

Grayscale Medical Image Segmentation Method Based on 2D&3D Object Detection with Deep Learning

Yunfei Ge

Tongji University

Qing Zhang

Tongji University

Yuantao Sun (✉ sun1979@sina.com)

Tongji University

Yidong Shen

The First people's Hospital of Yancheng

Xijiong Wang

Shanghai Bojin Electric Instrument & Device Co., Ltd

Research Article

Keywords: Grayscale medical image, Image segmentation, Deep learning, Object detection, Point cloud

Posted Date: October 29th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-1018292/v1>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Version of Record: A version of this preprint was published at BMC Medical Imaging on February 27th, 2022. See the published version at <https://doi.org/10.1186/s12880-022-00760-2>.

Grayscale Medical Image Segmentation Method Based on 2D&3D Object Detection with Deep Learning

Yunfei Ge¹, Qing Zhang¹, Yuantao Sun^{1*}, Yidong Shen², Xijiong Wang³

¹ School of Mechanical Engineering, Tongji University, Shanghai, China

² Department of Orthopaedics, The First people's Hospital of Yancheng, Yancheng, China

³ Shanghai Bojin Electric Instrument & Device Co., Ltd, Shanghai, China

Abstract

Background: Grayscale medical image segmentation is the key step in clinical computer-aided diagnosis. Model-driven and data-driven image segmentation methods are widely used for their less computational complexity and more accurate feature extraction. However, model-driven methods like thresholding usually suffer from wrong segmentation and noises regions because different grayscale images have distinct intensity distribution property thus pre-processing is always demanded. While data-driven methods with deep learning like encoder-decoder networks always are always accompanied by complex architectures which require amounts of training data.

Methods: Combining thresholding method and deep learning, this paper presents a novel method by using 2D&3D object detection technologies. First, interest regions contain segmented object are determined with fine-tuning 2D object detection network. Then, pixels in cropped images are turned as point cloud according to their positions and grayscale values. Finally, 3D object detection network is applied to obtain bounding boxes with target points and boxes' bottoms and tops represent thresholding values for segmentation. After projecting to 2D images, these target points could composite the segmented object.

Results: Three groups of grayscale medical images are used to evaluate the proposed image segmentation method. We obtain the IoU (DSC) scores of 0.92 (0.96), 0.88 (0.94) and 0.94 (0.94) for segmentation accuracy on different datasets respectively. Also, compared with five state of the arts and clinically performed well models, our method achieves higher scores and better performance.

Conclusions: The prominent segmentation results demonstrate that the built method based on 2D&3D object detection with deep learning is workable and promising for segmentation task of grayscale medical images.

Keywords: Grayscale medical image, Image segmentation, Deep learning, Object detection, Point cloud.

*Corresponding Author Yuantao Sun, E-mail: sun1979@sina.com

1 Background

Medical imaging plays the key role in diagnosis or disease treatment by revealing internal structures with technologies mainly of computer tomography (CT), magnetic resonance imaging (MRI), ultrasound, and especially X ray radiography [1]. Due to different absorption capability of various organs or tissues for radiations, waves, and etc., pixels belong to various object in grayscale medical images have diverse grayscale values usually from 0-255 [2] and meanwhile values of pixels of the same object always gather within a range.

38 Medical image segmentation has been widely applied to make images clearer with anatomical or
39 pathological structures changes [3], such as bone segmentation [4], lung segmentation [5,6], heart
40 fat segmentation [7] and liver or liver-tumor segmentation [8,9], etc. They could be considered to
41 divide origin images into several sub regions for picking up some crucial objects and extracting
42 interesting features which improve the computer aided diagnostic efficiency. There has raised
43 enormous approaches and they could be classified into two categories: model-driven techniques
44 and data-driven techniques. [5,10]

45 Many model-driven methods for medical image segmentation, including thresholding, clustering,
46 and region growing, were presented in particular before the widespread application of deep
47 learning. [10] Thresholding was one of the most common used method in practice due to its
48 efficiency. [11] The basic working of thresholding was to determine specific threshold values and
49 each pixel in the image could be classified as the foreground or background depending on the
50 comparison between their intensity values and threshold values. [12-14] Traditional thresholding
51 methods always relied on single models for universal segmentation tasks which could lead to
52 incorrect results. Also, segmentation objects often occupied only parts of whole images and pixels
53 of different objects may share same intensity values, so noises could appear if image segmentation
54 was applied overall.

55 With the era of big data coming, emerging data-driven technologies with deep learning have
56 remarkably demonstrated in variety medical image segmentation task. Supervised learning
57 methods and especially some CNN (Convolutional neural network) based encoder-decoder
58 structures such as FCN (Fully convolutional networks) [15], U-Net [16], DeepLab [17] has
59 practically proved [5]. Compared with traditional methods, deep learning could help analyze
60 medical images more effectively and extract more detailed features.

61 Although these end-to-end structures was pragmatic for medical images semantic segmentation,
62 the segmentation accuracy always relied on a large amount of training dataset. But medical image
63 annotation could be time-consuming and quite expensive, thus transfer learning was used to solve
64 the problem of limited labeled data and pre-trained networks on natural images as ImageNet [18]
65 were often adopted for image segmentation. [19,20] However, considering these datasets were
66 mainly designed to train models for object detection or classification, they may be more suitable
67 to pretrain networks for object detection. This inspired us to segment images with object detection.
68 We find that grayscale images could be segmented according to the comparison of thresholding
69 values with values of pixels in images and these pixels could be turned into 3D point cloud
70 according to their positions and grayscale values. Thus, by applying 3D object detection in the
71 point cloud, we could achieve groups of points within 3D bounding boxes. The top and the bottom
72 of boxes represent the thresholding values for segmentation and after mapping these points into
73 2D images, corresponding pixels could compose segmented results. Besides, 2D object detection
74 could determinate regions of interest (ROI) in grayscale medical images to reduce noises.
75 Therefore, according to above strategy, we propose the grayscale medical image segmentation
76 method based on 2D&3D object detection.

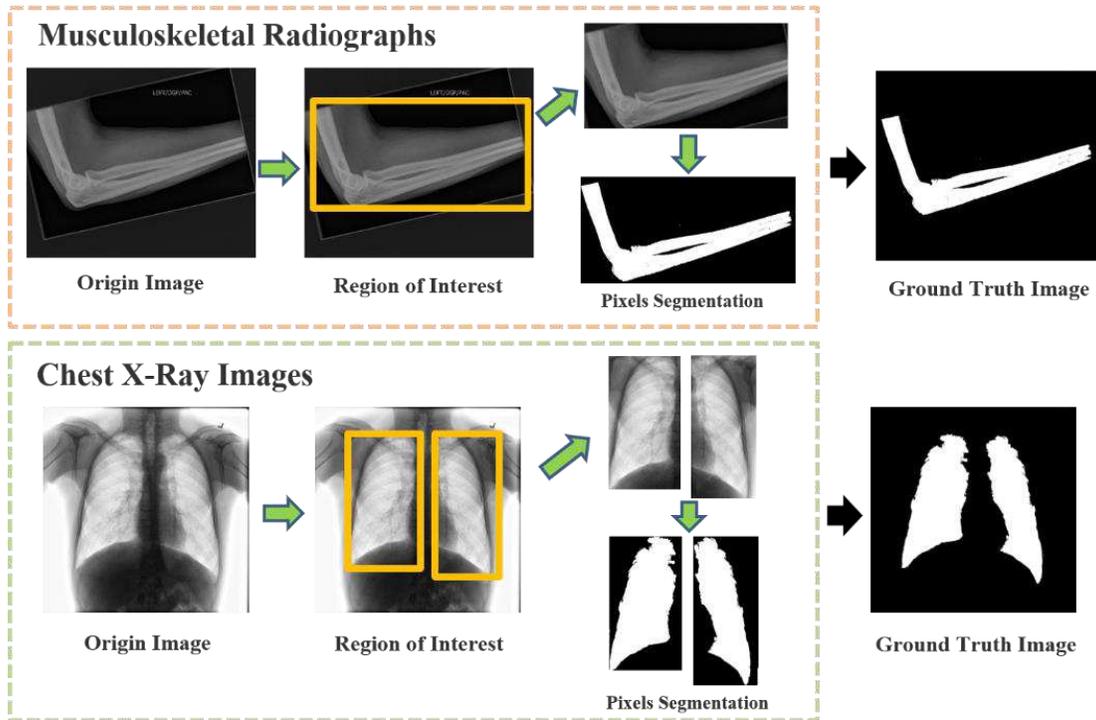
77 The remainder of this paper is organized as following: Section 2 introduces the applied medical
78 Image datasets and describes details of proposed technologies, while in section 3 the obtained
79 results are displayed and the discussion is provided. Finally, section 4 presents the conclusions as
80 well as future work suggestions.

