

# Pose-Aware Instance Segmentation Framework from Cone Beam CT Images for Tooth Segmentation

Minyoung Chung, Minkyung Lee, Jioh Hong, Sanguk Park, Jusang Lee, Jingyu Lee, Jeongjin Lee\*, and Yeong-Gil Shin

**Abstract**—Individual tooth segmentation from cone beam computed tomography (CBCT) images is an essential prerequisite for an anatomical understanding of orthodontic structures in several applications, such as tooth reformation planning and implant guide simulations. However, the presence of severe metal artifacts in CBCT images hinders the accurate segmentation of each individual tooth. In this study, we propose a neural network for pixel-wise labeling to exploit an instance segmentation framework that is robust to metal artifacts. Our method comprises of three steps: 1) image cropping and realignment by pose regressions, 2) metal-robust individual tooth detection, and 3) segmentation. We first extract the alignment information of the patient by pose regression neural networks to attain a volume-of-interest (VOI) region and realign the input image, which reduces the inter-overlapping area between tooth bounding boxes. Then, individual tooth regions are localized within a VOI realigned image using a convolutional detector. We improved the accuracy of the detector by employing non-maximum suppression and multiclass classification metrics in the region proposal network. Finally, we apply a convolutional neural network (CNN) to perform individual tooth segmentation by converting the pixel-wise labeling task to a distance regression task. Metal-intensive image augmentation is also employed for a robust segmentation of metal artifacts. The result shows that our proposed method outperforms other state-of-the-art methods, especially for teeth with metal artifacts. Our method demonstrated 5.68% and 30.30% better accuracy in the F1 score and aggregated Jaccard index, respectively, when compared to the best performing state-of-the-art algorithms. The primary significance of the proposed method is two-fold: 1) an introduction of pose-aware VOI realignment followed by a robust tooth detection and 2) a metal-robust CNN framework for accurate tooth segmentation.

**Index Terms**—Cone beam computed tomography image segmentation, pose-aware tooth detection, pose regression neural network, tooth instance segmentation.

## I. INTRODUCTION

**D**IGITIZED orthodontic applications in dentistry have increased based on the development of cone beam computed tomography (CBCT) imaging. CBCT is a widely used medical imaging technique that provides high resolution 3D volumetric data. To build an effective computer-aided diagnosis system in orthodontic applications, such as oral treatment planning for tooth reformation and implant guide simulation, automatic segmentation of individual teeth from the CBCT images is an essential prerequisite (Fig. 1). However, accurate

segmentation of an individual tooth from a CBCT image is a challenging task owing to heterogeneous intensity distribution, unclear boundaries between the tooth root and alveolar bone (Fig. 2a), and diverse shapes and poses. Moreover, the majority of CBCT images contain severe metal artifacts that hinder the accurate segmentation of teeth (Fig. 2b). In this study, we propose a fully automated instance segmentation framework using 3D images of teeth that is robust to several challenging conditions of the teeth, such as, dynamic poses, missing teeth, inter-tooth proximity, and presence of severe metal artifacts.

An extensive amount of literature on individual tooth segmentation were proposed in the last few decades. Classical image processing methods that exploited region growing [1], morphological operations [2], and watershed algorithm [2], [3] were studied. Several works employed contour-based level-set methods [4]–[7] or shape-based registration methods [8], [9]. However, all the classical algorithms demonstrated limitations while handling the aforementioned challenging conditions, such as heterogeneous intensities, unclear boundaries, diverse anatomical poses, and presence of metal artifacts. Moreover, classical algorithms typically require manual seed points to perform tooth segmentation, which results in a semiautomated application. More recently, a few studies [10], [11] applied convolutional neural network (CNN) architectures to resolve individual tooth segmentation. Although the proposed CNN-based methods showed promising results over previous approaches, the metal artifact condition is still difficult to overcome, implying that accurate detection and segmentation is still a challenging task. Note that the teeth including severe metal artifacts were not previously researched, either in the methodological or in the experimental perspectives.

