the limitations of human-level performance for annotation and evaluation on the Benchmark dataset.  
608 We suggest that the interventionist may take the errors into account in assessing the ablation zone. We  
609 acknowledge that the time cost of the proposed method, which is about 120 seconds, is higher than  
610 automatic segmentation methods that can produce results in a matter of seconds. This difference in  
611 time cost is indeed a consideration and can be seen as a limitation when comparing our method to fully  
612 automated approaches. However, the proposed method allows for higher accuracy and customization, as  
613 users can iteratively refine the segmentation. This is particularly advantageous in cases where automatic  
614 methods might struggle with complex or ambiguous regions. Finally, in this study, we mainly focus on  
615 developing methods for precise ablation zone segmentation without further investigating the effect of  
616 the ablation zone segmentation on the clinical outcome. Nevertheless, Lin et al. (2023)<sup>46</sup> suggested in  
617 their recent study that precise ablation zone segmentation has clinical benefits.

618 The use of deep learning for medical image analysis is massively expanding at present, especially  
619 for image segmentation applications<sup>51,52</sup>. A major drawback of deep learning is that it requires a  
620 sufficiently large amount of data for effective training of the models. However, it is frequently difficult  
621 to acquire a sufficient amount of medical images with labels that are appropriate for a specific application,  
622 potentially resulting in sub-optimal performance. For fully automatic CNN-based segmentation methods,  
623 the predicted segmentation may therefore contain segmentation errors. However, with simple interactive  
624 corrections using the proposed semi-automatic CNN-based method, the accuracy of segmentation can  
625 be improved significantly. Therefore, we expect that using the proposed approach, other segmentation  
626 problems may be similarly addressed without requiring large amounts of training data.

## 627 V. Conclusions

628 This study has proposed a semi-automatic approach for ablation zone segmentation in thermal treat-  
629 ments of liver cancer. An accurate segmentation is obtained by combining automatic CNN-based seg-  
630 mentation and click-based CNN segmentation methods. Regarding segmentation accuracy, the proposed  
631 method is superior to the well-known CNNs in almost all metrics, achieving comparable performance to

---

632 manual segmentation of human experts on a benchmark dataset, yielding a mean *DSC* score of 87.8%  
633 on average. The obtained segmentation accuracy scores of the proposed approach are also better than  
634 those of the other methods when applied to the internal dataset, achieving state-of-the-art performance  
635 in accuracy (*DSC* score of 94.0% on average), and the method is sufficiently fast for the use in clinical  
636 practice. In conclusion, this study has shown the potential of the semi-automatic approach in supporting  
637 the interventionist in assessing the treatment outcome of thermal ablation for liver cancer treatment.

---

## 638 Ethical statement

639 The local medical research ethics committee decided that the Medical Research Involving Human Sub-  
640 jects Act does not apply to this study. The Benchmark dataset is publicly available for research purpose.

## 641 Conflicts of interest

642 None.

## 643 Acknowledgments

644 This research was funded by Vietnam National Foundation for Science and Technology Development  
645 (NAFOSTED) under grant number 102.01-2018.316. We would like to thank NVIDIA for supporting  
646 the RTX 8000 GPU for this study. We would like to thank Mr. Yaro Roodenburg from TU Delft for  
647 the early experiment on ablation zone segmentation using nnU-Net. We also would like to thank Mr.  
648 Quang Van Nguyen for evaluating the developed tool.