81 **2 Methods**

82 *2.1 Image Datasets*

83 Considering roles of grayscale images in medical field, two typical sets of available datasets are
84 prepared including musculoskeletal radiographs, and chest radiographs. Musculoskeletal
85 radiographs dataset (MURA (Musculoskeletal radiographs) & LERA (Lower extremity
86 radiographs)) contains bone X ray images of upper and lower extremity. [21,22] Chest radiographs
87 dataset CheXpert (Chest radiography) has chest X ray images. [23,24]

88 The proposed grayscale medical image segmentation method is based on the supervised artificial
89 intelligence techniques, and labels are performed manually in two types medical images for model
90 training. Fig.1. shows origin images, and their respective Ground Truth (GT) images in different
91 datasets.



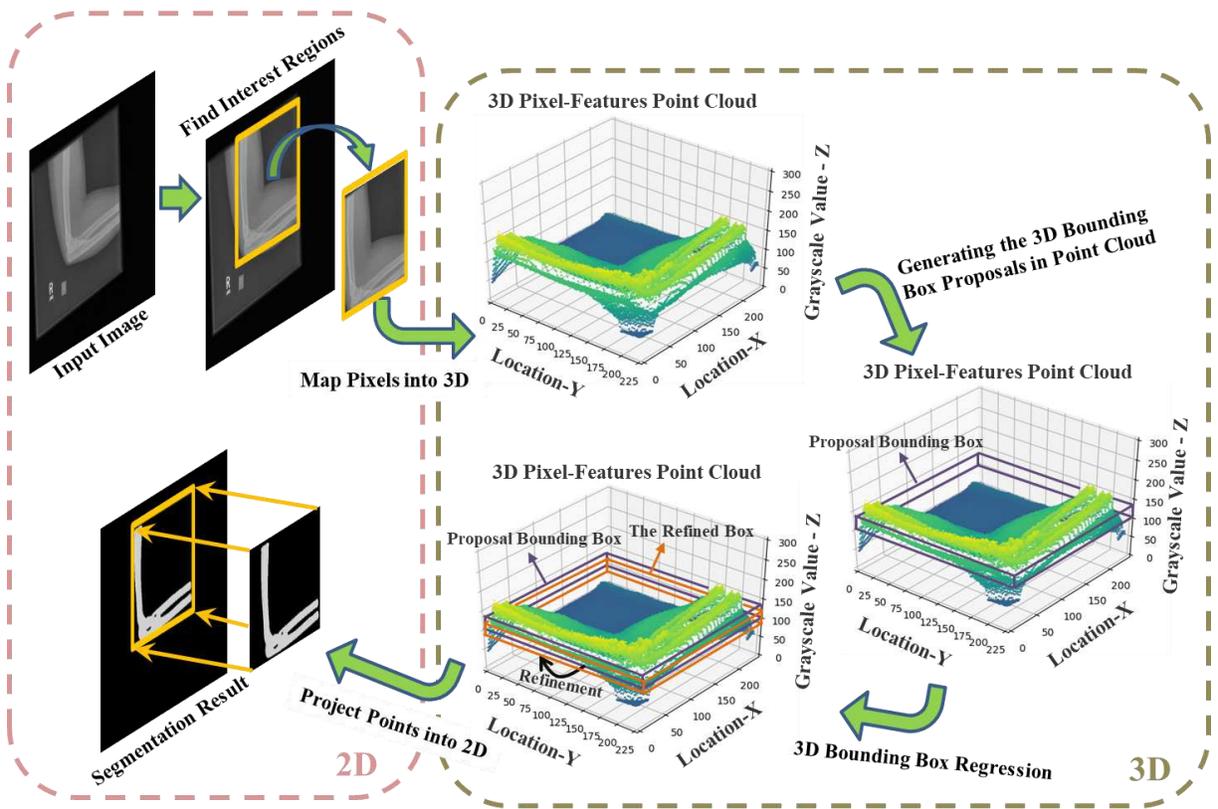
92

93

Fig.1 Examples of medical images in two datasets and manual segmentation results.

94 2.2 Grayscale Image Segmentation Framework

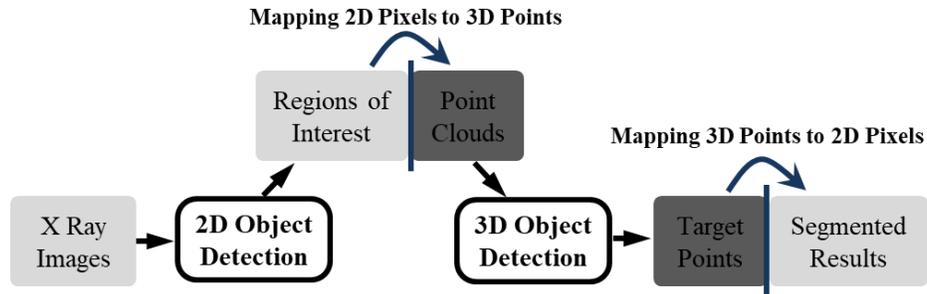
95 The proposed image segmentation method maps each pixel in the medical grayscale image to 3D
96 coordinates as the pixel-features point cloud, according to their positions and gray values. By
97 acquisition of foreground points and their corresponding bounding box using 3D object detection
98 method, we could achieve threshold values and the segmentation result of the corresponding
99 grayscale image. The whole pipeline and the implementation flow of this method are shown in
100 Fig.2 and Fig.3 respectively. Given a grayscale medical image, after (1) obtaining interest regions
101 of associated segmentation objects in the image, (2) generating 3D bounding box proposals in
102 point cloud and (3) the regression of their locations and scales, the refined boxes could be achieved.
103 The projection of points in refined bounding box into the 2D image is the segmentation result.



104

105

Fig.2 The pipeline of the proposed grayscale medical image segmentation method.



106
107

Fig.3 The implementation flow of the proposed method.