To employ a CNN for individual tooth segmentation, we chose the instance segmentation framework. Unlike other organs in medical imaging, such as liver or spleen, there are multiple instances of a tooth in a single CBCT image. Thus, individual tooth segmentation requires an instance segmentation technique rather than semantic segmentation, which is performed through preceding object detection. The primary challenges of individual tooth detection and segmentation arise from: 1) high overlapping ratio between instances (i.e., teeth) and 2) the existence of severe metal artifacts in CBCT images. These two primary issues degrade the accuracy of the individual tooth segmentation task.

In this study, we propose a fully automatic, hierarchical method that performs accurate individual tooth detection followed by a CNN based on single pixel-wise labeling to segment the tooth (Fig. 3). We first extracted the volume-

*Asterisk indicates corresponding author.*

M. Chung, M. Lee, J. Hong, S. Park, J. Lee, J. Lee, and Y.-G. Shin are with the Department of Computer Science and Engineering, Seoul National University, Republic of Korea (e-mail: chungmy@snu.ac.kr).

\*J. Lee is with the Department of Computer Science and Engineering, Soong-sil University, Republic of Korea (e-mail: profjlee@naver.com).

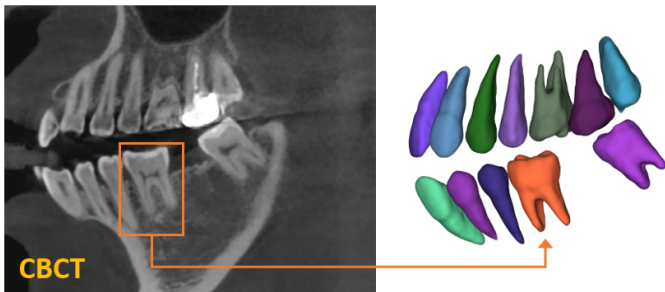


Fig. 1: Individual tooth segmentation from a cone beam computed tomography (CBCT) image. Example of surface visualization of tooth segmentation results.

of-interest (VOI) region by pose regression using neural networks. Then, the extracted VOI region was realigned based on the pose (i.e., axes). Subsequently, we detected individual teeth and performed individual tooth segmentation. The key components of our method are the accurate tooth detection framework through VOI realignment, which reduces the inter-overlapping area between boxes, multiclass classification within a detector, which boosts the accuracy of metal-tooth detection, and the metal-robust CNN for accurate tooth segmentation.

The remainder of this paper is structured as follows. In Section II, we review the related works on tooth segmentation methods and CNN architectures that perform object detection and segmentation. Further, we describe our proposed method in Section III. Section IV demonstrates experimental results and Section V presents the discussion and conclusion.

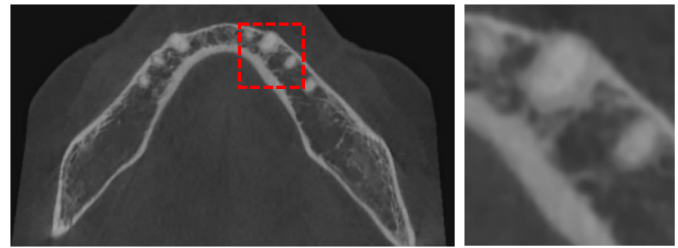
## II. RELATED WORKS

In this section, we first review the literature on tooth segmentation methods that were proposed until recently. In the following subsections, the CNN architectures for object detection and medical image segmentation will be highlighted.

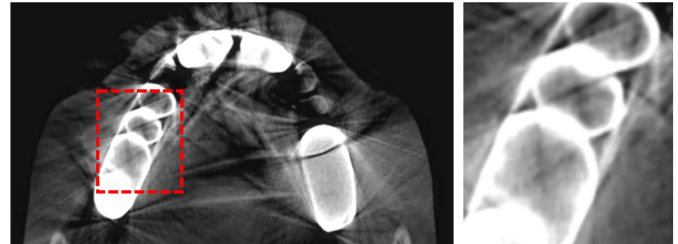
### A. Literature on Tooth Segmentation

1) *Classical Methods*: Classical image processing methods have been widely studied to achieve tooth segmentation [2], [3], [12], [13]. Several methods including region growing [1], watershed algorithm [2], [3], morphological operators [2], graph-cut-based segmentation [12], template-based registration [8], [9], and random forest classification [13] were implemented. Semiautomatic algorithms with manually annotated cues for easy implementation have also gained popularity [3], [14], [15]. The primary drawbacks of the classical methods are their reliance on the intensity or anatomical heuristics (i.e., assumptions). The algorithms break down in many cases where patient conditions do not satisfy the assumptions, such as intensity dynamics on metal-teeth or large anatomical shape variations (e.g., missing or misalignment of teeth).