## 650 References

- 649
- 651 <sup>1</sup> M. J. Rutherford et al., Comparison of liver cancer incidence and survival by subtypes across seven  
652 high-income countries, *International Journal of Cancer* **149**, 2020–2031 (2021).
  - 653 <sup>2</sup> H. Rumgay, M. Arnold, J. Ferlay, O. Lesi, C. J. Cabasag, J. Vignat, M. Laversanne, K. A. McGlynn,  
654 and I. Soerjomataram, Global burden of primary liver cancer in 2020 and predictions to 2040, *Journal*  
655 *of Hepatology* **77**, 1598–1606 (2022).
  - 656 <sup>3</sup> B. Cadier et al., Early detection and curative treatment of hepatocellular carcinoma: a cost-  
657 effectiveness analysis in France and in the United States, *Hepatology* **65**, 1237–1248 (2017).
  - 658 <sup>4</sup> J. Han, Y.-c. Fan, and K. Wang, Radiofrequency ablation versus microwave ablation for early stage  
659 hepatocellular carcinoma: A PRISMA-compliant systematic review and meta-analysis, *Medicine* **99**  
660 (2020).
  - 661 <sup>5</sup> D. Li, J. Kang, B. J. Golas, V. W. Yeung, and D. C. Madoff, Minimally invasive local therapies for  
662 liver cancer, *Cancer biology & medicine* **11**, 217 (2014).
-

- 663 <sup>6</sup> M. B. Glassberg, S. Ghosh, J. W. Clymer, R. A. Qadeer, N. C. Ferko, B. Sadeghirad, G. W. Wright,  
664 and J. F. Amaral, Microwave ablation compared with radiofrequency ablation for treatment of  
665 hepatocellular carcinoma and liver metastases: a systematic review and meta-analysis, *OncoTargets*  
666 and therapy **12**, 6407 (2019).
- 667 <sup>7</sup> Z. Yu, G. Li, N. Yuan, and W. Ding, Comparison of ultrasound guided versus CT guided radiofre-  
668 quency ablation on liver function, serum PIVKA-II, AFP levels and recurrence in patients with  
669 primary hepatocellular carcinoma, *American Journal of Translational Research* **13**, 6881 (2021).
- 670 <sup>8</sup> J. H. Kim, H. J. Won, Y. M. Shin, P. N. Kim, S.-G. Lee, and S. Hwang, Radiofrequency ablation for  
671 recurrent intrahepatic cholangiocarcinoma after curative resection, *European journal of radiology*  
672 **80**, e221–e225 (2011).
- 673 <sup>9</sup> G. Antoch, F. M. Vogt, P. Veit, L. S. Freudenberg, N. Blehschmid, O. Dirsch, A. Bockisch,  
674 M. Forsting, J. F. Debatin, and H. Kuehl, Assessment of liver tissue after radiofrequency ablation:  
675 findings with different imaging procedures, *Journal of Nuclear Medicine* **46**, 520–525 (2005).
- 676 <sup>10</sup> S. Yedururi, S. Terpenning, S. Gupta, P. Fox, S. Martin, C. Conrad, and E. M. Loyer, Radio-  
677 frequency Ablation of Hepatic Tumor: Subjective Assessment of the Perilesional Vascular Network  
678 on Contrast Enhanced CT Before and After Ablation can Reliably Predict the Risk of Local Recur-  
679 rence, *Journal of computer assisted tomography* **41**, 607 (2017).
- 680 <sup>11</sup> P. Bilic et al., The liver tumor segmentation benchmark (lits), *Medical Image Analysis* **84**, 102680  
681 (2023).
- 682 <sup>12</sup> J. Egger et al., Interactive volumetry of liver ablation zones, *Scientific Reports* **5**, 15373 (2015).
- 683 <sup>13</sup> P.-h. Wu, M. Bedoya, J. White, and C. L. Brace, Feature-based automated segmentation of  
684 ablation zones by fuzzy c-mean clustering during low-dose computed tomography, *Medical physics*  
685 **48**, 703–714 (2021).
- 686 <sup>14</sup> K. He, X. Liu, R. Shahzad, R. Reimer, F. Thiele, J. Niehoff, C. Wybranski, A. C. Bunck, H. Zhang,  
687 and M. Perkuhn, Advanced deep learning approach to automatically segment malignant tumors and  
688 ablation zone in the liver with contrast-enhanced CT, *Frontiers in Oncology* **11**, 669437 (2021).
- 689 <sup>15</sup> B. M. Anderson, B. Rigaud, Y.-M. Lin, A. K. Jones, H. C. Kang, B. C. Odisio, and K. K. Brock,  
690 Automated segmentation of colorectal liver metastasis and liver ablation on contrast-enhanced CT  
691 images, *Frontiers in Oncology* **12** (2022).
-