108 *2.2.1 Related Work*

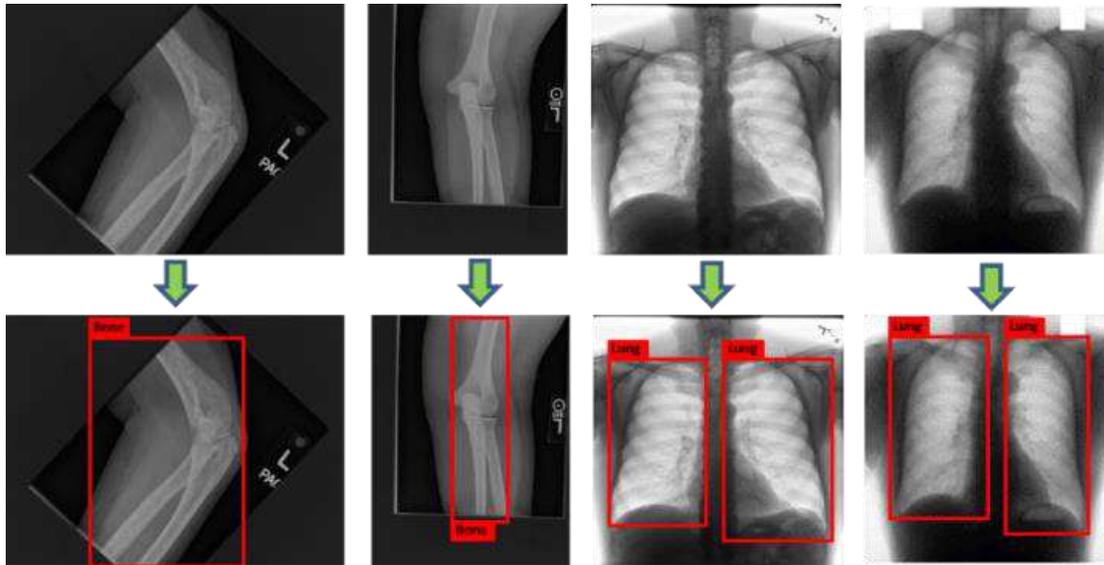
109 According to the proposed strategy and above pipeline, object detections play the central roles at
 110 each block of our method. Many researches about 2D&3D object detection has raised ever and
 111 they could perform well especially those with deep learning.

112 The current mainstream 2D object detection methods based on deep learning could be generally
 113 classified into two-stage and one-stage methods. [25] With two-stage methods, proposal bounding
 114 boxes are generated firstly and the further refinement of proposals and confidences is obtained in
 115 the second stage. [26] While using the one-stage methods [27,28], the location and the
 116 classification of object bounding boxes could be estimated directly without refinement which
 117 means one-stage methods are usually faster than two-stage ones but have lower object detection
 118 accuracy. [29]

119 The widespread application of 3D geometric data spurs the development of 3D object detection
 120 and it could be categorized into monocular/stereo image-based, point cloud-based and multimodal
 121 fusion-based methods in terms of the modality of input data. [30] Due to point clouds are the most
 122 regular data which could be achieved with different sensors, enormous researches of point cloud-
 123 based methods have raised. [31-33] Among these method, different data format like raw point
 124 clouds or 3D voxel grids transformed from points could be feed into deep net architectures to find
 125 targets with bounding boxes and their classes. [34]

126 2.2.2 Achievement of interest regions in image

127 In a medical grayscale image, pixels of the segmentation object always just take up a part of the
128 entire image and there may exists noisy pixels with the same gray values in irrelevant regions.
129 Therefore, 2D object detection is adopted as the pre-processing procedure to identify the specially
130 interest regions with segmentation objects and reduce noisy pixels as shown in Fig.4.

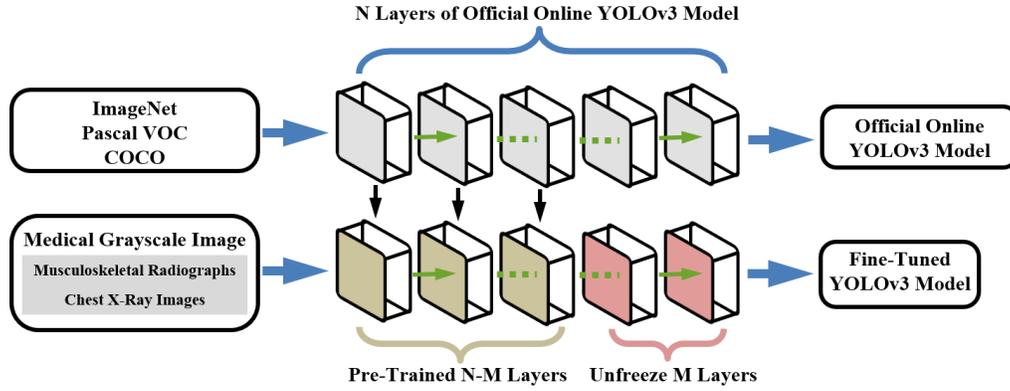


131

132 **Fig.4** Achievement of interest regions in 2D images.

133 Compared with the accuracy, the proposed pre-processing procedure cares more about the
134 detection speed, so we adopt the one-stage method YOLOv3 [35,36] as the backbone network.
135 And considering the scarcity of labeled medical grayscale images, we apply the fine tuning - a
136 transfer learning method [29] to migrate most layers of the backbone model which was pretrained
137 on ImageNet, Pascal VOC (Pattern analysis, statistical modeling and computational learning visual
138 object classes) and MS COCO (Microsoft common objects in context) datasets. [37,38] As Fig.5.
139 shown, with fine tuning method, we could freeze N-M layers of pre-trained model and only train
140 the last M layers on local dataset. In order to retain the detection ability of pre-trained model as
141 much as possible, and ensure the stability of the loss change during the training process, the

142 proposed image segmentation pre-processing method only unfreeze the last 3 layers of pre-trained
 143 network for training.



144

145 **Fig.5** The proposed 2D object detection network with fine-tuning method.

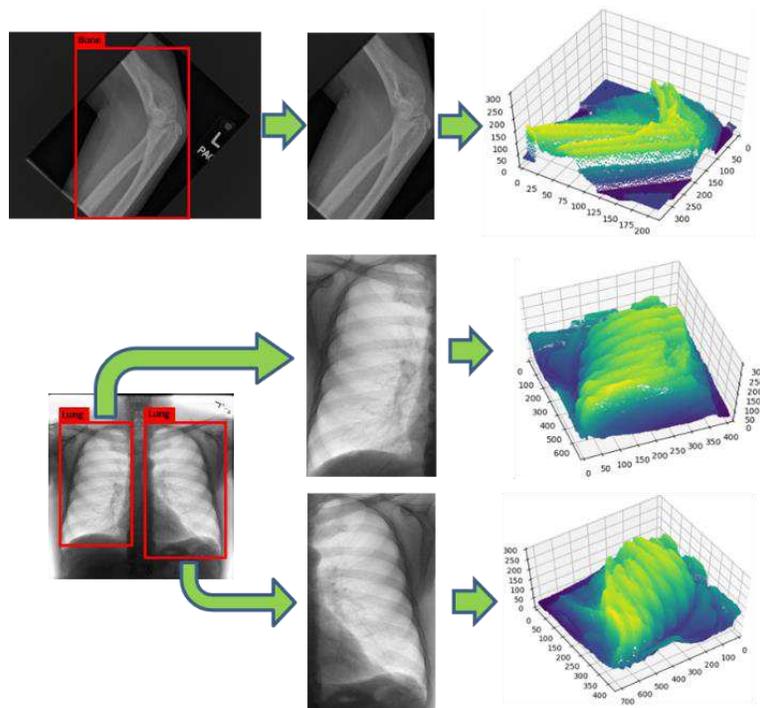
146 *2.2.3 Generation of proposal bounding box in pixel-features point cloud*

147 The grayscale value of each pixel in interest regions represents their brightness. [39] Pixels
 148 compose the same tissues in particular image always share the grayscale value ranges and we could
 149 recognize them manually. All values range from 0 to 255 (Typically zero is taken to be black, and
 150 255 is taken to be white). Darker pixels represent structures like soft tissues having less attenuation
 151 to the beam, while light ones represent structures like bones having high attenuation. Due to the
 152 lack of detailed gray values of pixels displayed on 2D images, it is hard to determinate their specific
 153 grayscale value ranges.