2) *Level-set Methods*: The application of a level-set method to the individual tooth segmentation was extensively studied [4]–[7], [16]. In [4], two different level sets were introduced to handle the dynamic root branching and crown parts. The authors applied a coupled level-set method to resolve the



(a) Boundaries between the tooth root and alveolar bone.



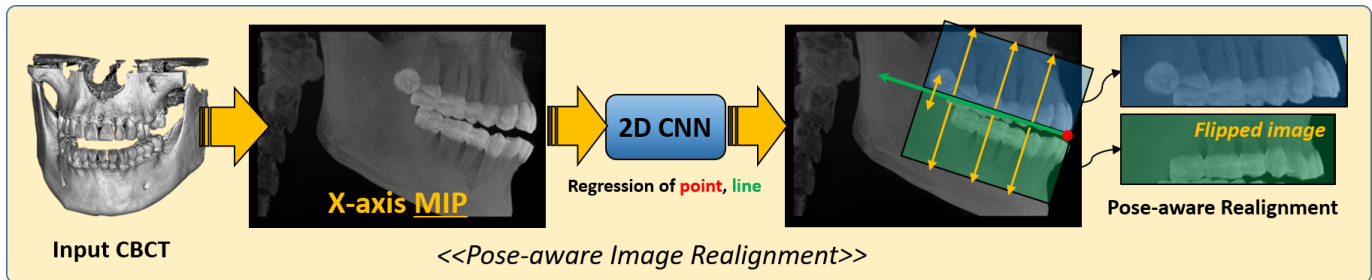
(b) Boundaries in metal artifact region.

Fig. 2: Unclear boundaries in CBCT images.

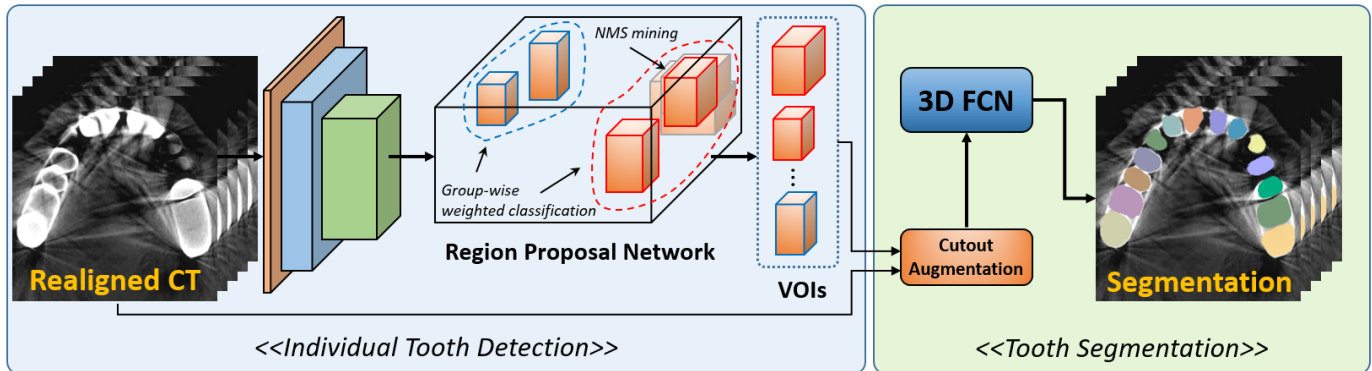
adjacency problem between neighboring crowns [4]. The intensity distribution was studied to enhance the performance of a level-set framework that models the inside and outside of the tooth [5]. To tackle the over-segmentation problem presented by the proximity of alveolar bone, simultaneous segmentation of tooth and alveolar bone was proposed [6]. The hybrid level-set model was presented in [7] by forming both local likelihood image fitting and prior shape constraint energy terms. The authors regularized the level-set functional by a reaction diffusion [7].