- 692 <sup>16</sup> O. Ronneberger, P. Fischer, and T. Brox, U-net: Convolutional networks for biomedical image  
693 segmentation, in *Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015:*  
694 *18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III 18,*  
695 pages 234–241, Springer, 2015.
- 696 <sup>17</sup> W. Xu, H. Liu, X. Wang, and Y. Qian, Liver segmentation in CT based on ResUNet with 3D  
697 probabilistic and geometric post process, in *2019 IEEE 4th International Conference on Signal and*  
698 *Image Processing (ICSIP)*, pages 685–689, IEEE, 2019.
- 699 <sup>18</sup> X. Li, H. Chen, X. Qi, Q. Dou, C.-W. Fu, and P.-A. Heng, H-DenseUNet: hybrid densely connected  
700 UNet for liver and tumor segmentation from CT volumes, *IEEE transactions on medical imaging*  
701 **37**, 2663–2674 (2018).
- 702 <sup>19</sup> F. Isensee, P. F. Jaeger, S. A. Kohl, J. Petersen, and K. H. Maier-Hein, nnU-Net: a self-configuring  
703 method for deep learning-based biomedical image segmentation, *Nature methods* **18**, 203–211  
704 (2021).
- 705 <sup>20</sup> A. Dosovitskiy et al., An image is worth 16x16 words: Transformers for image recognition at scale,  
706 arXiv preprint arXiv:2010.11929 (2020).
- 707 <sup>21</sup> A. Hatamizadeh, Y. Tang, V. Nath, D. Yang, A. Myronenko, B. Landman, H. R. Roth, and D. Xu,  
708 Unetr: Transformers for 3d medical image segmentation, in *Proceedings of the IEEE/CVF winter*  
709 *conference on applications of computer vision*, pages 574–584, 2022.
- 710 <sup>22</sup> J. Chen, Y. Lu, Q. Yu, X. Luo, E. Adeli, Y. Wang, L. Lu, A. L. Yuille, and Y. Zhou, Tran-  
711 sunet: Transformers make strong encoders for medical image segmentation, arXiv preprint  
712 arXiv:2102.04306 (2021).
- 713 <sup>23</sup> Y. Xie, J. Zhang, C. Shen, and Y. Xia, Cotr: Efficiently bridging cnn and transformer for 3d  
714 medical image segmentation, in *Medical Image Computing and Computer Assisted Intervention–*  
715 *MICCAI 2021: 24th International Conference, Strasbourg, France, September 27–October 1, 2021,*  
716 *Proceedings, Part III 24*, pages 171–180, Springer, 2021.
- 717 <sup>24</sup> J. Wallner, M. Schwaiger, K. Hocegger, C. Gsaxner, W. Zemann, and J. Egger, A review on  
718 multiplatform evaluations of semi-automatic open-source based image segmentation for cranio-  
719 maxillofacial surgery, *Computer methods and programs in biomedicine* **182**, 105102 (2019).
-