154 Thus, we turn pixels in 2D interest regions into the 3D representations as Fig.6. shown. In Fig.6.
 155 the first two dimensions represent pixels locations and the third dimension represents their
 156 grayscale values. The 3D data could be considered as the pixel-features point cloud and it is distinct
 157 and intuitive to obtain points which represent pixels belong to the same tissues. This helps us
 158 translate the 2D image segmentation task into the 3D object detection with point cloud. We only
 159 need to determine locations and widths of 3D bounding boxes which contain the foreground points

160 during the object detection. Then bottoms and tops of bounding boxes could represent the
161 segmentation required threshold values for 2D images.

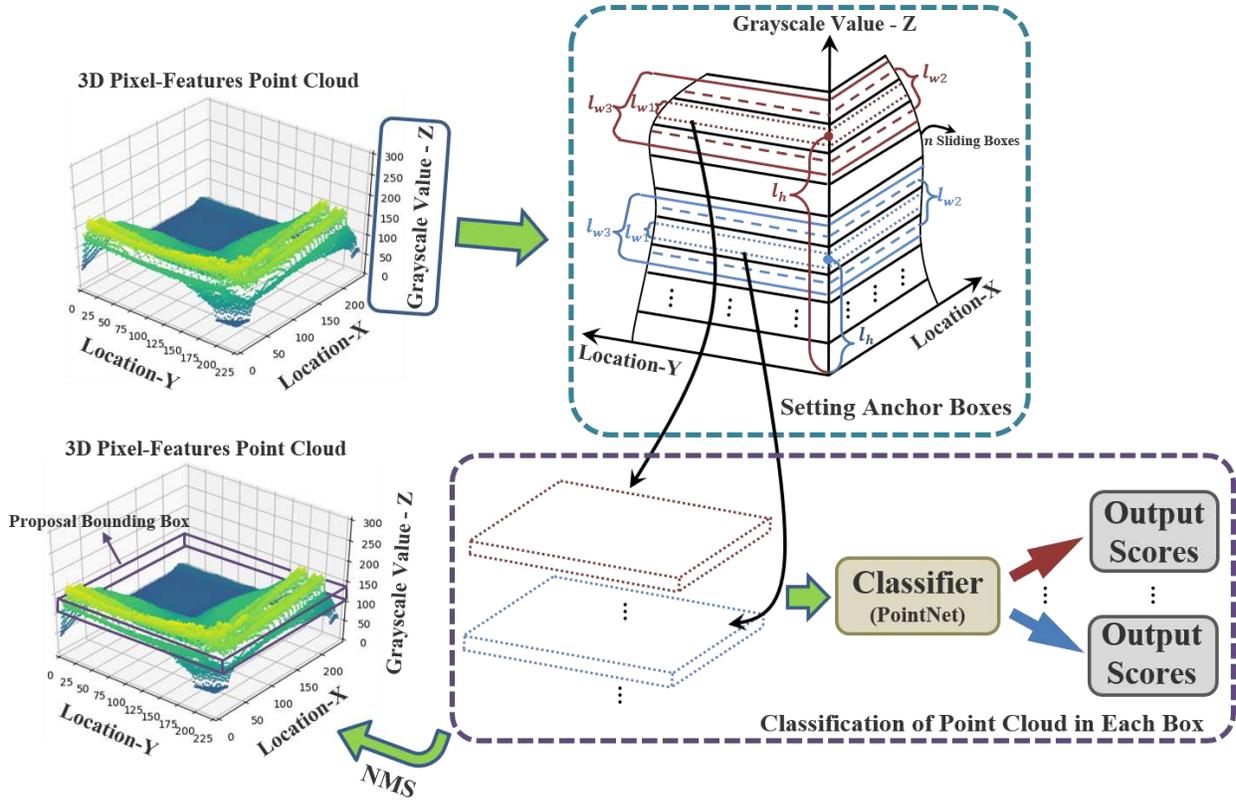
162 Inspired by two-stage 2D object detection methods, we present a novel two-stage 3D object
163 detection method, which is operated on pixel-features point cloud. In the first stage of existing
164 popular two-stage 2D object detection method, the proposal bounding boxes with its classification
165 scores are generated with convolutional neural network and the refinements of those boxes are
166 obtained in the following stage after the Non-Maximum Suppression (NMS). While in our
167 proposed 3D object detection method, based on two-stage strategy, the proposal 3D bounding
168 boxes with the classification scores of points inside them are estimated firstly and these proposals
169 are refined with regression in second stage.



170
171 **Fig.6** Turning pixels in interest regions into the pixel-features point cloud.

172 The generation of proposal bounding boxes in pixel-features point cloud has three modules. As
173 shown in Fig.7, These modules include localization of anchor boxes, classification of points inside

174 boxes utilizing PointNet [34] as backbone network and Non-Maximum Suppression with 3D
 175 Intersection-over-Union (IoU).



176
 177 **Fig.7** The generation of proposal bounding boxes in pixel-features point cloud.

178 *2.2.3.1 Anchor boxes*

179 Proposal bounding boxes generation takes the $l_x \times l_y \times 255$ point cloud representation as input
 180 where l_x and l_y respectively indicate the length and width of 2D interest region. In order to avoid
 181 high overlap rate of predict boxes and the low search efficiency using selective search as Region
 182 Convolutional Neural Network (RCNN) method, inspired by the Region Proposal Networks
 183 (RPN) in Faster RCNN, we apply the anchor boxes method for electing predict boxes.
 184 To generate proposals, we slide a small network over the input by a shared 3D convolutional layer
 185 referred to RPN and Single Shot MultiBox Detector (SSD) method as Fig.7. shown. At each

186 sliding-box location, we could predict multiple proposals simultaneously, and we denote the
187 maximum number of possible proposals as k . These proposals are parameterized relative to k 3D
188 anchor boxes. Each anchor is centered at its corresponding sliding box and is associated with a
189 scale. Each anchor is defined with coordinates (l_h, l_w) where l_h and l_w represent its location and
190 scale. We apply 3 scales by default, deciding $k = 3$ anchors at each sliding box and $n \times k$ anchors
191 in total.

192 *2.2.3.2 Classification of point cloud*

193 Anchor boxes with different scales share the same box-length l_x and box-width l_y , and they are
194 distinguished by their center locations and box-heights. In order to determine the proposal
195 bounding box from numerous anchor boxes, we utilize the PointNet as our backbone network and
196 apply the fine-tuning method for training our classification module.

197 The classification network in Fig.7. indicates that raw point clouds are directly taken as the input
198 and each point is processed independently at the initial stage. Due to point clouds could be easily
199 applied rigid or affine transformations, input points are sorted into a canonical order with the first
200 affine transformation by a mini-net (T-net) and moreover, after points features extraction with
201 multi-layer perceptron (mlp), features from different points could also be aligned using another
202 alignment network by feature transformation matrix. Then, the max pooling layer aggregates all
203 points features extracted from the second mlp and outputs the global features. The final fully
204 connected layers set the global feature as input and outputs k scores for all the k candidate classes.