Despite all the appealing features from the contour propagating scheme, level-set methods demonstrate a common limitation, i.e., it is difficult to delineate unclear boundaries (Fig. 2). In a CBCT, the boundaries between the tooth root and the alveolar bone have weak edge characteristics (i.e., a low magnitude of gradient; Fig. 2a) and the presence of severe metal artifacts (Fig. 2b), which hinder the accurate propagation of a contour-based level-set method. Furthermore, level-set methods require tedious manual user interaction to define the initial contour [4], [5], [7], indicating a semiautomated algorithm.

3) *Convolutional Neural Network (CNN)-based Methods*: The body of literature focusing on CNNs for the individual tooth segmentation task is relatively low. Recently, a full-CNN-based method was presented that performs individual tooth segmentation [11]. The authors employed instance segmentation architecture using an end-to-end CNN framework similar to the Mask region-based CNN (Mask-RCNN) [17]. The singularities of the proposed method were the employment of edge map, similarity matrix, and combined tooth identification to boost the performance. However, metal artifacts, which are commonly used in dental clinics, were not considered in their proposed method.



(a) Architecture of the volume-of-interest (VOI) extraction and realignment. The input CBCT is first projected to a 2D image based on the maximum intensities. Point and line pairs are regressed by 2D convolutional neural network (CNN) for pose extraction. Finally, the two VOIs are realigned by fixed margins. We applied y-axis flipping for the lower part to align the teeth. A single point and line pair are visualized for simplicity.



(b) Our proposed instance segmentation framework. We adopted the backbone of the faster region-based CNN [18] encoder part and region proposal network. Non-maximum-suppression-based sampling and group-wise multiclass classification are applied to train the detector. The original CBCT image is then cropped to the detected boxes and fed to the 3D fully convolutional network (FCN) (Fig. 4) through cutout augmentation to perform individual tooth segmentation.

Fig. 3: Overall architecture of the proposed method: (a) Workflow of pose-aware image realignment, (b) Subsequent instance segmentation, which performs individual tooth detection and segmentation (Best viewed in color).

### B. Object Detection and Instance Segmentation

There are many elaborate studies on object detection based on CNNs [18]–[21]. The two major approaches are classification- and regression-based approaches. The classification-based object detection method attempts to identify boxes that contain an object by locally windowed classification [18], [21]. In contrast, the regression-based approach attempts to detect boxes through direct regression [20], [22], [23]. The regression-based approaches, such as YoLo algorithm [20], were designed for real-time applications (e.g., automobile, surveillance vision). The classification-based approaches, such as faster R-CNN [18], are superior to the regression-based detection in terms of accuracy primarily because of the fused, local classification and regression procedures. Object detection techniques were successfully employed in the instance segmentation frameworks [17], [21], [24]. The representative method is Mask-RCNN [17], which adapted the faster R-CNN [18] framework by extending the network with pixel-wise classification loss at the end (i.e., instance segmentation framework).

### C. CNNs in Medical Image Segmentation

A fully CNN is the most successful method in medical image segmentation tasks [25]–[35]. Especially, a full 3D-CNN-based approaches [27], [28], [31], [34] outperform other

2D-based methods [25] that are based on 3D convolutions. In [26], [28], a full 3D-CNN-based U-net-like architecture was reported to segment volumetric medical images. A dice coefficient loss metric was proposed to overcome the class imbalance issue in [28]. A voxelwise residual network, called VoxResNet, performed brain tissue segmentation [27]. A residual learning mechanism [36] was used to classify each voxel in the VoxResNet [27]. More recently, training through a shape-prior method was proposed in a densely connected V-net-like structure [35].

## III. METHODOLOGY

The proposed individual tooth segmentation method comprises of three steps: VOI realignment, tooth detection, and individual segmentation. The VOI realignment is primarily performed by a pose regression in a projected 2D space (i.e., x-axis-directional projection). The realignment of VOI not only crops the interested region but also aligns the primary axes of the teeth to be better aligned to the image axes. Subsequently, tooth detection and segmentation are performed using CNNs. The details of the algorithm are described in the following subsections.