- 720 <sup>25</sup> R. M. Haralick and L. G. Shapiro, Image segmentation techniques, *Computer vision, graphics, and*  
721 *image processing* **29**, 100–132 (1985).
- 722 <sup>26</sup> T. F. Chan and L. A. Vese, Active contours without edges, *IEEE Transactions on image processing*  
723 **10**, 266–277 (2001).
- 724 <sup>27</sup> C. Rother, V. Kolmogorov, and A. Blake, "GrabCut" interactive foreground extraction using  
725 iterated graph cuts, *ACM transactions on graphics (TOG)* **23**, 309–314 (2004).
- 726 <sup>28</sup> Y. Gao, R. Kikinis, S. Bouix, M. Shenton, and A. Tannenbaum, A 3D interactive multi-object  
727 segmentation tool using local robust statistics driven active contours, *Medical image analysis* **16**,  
728 1216–1227 (2012).
- 729 <sup>29</sup> N. Otsu, A threshold selection method from gray-level histograms, *IEEE transactions on systems,*  
730 *man, and cybernetics* **9**, 62–66 (1979).
- 731 <sup>30</sup> V. Caselles, R. Kimmel, and G. Sapiro, Geodesic active contours, *International journal of computer*  
732 *vision* **22**, 61 (1997).
- 733 <sup>31</sup> M. Rajchl et al., Deepcut: Object segmentation from bounding box annotations using convolutional  
734 neural networks, *IEEE transactions on medical imaging* **36**, 674–683 (2016).
- 735 <sup>32</sup> D. Lin, J. Dai, J. Jia, K. He, and J. Sun, Scribblesup: Scribble-supervised convolutional networks  
736 for semantic segmentation, in *Proceedings of the IEEE conference on computer vision and pattern*  
737 *recognition*, pages 3159–3167, 2016.
- 738 <sup>33</sup> G. Wang et al., DeepGeoS: a deep interactive geodesic framework for medical image segmentation,  
739 *IEEE transactions on pattern analysis and machine intelligence* **41**, 1559–1572 (2018).
- 740 <sup>34</sup> G. Wang et al., Interactive medical image segmentation using deep learning with image-specific  
741 fine tuning, *IEEE transactions on medical imaging* **37**, 1562–1573 (2018).
- 742 <sup>35</sup> X. Luo, G. Wang, T. Song, J. Zhang, M. Aertsen, J. Deprest, S. Ourselin, T. Vercauteren, and  
743 S. Zhang, MIDeepSeg: Minimally interactive segmentation of unseen objects from medical images  
744 using deep learning, *Medical image analysis* **72**, 102102 (2021).
- 745 <sup>36</sup> L. Sun, Z. Tian, Z. Chen, W. Luo, and S. Du, An efficient interactive segmentation framework for  
746 medical images without pre-training, *Medical Physics* (2022).
-

- 747 <sup>37</sup> K. Sofiiuk, I. A. Petrov, and A. Konushin, Reviving iterative training with mask guidance for  
748 interactive segmentation, in *2022 IEEE International Conference on Image Processing (ICIP)*,  
749 pages 3141–3145, IEEE, 2022.
- 750 <sup>38</sup> Ö. Çiçek, A. Abdulkadir, S. S. Lienkamp, T. Brox, and O. Ronneberger, 3D U-Net: learning dense  
751 volumetric segmentation from sparse annotation, in *Medical Image Computing and Computer-  
752 Assisted Intervention–MICCAI 2016: 19th International Conference, Athens, Greece, October 17-  
753 21, 2016, Proceedings, Part II 19*, pages 424–432, Springer, 2016.
- 754 <sup>39</sup> H. Shan, A. Padole, F. Homayounieh, U. Kruger, R. D. Khera, C. Nitiwarangkul, M. K. Kalra, and  
755 G. Wang, Competitive performance of a modularized deep neural network compared to commercial  
756 algorithms for low-dose CT image reconstruction, *Nature Machine Intelligence* **1**, 269–276 (2019).
- 757 <sup>40</sup> R. Benenson, S. Popov, and V. Ferrari, Large-scale interactive object segmentation with human  
758 annotators, in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*,  
759 pages 11700–11709, 2019.
- 760 <sup>41</sup> H. M. Luu, W. Niessen, T. Van Walsum, C. Klink, and A. Moelker, An automatic registration  
761 method for pre-and post-interventional CT images for assessing treatment success in liver RFA  
762 treatment, *Medical Physics* **42**, 5559–5567 (2015).
- 763 <sup>42</sup> F. Isensee et al., nnu-net: Self-adapting framework for u-net-based medical image segmentation,  
764 arXiv preprint arXiv:1809.10486 (2018).
- 765 <sup>43</sup> M. J. Cardoso et al., MONAI: An open-source framework for deep learning in healthcare, arXiv  
766 preprint arXiv:2211.02701 (2022).
- 767 <sup>44</sup> T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick,  
768 Microsoft coco: Common objects in context, in *Computer Vision–ECCV 2014: 13th European  
769 Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pages 740–755,  
770 Springer, 2014.
- 771 <sup>45</sup> A. Gupta, P. Dollar, and R. Girshick, Lvis: A dataset for large vocabulary instance segmentation,  
772 in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages  
773 5356–5364, 2019.
-