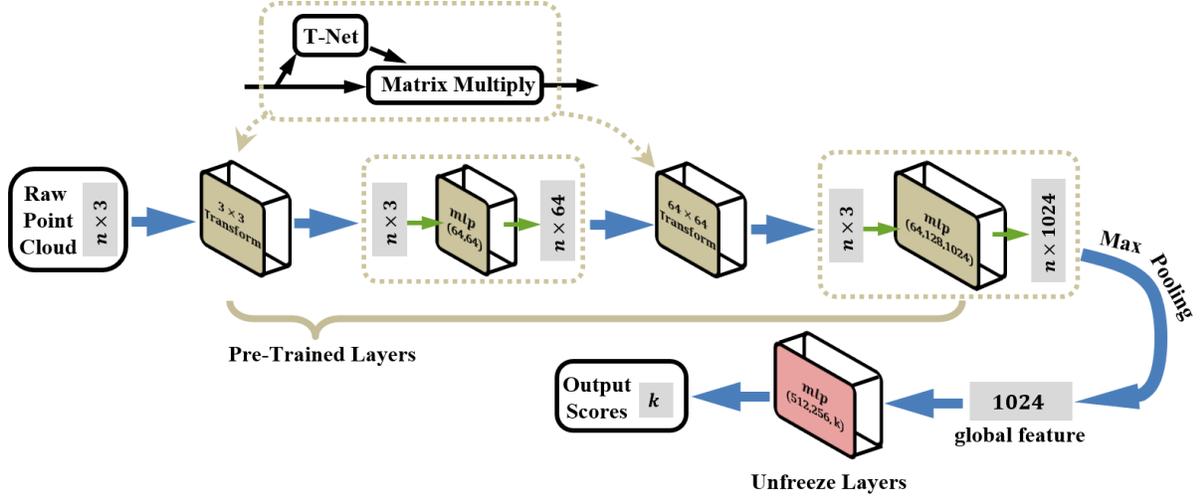


Fig.8 The proposed point cloud classification network with fine-tuning method.

205

206

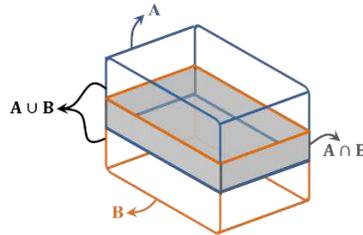
207 It should be noted that models-based point clouds datasets which mapped from grayscale medical
 208 images is scarce, thus we apply the fine-tuning method again. With the migration of PointNet
 209 model pretrained on ModelNet40 [40], we freeze most layers of the network except the final fully
 210 connected layers as shown in Fig.8.

211 2.2.3.3 NMS with 3D IoU

212 After the above module, the classification results of point cloud in each anchor box could be
 213 achieved with scores. But as many 2D object detection method, there exists some repeated
 214 proposals of one object. They belong to the same candidate class and overlap with the local
 215 highest-score box. For reducing the redundancy, we adopt the non-maximum NMS on these
 216 proposals with 3D intersection over union (3D IoU). Different from the IoU computation for 2D
 217 based on the relationships of areas between box A and B [41], like Fig.9. shows, volumes of two
 218 boxes are applied for 3D IoU calculation [42] which could be formulated as:

$$219 \quad 3D \text{ IoU}(A, B) = \frac{A_v \cap B_v}{A_v \cup B_v} = \frac{A_v \cap B_v}{|A_v| + |B_v| - A_v \cap B_v} \quad (1)$$

220 Through the setting of 3D IoU threshold for NMS and ranking with classification scores, it remains
 221 only one box for each candidate class which could be considered as the proposal bounding box.



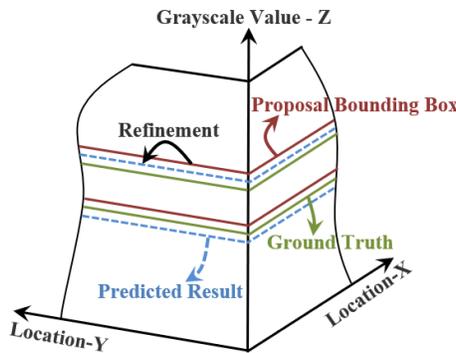
222

223 **Fig.9** IoU computation for 3D. The intersection volume is highlighted in gray.

224 *2.2.4 Refinement of proposal bounding box*

225 Even though high classification scores of the proposal bounding boxes, the location and scale
 226 errors between them and ground truth exist. We train and implement a class-specific bounding box
 227 linear regression model to reduce errors and improve detection performance.

228 On the assumption that we achieve one proposal bounding box P^i and its nearby ground-truth box
 229 G^i as shown in Fig.10, where $P^i = (P_{l_h}^i, P_{l_w}^i)$ specifies height l_h of the center of proposal
 230 bounding box together with its width l_w . Meanwhile, the ground-truth bounding box G^i is
 231 specified in the same way: $G^i = (G_{l_h}^i, G_{l_w}^i)$. The goal of the bounding box regressor is to learn a
 232 transformation which could map each proposal bounding box P to the ground-truth box G .



233

234 **Fig.10** Refinement of proposal bounding box.

235 The transformation could be parameterized in terms of two functions $d_{l_h}(P)$ and $d_{l_w}(P)$. The first
 236 function specifies the translation of bounding box P 's center which is scale-invariant, while the
 237 second specifies the log-space translation of its width. By applying the transformation as following
 238 equations, an input proposal bounding box P could be transformed into a predicted ground-truth
 239 box \hat{G} .

$$240 \quad \hat{G}_{l_h} = P_{l_w} \times d_{l_h}(P) + P_{l_h} \quad (2)$$

$$241 \quad \hat{G}_{l_w} = P_{l_w} \times \exp(d_{l_w}(P)) \quad (3)$$

242 Inspired by the 2D object detection, the bounding box regression of our method is performed on
 243 global features which is max pooled from PointNet model. Above two functions $d_{l_h}(P)$ and
 244 $d_{l_w}(P)$ could be modeled as linear functions of the global features of proposal bounding box P ,
 245 denoted as $f_{mp}(P)$. Therefore, we have $d_*(P) = T_* \times f_{mp}(P)$, where $*$ represents l_h or l_w , and T_*
 246 is a vector composed of learnable model parameters.

247 The transformation targets t_* between proposal bounding box P and the real ground-truth box G
 248 could be defined as:

$$249 \quad t_{l_h} = \frac{G_{l_h} - P_{l_h}}{P_{l_w}} \quad (4)$$

$$250 \quad t_{l_w} = \log\left(\frac{G_{l_w}}{P_{l_w}}\right) \quad (5)$$

251 Thus, after setting the loss function and by optimizing the regularized least squares objective as
 252 following, we could learn T_* and achieve the transformation to refine the proposal bounding box.

$$253 \quad \mathbf{Loss} = \sum_i^N \left(t_*^i - \hat{T}_* \times f_{mp}(P^i) \right)^2 \quad (6)$$

$$254 \quad T_* = \operatorname{argmin}_{\hat{T}_*} \mathbf{Loss} + \lambda \|\hat{T}_*\|^2 \quad (7)$$

255 2.2.5 Training strategy

256 The proposed grayscale medical image segmentation method follows a three-stage training
257 strategy. First, we obtain interest regions from raw grayscale images with fine-tuning YOLOv3
258 model. Second, by training the pixel-features point cloud classification model based on PointNet,
259 proposal 3D bounding boxes could be achieved from the point cloud representations of pixels in
260 interest regions. By training the linear regressor, proposal bounding boxes are refined with location
261 and scale transformation. Three independent modules including regions extractor, point cloud
262 classifier and bounding box regressor in three stages compose our method.

263 2.3 Performance assessment

264 In this study, we evaluate the segmentation performance by following four metrics: Dice similarity
265 coefficient (DSC) scores [6], intersection over union (IoU), False negative (FN) and False positive
266 (FP) [7]. Ranges of DSC and IoU are between 0 and 1, higher values of them and lower values of
267 FN and FP indicate the higher accuracy. The calculation formula of DSC is defined as:

$$268 \quad \text{DSC} = \frac{2|T \cap G|}{|T| + |G|} \quad (8)$$

269 where T is the detected region and G is the ground truth region.

270 3 Results

271 We conduct experiments by the proposed grayscale image segmentation method on above
272 mentioned datasets including musculoskeletal radiographs dataset and chest radiographs dataset.
273 Moreover, our prepared phalanx and forearm X ray images obtained with the portable X ray
274 machine as Fig.11. shown are also adopted for model training and validation.