### A. Pose-Aware Volume-of-Interest Realignment

In the first step, we reduced the dimension of the original 3D image to 2D for a robust detection of the VOI region,

as proposed in [37]. Let us consider a computed tomography (CT) image  $I(\mathbf{x}) \rightarrow \mathbb{R}$  where  $\mathbf{x} \in \Omega$  ( $\Omega \in \mathbb{R}^3$ ). A maximum intensity projection image corresponding to the x-axis direction,  $I_p$ , is generated from  $I$  (Fig. 3a). The image is then normalized to the range of [0-1]. Finally, we used the trained CNN model to acquire the corresponding point and line pairs. Figure 3a shows an example of a single point and line pair for the upper jaw. After the regression of a point and line pair, we cropped the original CT image by a fixed depth value of 12mm corresponding to the line regression (Fig. 3a) to include the entire root area. We applied an additional 2mm of margin in the opposite direction to compensate for potential errors. In the case of the lower jaw, we flipped the realigned image to synchronize the directions of all the teeth to ease the complexity of shape variations.

For training, we used 100 CT images. We manually annotated (i.e., a point and the angle pairs of the lines) images for  $I_p$ . The overall loss is formulated as follows.

$$\chi(X, \mathbf{p}_i, \theta_i; W) = \sum_{i=0}^2 \|\mathbf{p}_i - \mathbf{y}_i\|_2 + \alpha \sum_{i=0}^2 \|\theta_i - \phi_i\|_2 + \beta \|W\|_2^2, \quad (1)$$

where  $X$ ,  $\mathbf{p}_i$ , and  $\theta_i$  are the input 2D image, ground-truth 2D points, and angles of the lines, respectively.  $W$  represented the weights of the network,  $\mathbf{y}_i$  and  $\phi_i$ , as the network outputs. The network is trained according to the weighting parameters  $\alpha$  and  $\beta$ . For training and inference, we used the traditional VGG-16 model developed by the Visual Geometry Group [38] with a minor modification in the final layer to output a 6D tensor (i.e., 2D point and angle pairs). ‘‘Xavier’’ initialization [39] was used for initializing all the weights of the network. While training the network, we fixed the loss parameters as  $\alpha = \beta = 0.1$ . We used the Adam optimizer [40] with a batch size of 32 and set the learning rate to 0.001. We decayed the learning rate by multiplying 0.1 for every 20 epochs. We trained the network for 100 epochs using an Intel i7-7700K desktop system with 4.2 GHz processor, 32 GB of memory, and Nvidia Titan XP GPU machine. It took 1h to complete all the training procedures.

## B. Tooth Detection

Once VOI was realigned, we performed individual tooth detection. We considered a cropped and realigned image,  $I_{\mathbf{p},\theta,d}$ , where  $\mathbf{p}, \theta$  are a point and an angle regressed by a pose regression step, respectively, and  $d$  indicates the depth of the VOI region with respect to the pose aligning vector (Fig. 3a). We employed a depth of 14mm to include all the teeth, based on to the typical size of the anatomical human tooth presented in [41]. We adopted the faster R-CNN framework [18] as a baseline and applied a few modifications: 1) sampling of the anchors for training the classifier in the region proposal network (RPN) by mining true samples through non-maximum suppression (NMS) and 2) transforming the single-class classification task to a multiclass by anatomical grouping. The former anchor classification in the RPN module is critical for the initial box proposals. To resolve the hard example mining

problem while training the classifier, we applied the NMS technique while training the RPN module. That is, NMS was applied both in the training procedure of the RPN module and in the localizing step for the final output. In the latter grouped classification, we transformed a given single-class classification problem into a multiclass form by grouping the teeth based on the anatomical shapes. We used three classes, i.e., metal, canines/premolars (i.e., one rooted; the identified numbers were 11-13, 21-23, 31-33, 41-43), and the others (i.e., two or more rooted). The group-wise weighted classification aided the RPN proposals of metal-teeth, and thus, improved the final accuracy.

Finally, we applied a 2mm margin, i.e., dilation, to each axis of the output boxes to compensate for a possible inaccurate detection. We resized the realigned image to  $224 \times 224 \times 112$  for all inputs of the network. The performance of the tooth detector can be improved significantly by image realignment and VOI cropping to reduce the overlapping ratio of an object. It is clear that reducing the inter-overlapping area boosts the NMS performance for true example mining (i.e., sampling), which leads to accurate region detection. Moreover, the tooth-to-volume ratio significantly increased through VOI realignment; thus, our proposed framework runs without a patch-wise input cropping procedure [11], which was previously proposed to resolve the problems that arise from anchor-based methods for small objects. The comparative experiments and ablation studies are described in Section IV.