- 774 <sup>46</sup> Y.-M. Lin, I. Paolucci, C. S. O'Connor, B. M. Anderson, B. Rigaud, B. M. Fellman, K. A. Jones,  
775 K. K. Brock, and B. C. Odisio, Ablative margins of colorectal liver metastases using deformable  
776 CT image registration and autosegmentation, *Radiology* , 221373 (2023).
- 777 <sup>47</sup> A. Diaz-Pinto et al., Monai label: A framework for ai-assisted interactive labeling of 3d medical  
778 images, arXiv preprint arXiv:2203.12362 (2022).
- 779 <sup>48</sup> F. Izzo, V. Granata, R. Grassi, R. Fusco, R. Palaia, P. Delrio, G. Carrafiello, D. Azoulay, A. Petrillo,  
780 and S. A. Curley, Radiofrequency ablation and microwave ablation in liver tumors: an update, *The*  
781 *oncologist* **24**, e990–e1005 (2019).
- 782 <sup>49</sup> P. A. Yushkevich, Y. Gao, and G. Gerig, ITK-SNAP: An interactive tool for semi-automatic seg-  
783 mentation of multi-modality biomedical images, in *2016 38th annual international conference of*  
784 *the IEEE engineering in medicine and biology society (EMBC)*, pages 3342–3345, IEEE, 2016.
- 785 <sup>50</sup> R. Kikinis, S. D. Pieper, and K. G. Vosburgh, 3D Slicer: a platform for subject-specific image  
786 analysis, visualization, and clinical support, in *Intraoperative imaging and image-guided therapy*,  
787 pages 277–289, Springer, 2013.
- 788 <sup>51</sup> G. Litjens, T. Kooi, B. E. Bejnordi, A. A. A. Setio, F. Ciompi, M. Ghafoorian, J. A. Van Der Laak,  
789 B. Van Ginneken, and C. I. Sánchez, A survey on deep learning in medical image analysis, *Medical*  
790 *image analysis* **42**, 60–88 (2017).
- 791 <sup>52</sup> X. Chen, X. Wang, K. Zhang, K.-M. Fung, T. C. Thai, K. Moore, R. S. Mannel, H. Liu, B. Zheng,  
792 and Y. Qiu, Recent advances and clinical applications of deep learning in medical image analysis,  
793 *Medical Image Analysis* , 102444 (2022).

## 794 Notes

796 † Deceased 31 July 2023

797 <sup>a</sup><https://github.com/MIC-DKFZ/nnUNet>

798 <sup>b</sup><https://github.com/Project-MONAI>

799 <sup>c</sup><https://github.com/YtongXie/CoTr>

800 <sup>d</sup> [https://github.com/SamsungLabs/ritm\\_interactive\\_segmentation](https://github.com/SamsungLabs/ritm_interactive_segmentation)

801 <sup>e</sup>[https://github.com/lqanh11/Interactive\\_AblationZone\\_Segmentation](https://github.com/lqanh11/Interactive_AblationZone_Segmentation)

---