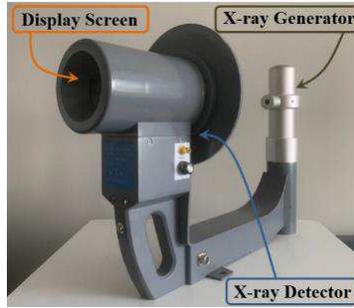
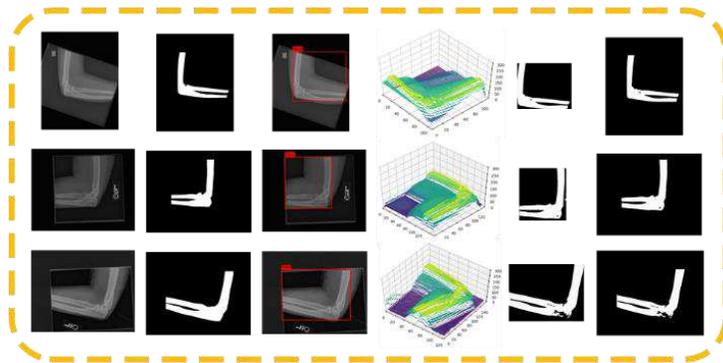


Fig.11 The portable X ray machine applied in experiments.

275
276
277
278
279
280
281
282
283
284
285
286
287
288

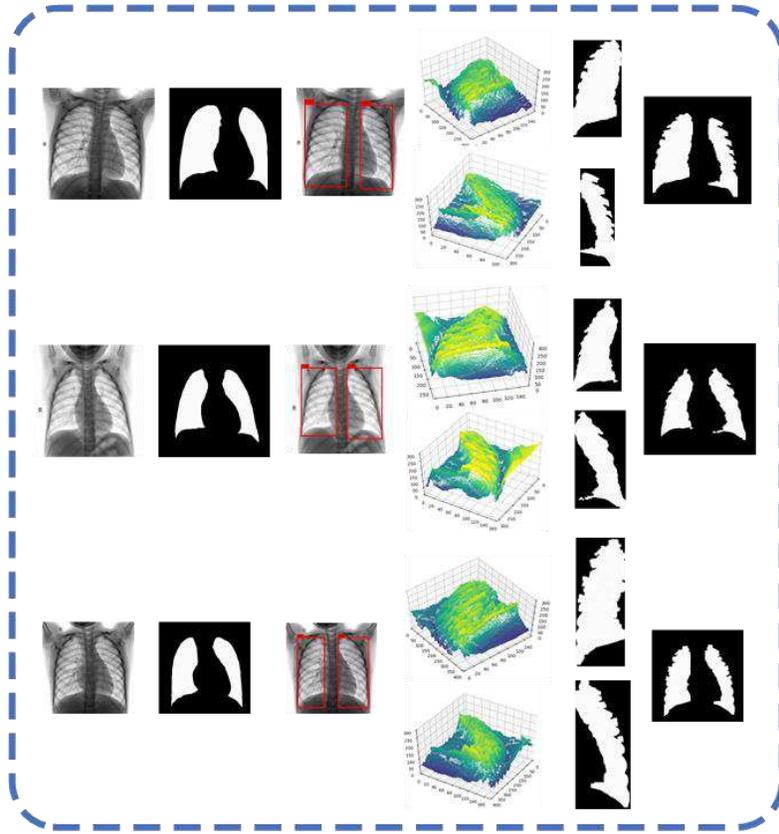
Our model is implemented with Pytorch [43] and its entire training process is performed on a computer with Windows 10 operating system, Intel Core i7 processor with 3.0 GHz, 64GB of RAM and a single NVIDIA GPU (Quadro RTX 4000). The 2D object detection model is trained with 50 epochs for achieving interesting regions and it takes 1.75 h, while the training of the 3D object detection model for generating proposal bounding boxes spends 2.5 h on 200 epochs.

After training process, by applying the proposed method with the given grayscale medical images input and following the method pipeline as Fig.2. shown, regions of target issues could be segmented. Each block in Fig.12. presents several examples of segmentation performance from different kinds of datasets, as well as processing results after each stage, where white represents true positive pixels and black is for true negatives pixels. Moreover, according to evaluation criteria, Table 1 shows four metrics including IoU, DSC, FN and FP to assess the segmentation performance of images in different datasets.

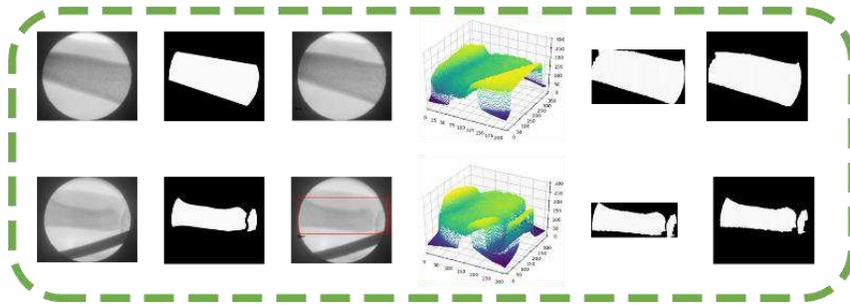


289
290

(a) Examples of segmentation performance in musculoskeletal radiographs dataset.



(b) Examples of segmentation performance in chest radiographs dataset.



(c) Examples of segmentation performance in X ray images with the portable X ray machine.

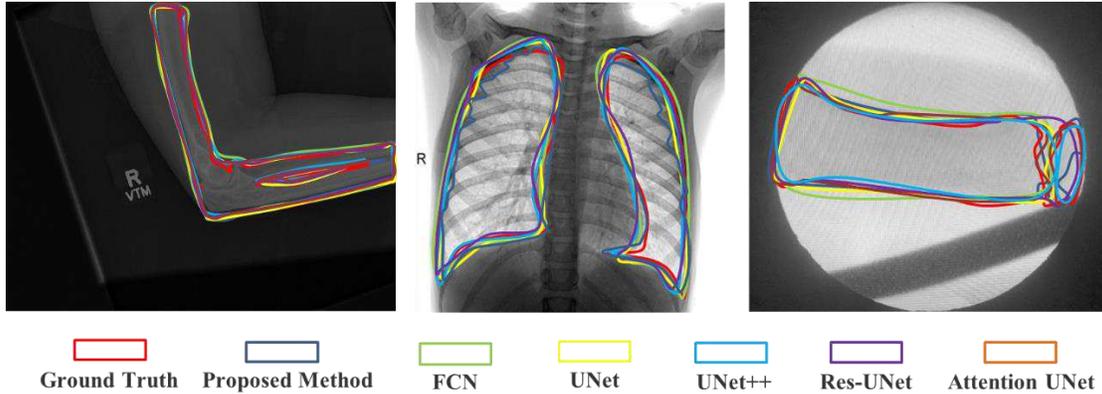
Fig.12 Segmentation results from different kinds of datasets. From the first to the last column are origin images, ground truth, achievements of interest regions, representations of pixel-feature point cloud, local segmentation results, and segmentation results in original image size, respectively.

Table 1 The values of evaluative metrics from experiments in different datasets.

Datasets	IoU	DSC	FN	FP
Musculoskeletal radiographs	0.92	0.96	0.05	0.02
Chest radiographs	0.88	0.93	0.11	0.15
Images from X ray machine	0.94	0.94	0.06	0.08

Table 2 Comparison between segmentation performance (IoU) of the proposed approach with other methods.

Datasets	Proposed	FCN	UNet	UNet++	Res-UNet	Attention Unet
Musculoskeletal radiographs	0.92	0.82	0.85	0.84	0.91	0.90
Chest radiographs	0.88	0.76	0.81	0.83	0.88	0.86
Images from X ray machine	0.94	0.72	0.82	0.87	0.85	0.91



300

301

Fig.13 Performance comparison of grayscale medical image segmentation with different methods.