## C. Individual Tooth Segmentation

The individual tooth segmentation was performed by a single CNN. We adopted the base architecture of the 3D U-net [26] which is a popular network for medical image segmentation. The proposed network has three significant differences from the 3D U-net [26], i.e.,: 1) the replacement of a single convolution layer by a skip-connected block that is parameter-efficient, 2) the employment of cutout augmentation [42], [43], and 3) the modification of the final loss function (i.e., voxel-wise classification) to a distance map regression. The former modification was for the efficiency of the network and the latter two modifications were designed to overcome the presence of severe metal artifacts and inter-tooth proximity on CBCT images, respectively.

1) *Base Architecture*: Our proposed network (Fig. 4) architecture has an analysis path and a synthesis path, as in the standard U-net [26], each with four resolution steps. The down-sampling of the dimension was performed by a  $2^3$  max pooling layer with strides of two, and the up-sampling by a  $2^3$  transposed convolution (i.e., de-convolution) with strides of two. The network was implemented similar to the 3D U-net [26], except for the use of SkipBlocks (Fig. 5). We replaced each single set of the convolution layer, batch normalization, and rectified linear unit with a single SkipBlock. The SkipBlock contained two series of non-linearities and a skip connection. The SkipBlock applied a separable convolution technique [44], which is known to reduce the number of parameters and improve the performance of generalization. We used four separable groups in all the experiments.

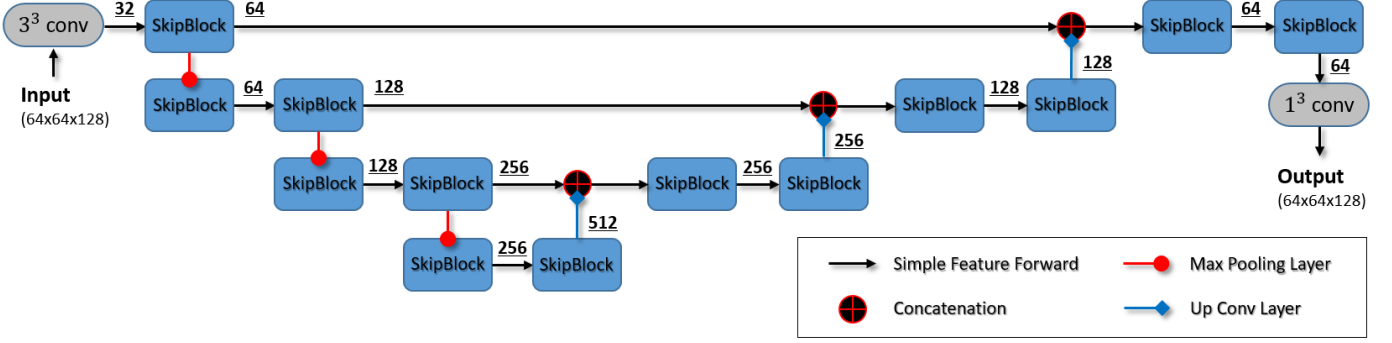


Fig. 4: Proposed 3D volumetric FCN architecture. The network has the same architecture as 3D U-net [26]; however, each nonlinear layer (i.e., convolution layer, batch normalization, and rectified linear unit activation function) is replaced by a SkipBlock (Fig. 5). The blue (squared) and red (circled) arrows indicate the up-convolution and max-pooling layer, respectively, as in the original paper [26]. The numbers demonstrate the output features for each SkipBlock unit.