302

As shown in Fig.12. and Table 1, we could obtain high IoU and DSC scores with satisfied

303

segmentation results on different datasets. This indicates that based on the proposed method, 2D

304

interest regions and 3D bounding boxes containing target pixel-features point cloud during the

305

processing could be successfully achieved.

306

4 Discussion

307

In this section, we compare the image segmentation performance of the proposed method with

308

multiple famous and clinically performed well models. As well known, CNN based models are

309

among the most successful and widely used for medical image processing. Besides the milestone

310

FCN model, UNet built on top of the fully convolutional networks with a U-shaped architecture

311

to capture context information, and based on it, Res-UNet [44] improved the segmentation results

312

using residual blocks as the building block and UNet++ [45] enhanced segmentation quality of

313

varying-size objects. Also, Attention UNet [46] achieved the better performance with the attention

314

gate. We train these models in the same dataset as our proposed method and Table 2 presents the

315 comparison results. Meanwhile, Fig.13. shows results by visualization. It indicates that compared
316 with other models, our proposed approach improves the segmentation performance and it obtains
317 the highest IoU scores of 0.92, 0.88 and 0.94 with three datasets respectively. In our approach, 2D
318 and 3D object detection models could be both trained with transfer learning method which makes
319 it possible to achieve a quite accurate image segmentation model with small training datasets.
320 While other semantic segmentation methods may be sensitive to the scale of datasets because the
321 pre-trained model could only help simplify the downsample training procedure, and the training
322 of upsample still requires a number of datasets. This indicates that it is impossible to adapt them
323 for every application task well because training data is scarce especially in medical image field.
324 Moreover, in grayscale images, grayscale values of pixels are important features to distinguish
325 different objects, and the intuitive logic of grayscale image segmentation could be considered as
326 the collection of pixels with similar grayscale values. So, the proposed image segmentation model
327 which obtains the purpose ranges of grayscale values with 3D object detection have better
328 explicability and segmentation effect.

329 Under different medical imaging devices and environment in clinical, ranges of grayscale values
330 of pixels which compose the same segmentation target in different medical images are always
331 different. But our proposed method could settle this and we could obtain thresholding values (top
332 and bottom of 3D bounding boxes) by mapping pixels in 2D images into 3D point clouds and
333 adopting 3D object detection with features of pixels.

334 **5 Conclusions**

335 In this paper, we present a new method for grayscale medical image segmentation only with two
336 object detection models. The method applies 2D object detection model for location identification
337 of segmentation objects. It could crop the origin images and increase the efficiency of further

338 detailed segmentation. Pixels in interest regions are mapped as point cloud according to their
339 positions and grayscale values. Using 3D object detection methods, we achieve bounding boxes
340 which contain target pixels-feature points. After projecting these points to 2D images, they could
341 composite the segmentation results. The effectiveness of the proposed image segmentation method
342 is proven by several experiments in different image datasets and the comparison with other famous
343 approaches and it indicates the proposed method could perform better in grayscale image
344 segmentation tasks. In further research, we will concentrate on multi-oriented objects detection
345 technologies for more fine segmentation results.

346 **Abbreviations**

347 CNN(s): Convolutional neural network(s); FCN: Fully convolutional networks; MURA:
348 Musculoskeletal radiographs; LERA: Lower extremity radiographs; CheXpert : Chest radiography;
349 GT: Ground Truth; YOLO: You only look once; PASCAL VOC: Pattern analysis, statistical
350 modeling and computational learning visual object classes; MS COCO: Microsoft common objects
351 in context; NMS: Non-maximum suppression; IoU: Intersection-over-Union; RCNN: Region
352 convolutional neural network; RPN: Region proposal networks; SSD: Single shot multiBox
353 detector; mlp: multi-layer perceptron; DSC: Dice similarity coefficient; FN: False negative; FP:
354 False positive.

355 **Declarations**

356 **Ethics approval**

357 We declare that all of us obey the principles of the Declaration of Helsinki. In other words, all
358 experiments and methods in this paper are in accordance with these principles. The study was
359 approved by the Ethics Committee of the First people's Hospital of Yancheng.

360 **Consent to participate**

361 The fully anonymized phalanx and forearm X ray images were received by authors on 2 April,
362 2021 and the requirement for informed consent was waived for this study because of the
363 anonymous nature of the data.

364 **Consent for publication**

365 Not applicable for this paper

366 **Availability of data and materials**

367 Musculoskeletal radiographs and chest radiographs which support our research are available from
368 Stanford ML Group. But restrictions apply to the availability of these data, which were used under
369 license for the current study, and so are not publicly available. Data are however available from
370 the authors upon reasonable and with permission of Stanford ML Group. While phalanx and
371 forearm X ray images are available only upon request by emailing authors due to the ethical
372 restrictions on sharing these data which could contain potentially sensitive information of patients.

373 **Competing interests**

374 All authors declare that they have no interest conflicts or competing interests.

375 **Founding**

376 This work was supported by the project of Tongji University Sheng Feiyun College Student
377 Science and Technology Innovation Practice Found.

378 **Authors' contributions**

379 Qing Zhang conceived the research. Yunfei Ge and Yidong Shen analyzed the clinical and imaging
380 data. Yuantao Sun, Yunfei Ge, and Xijiong Wang designed the study. Yunfei Ge and Yidong Shen
381 performed the experiments and collected the results. Yunfei Ge and Yuantao Sun drafted the
382 manuscript. Qing Zhang reviewed the final manuscript. All authors read and approved the final
383 manuscript.

384 **Acknowledgements**

385 Not applicable.

386 **References**

- 387 1. Justine Wallyn, Anton Nicolas, Akram Salman, et al. Biomedical imaging: principles, technologies,
388 clinical aspects, contrast agents, limitations and future trends in nanomedicines. *Pharmaceutical*
389 *Research*. 2019; 36(6):78-108.
- 390 2. Yeo W K, Yap D F W, et al. Grayscale medical image compression using feedforward neural networks.
391 2011 IEEE International Conference on Computer Applications and Industrial Electronics (ICCAIE).
392 2011; 633-638.
- 393 3. Lei Tao, et al. Medical Image Segmentation Using Deep Learning: A Survey. *arXiv*. 2020; 13120.
- 394 4. Rathnayaka K, Sahama T, Schuetz MA, et al. Effects of CT image segmentation methods on the
395 accuracy of long bone 3D reconstructions. *Medical Engineering & Physic*. 2011; 33(2): 226-233.