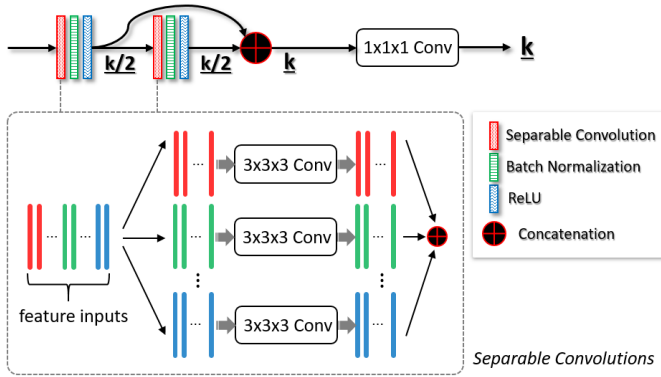


Fig. 5: Skip-connected separable convolution block (i.e., SkipBlock). The input features for each separable convolution are divided into groups. All groups are concatenated after separated convolutions. The two nonlinear layers (i.e., separable convolution, batch normalization, and ReLU activation function) were used in a single block. The output features are produced by a  $1^3$  convolution applied to concatenated features.  $k$  indicates the number of final output features.

2) *Metal-Intensive Augmentation*: We applied the cutout augmentation [42] to improve segmentation results in the metal artifact region. The inference within the metal artifact region, similarly to the inpainting manner, enhanced the final segmentation. The position of the cutout mask was not constrained to the boundaries. A randomly sized zero mask was applied in the range of  $L/5 \leq l \leq L/4$ , where  $l$  and  $L$  are the lengths of the mask and image in each dimension, respectively.

3) *Distance Map Regression Loss*: We employed a distance map to train the network. The ground-truth annotated label was first transformed to a distance map. The distance map is defined by assigning the distance to the closest point in the background to each pixel, as presented below:

$$I_D(\mathbf{i}) = \min_{\mathbf{b} \in B} \text{Dist}(\mathbf{i}, \mathbf{b}), \quad (2)$$

where  $B \subset \mathbf{I}$  is a set of background in the ground-truth label  $\mathbf{I}$ ,  $\text{Dist}$  is a distance function, and  $\mathbf{i} \in \mathbb{R}^3$ . We used

Chamfer distance [45] to approximate the Euclidean distance transformation. The overall loss function is defined by mean squared error (MSE) between the distance map and the final output of the network:

$$L(I_{P,\theta,d}, I_D; W) = \text{MSE}(y, I_D) + \alpha \|W\|_2^2, \quad (3)$$

where  $I_{P,\theta,d}$ ,  $I_D$ ,  $W$  are the cropped tooth image, distanced transformed ground-truth map (2), and the weights of the network, respectively.  $y$  is the final output of the network and  $\alpha$  is a weighting parameter that controls the impact of  $L_2$  regularization.

#### D. Learning the Network

In total, the manual annotation of 50 subjects (CT images) was acquired with the help of clinical experts in the field. In the dataset, the slice thickness values ranged from 0.2-0.4mm, and the pixel sizes ranged from 0.2-0.4mm. For the training dataset, the cropped individual tooth images were resampled into  $64 \times 64 \times 128$ . The input images were normalized in the range [0-1] for each voxel. On-the-fly random affine deformations were subsequently applied to the dataset for each iteration with 80% probability. Finally, the proposed cutout image augmentation was performed with an 80% probability.

“Xavier” initialization [39] was used for initializing all the weights of the proposed network. While training the network, the parameter  $\alpha$  was fixed to 0.1 in (3). The Adam optimizer was used with a batch size of four and a learning rate of 0.001. The learning rate was decayed by multiplying 0.1 for every 50 epochs. The network was trained for 150 epochs using an Intel i7-7700K desktop system with a 4.2 GHz processor, 32 GB memory, and Nvidia Titan XP GPU machine. It took 10 h to complete all the training procedures.

## IV. EXPERIMENTS

In this section, we present the evaluations of the proposed pose-aware (PA) tooth detector and the subsequent individual tooth segmentation. The experimental dataset included metal-intensive cases, which are commonly observed in dental clinics. The CT images were sourced from a multicenter including



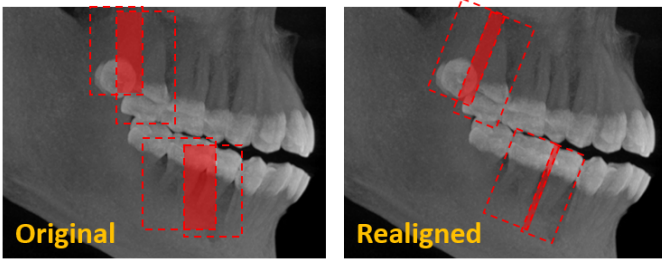


Fig. 6: Sample ground-truth box annotations. The realigned boxes (right) show lower overlaps between boxes.

four different centers. We used 100 independent subjects for training the pose regression network (1), and 50 subjects with each tooth defined for training the individual segmentation network (3). For testing, we used the other 25 subjects for all the quantitative evaluations.