- 396 5. Shuo Wang, Zhou Mu, Liu Zaiyi, et al. Central focused convolutional neural networks: Developing a
397 data-driven model for lung nodule segmentation. *Medical Image Analysis*. 2017; 40: 172-183.
- 398 6. Han Liu, Wang Lei, Nan Yandong, et al. SDFN: Segmentation-based deep fusion network for thoracic
399 disease classification in chest X ray images. *Computerized Medical Imaging and Graphics*. 2019; 75:
400 66-73.
- 401 7. de Albuquerque VHC, Rodrigues D A, Ivo RF, et al. Fast fully automatic heart fat segmentation in
402 computed tomography datasets. *Computerized Medical Imaging and Graphics*. 2020; 80: 101674.
- 403 8. Li Wen, et al. Automatic segmentation of liver tumor in CT images with deep convolutional neural
404 networks. *Journal of Computer and Communications*. 2015; 3(11): 146.
- 405 9. Vivanti R, Ephrat A, Joskowicz L, et al. Automatic liver tumor segmentation in follow-up CT studies
406 using convolutional neural networks. *Proc. Patch-Based Methods in Medical Image Processing*
407 *Workshop*. 2015; 2: 2.
- 408 10. Saleha Masood, Sharif Muhammad, Masood Afifa, et al. A survey on medical image segmentation.
409 *Current Medical Imaging*. 2015; 11(1): 3-14.
- 410 11. Khandare ST, Isalkar A D. A survey paper on image segmentation with thresholding. *International*
411 *Journal of Computer Science and Mobile Computing*. 2014; 3(1): 441-446.
- 412 12. Sezgin M, Sankur B. Survey over image thresholding techniques and quantitative performance
413 evaluation. *Journal of Electronic Imaging*. 2004; 13(1): 146-165.
- 414 13. Maolood I Y, Al-Salhi Y E A, Lu S. Thresholding for medical image segmentation for cancer using
415 fuzzy entropy with level set algorithm. *Open Medicine*. 2018; 13(1): 374-383.
- 416 14. Duo Hao, Li Qiuming, Li Chengwei. Histogram-based image segmentation using variational mode
417 decomposition and correlation coefficients. *Signal, Image and Video Processing*. 2017; 11(8): 1411-
418 1418.
- 419 15. Long J, Shelhamer E, Darrell T. Fully convolutional networks for semantic segmentation. *Proceedings*
420 *of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2015; 3431-3440.
- 421 16. Ronneberger O, Fischer P, Brox T. U-net: Convolutional networks for biomedical image segmentation.
422 *International Conference on Medical Image Computing and Computer-Assisted Intervention*
423 *(MICCAI)*. 2015; 234-241.
- 424 17. Chen LC, Papandreou G, Kokkinos I, et al. Deeplab: Semantic image segmentation with deep
425 convolutional nets, atrous convolution, and fully connected crfs. *IEEE Transactions on Pattern*
426 *Analysis and Machine Intelligence*. 2017; 40(4): 834-848.
- 427 18. Deng J, Dong W, Socher R, et al. Imagenet: A large-scale hierarchical image database. *Proceedings of*
428 *the IEEE Conference on Computer Vision and Pattern Recognition*. 2009; 248-255.
- 429 19. Kalinin A A, Iglovikov V I, Rakhlin A, et al. Medical image segmentation using deep neural networks
430 with pre-trained encoders. *Deep Learning Applications*. 2020; 39-52.
- 431 20. Pierre-Henri Conze, Brochard Sylvain, Burdin Val-E-Rie, et al. Healthy versus pathological learning
432 transferability in shoulder muscle MRI segmentation using deep convolutional encoder-decoders.
433 *Computerized Medical Imaging and Graphics*. 2020; 83: 101733.
- 434 21. Pranav Rajpurkar, Irvin Jeremy, Bagul Aarti, et al. Mura: Large dataset for abnormality detection in
435 musculoskeletal radiographs. *arXiv*. 2017; 1712.06957.
- 436 22. LERA - lower extremity radiographs. <https://aimi.stanford.edu/lera-lower-extremity-radiographs-2>.
- 437 23. Irvin J, Rajpurkar P, Ko M, et al. Chexpert: A large chest radiograph dataset with uncertainty labels
438 and expert comparison. *Proceedings of the AAAI Conference on Artificial Intelligence*. 2019; 33(01):
439 590-597.
- 440 24. Joseph-Paul Cohen, Morrison Paul, Dao Lan, et al. Covid-19 image data collection: Prospective
441 predictions are the future. *arXiv*. 2020; 2006.11988.
- 442 25. Jiao L, Zhang F, Liu F, et al. A survey of deep learning-based object detection. *IEEE Access*. 2019;
443 7:128837-128868.
- 444 26. Ross Girshick, Donahue Jeff, Darrell Trevor, et al. Rich feature hierarchies for accurate object detection
445 and semantic segmentation. *Proceedings of the IEEE Conference on Computer Vision and Pattern*
446 *Recognition*. 2014; 580-587.

- 447 27. Joseph Redmon, Divvala Santosh, Girshick Ross, et al. You only look once: Unified, real-time object
448 detection. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2016;
449 779-788.
- 450 28. Wei Liu, Anguelov Dragomir, Erhan Dumitru, et al. SSD: Single shot multibox detector. European
451 Conference on Computer Vision. 2016; 21-37.
- 452 29. Shin HC, Roth H R, Gao M, et al. Deep convolutional neural networks for computer-aided detection:
453 CNN architectures, dataset characteristics and transfer learning. IEEE Transactions on Medical
454 Imaging. 2016; 35(5): 1285-1298.
- 455 30. Qian R, Lai X, Li X. 3D Object Detection for Autonomous Driving: A Survey. arXiv. 2021;
456 2106.10823.
- 457 31. Zhou Y, Tuzel O. Voxelnet: End-to-end learning for point cloud based 3d object detection. Proceedings
458 of the IEEE conference on computer vision and pattern recognition. 2018: 4490-4499.
- 459 32. Chen Y, Liu S, Shen X, et al. Fast point r-cnn. Proceedings of the IEEE/CVF International Conference
460 on Computer Vision. 2019: 9775-9784.
- 461 33. Shi S, Wang X, Li H P. 3d object proposal generation and detection from point cloud. Proceedings of
462 the IEEE Conference on Computer Vision and Pattern Recognition, Long Beach, CA, USA. 2019: 16-
463 20.
- 464 34. Qi CR, Su H, Mo K, et al. Pointnet: Deep learning on point sets for 3d classification and segmentation.
465 Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2017; 652-660.
- 466 35. Redmon J, Farhadi A. Yolov3: An incremental improvement. arXiv. 2018; 1804.02767.
- 467 36. Rasmus Rothe, Guillaumin Matthieu, Van Gool Luc. Non-maximum suppression for object detection
468 by passing messages between windows. Asian Conference on Computer Vision. 2014; 290-306.
- 469 37. Everingham M, Van Gool L, Williams C K, et al. The pascal visual object classes (voc) challenge: A
470 Retrospective. International Journal of Computer Vision. 2014; 111: 98-136.
- 471 38. Lin T Y, Maire M, Belongie S, et al. Microsoft COCO: Common Objects in Context. European
472 Conference on Computer Vision. 2014; 740-755.
- 473 39. Tan L, Jiang J. Digital signal processing: fundamentals and applications. Academic Press; 2019.
- 474 40. Zhirong Wu, Song Shuran, Khosla Aditya, et al. 3d shapenets: A deep representation for volumetric
475 shapes. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2015; 1912-
476 1920.
- 477 41. Hamid Rezatofighi, Tsoi Nathan, Gwak JunYoung, et al. Generalized intersection over union: A metric
478 and a loss for bounding box regression. Proceedings of the IEEE Conference on Computer Vision and
479 Pattern Recognition. 2019; 658-666.
- 480 42. Zhou D, Fang J, Song X, et al. Iou loss for 2d/3d object detection. International Conference on 3D
481 Vision (3DV). 2019; 85-94.
- 482 43. Adam Paszke, Gross Sam, Massa Francisco, et al. Pytorch: An imperative style, high-performance deep
483 learning library. Advances in Neural Information Processing Systems. 2019; 32: 8026-8037
- 484 44. X. Xiao, S. Lian, Z. Luo and S. Li. Weighted Res-UNet for High-Quality Retina Vessel Segmentation.
485 2018 9th International Conference on Information Technology in Medicine and Education (ITME).
486 2018; 327-331.
- 487 45. Z. Zhou, M. M. R. Siddiquee, N. Tajbakhsh and J. Liang. UNet++: Redesigning Skip Connections to
488 Exploit Multiscale Features in Image Segmentation. IEEE Transactions on Medical Imaging. 2020;
489 39(6): 1856-1867.
- 490 46. Ozan Oktay, Jo Schlemper, et al. Attention U-Net: Learning Where to Look for the Pancreas. arXiv.
491 2018; 1804.03999.