#### A. Performance of Pose-Aware Detector

The overlapping ratio (OR) of the ground-truth boxes was significantly reduced owing to the VOI realignment. The overlapping ratio can be calculated by

$$OR = \frac{A \cap \bar{A}}{A}, \quad (4)$$

where  $A$  indicates the area of the ground-truth box and  $\bar{A}$  indicates the overlapping neighbor boxes. The *mean overlapping ratio* of the ground-truth boxes in the original image alignment was 0.32 while that in the PA VOI realigned image was 0.19. The realignment showed a 38.66% reduction in the overlapping area between boxes. To evaluate the accuracy of the detection, we newly defined a ratio called object include ratio (OIR), which is a ratio (%) of the foreground region of an included object in a detected box to the ground-truth foreground:

$$OIR = \frac{A_o \cap \{A_o \subset \tilde{A}\}}{A_o}, \quad (5)$$

where  $A_o$  indicates the area of an object inside the ground-truth box ( $A$ ) and  $\tilde{A}$  is the detected box. OIR is a critical measurement of whether the individual tooth region is included inside the box that is detected by a neural network. Table I shows the average precision (AP) [46] with threshold value of 0.5 for an intersection over union (i.e.,  $AP_{50}$ ) and the OIR results of the proposed tooth detection (TRCNN), PATRCNN (i.e., pose-aware TRCNN), which employed the realignment, showed significant improvement in the  $AP_{50}$ . By applying a marginal expansion of the boundaries, the OIR reached up to  $0.9951 \pm 0.0023$ , which indicates that the majority of the individual tooth region is included in the detected boxes.

#### B. Segmentation Performance

1) *Evaluation Metric*: The segmentation results were evaluated using the F1 score, aggregated Jaccard index (AJI), precision, sensitivity, Hausdorff distance (HD), and average

TABLE I  
AVERAGE PRECISION AND MEAN  
OBJECT-INCLUDE-RATIO OF THE TOOTH DETECTION

Methods	AP <sub>50</sub>	OIR
TRCNN	67.88	0.83 ± 0.03
PATRCNN	<b>89.75</b>	0.90 ± 0.02
PATRCNN with margin	N/A	<b>1.00 ± 0.00</b>

symmetric surface distance (ASSD). The F1 score is defined as follows:

$$F_1 = 2 \times \frac{Precision \times Sensitivity}{Precision + Sensitivity}. \quad (6)$$

Precision and sensitivity are defined by  $P = \frac{TP}{TP+FP}$  and  $S = \frac{TP}{TP+FN}$ , respectively, where TP, FN, and FP are the numbers of true positive, false negative, and false positive voxels, respectively. The F1 score is equivalent to the dice coefficient [28]. The AJI metric is a per-object metric, as presented in [47]. The one-to-one correspondences were first matched between the ground-truth box and the one detected by maximizing the Jaccard index. Then, the AJI was calculated similar to the Jaccard index; however, the falsely detected components were added to the denominator [48]:

$$AJI = \frac{\sum_i |G_i \cap S_k^*(i)|}{\sum_i |G_i \cup S_k^*(i)| + \sum_{l \in U} |S_l|}, \quad (7)$$

where  $|\cdot|$  is the cardinality of a set.  $G_i$  and  $S_k^*(i)$  are the ground-truth individual object voxels and the object voxels by the corresponding box, which are detected by maximizing the Jaccard index (i.e.,  $S_k^*(i) = \{\arg \max_{S_i} \frac{|G_i \cap S_i|}{|G_i \cup S_i|}\}$ ), respectively, and  $U$  is the set of indices of false detected components that were retained by matching the set,  $S_k^*(i)$ .

The surface distance metrics were evaluated on the integrated teeth basis. Let  $\mathbf{S}_X$  be a set of surface voxels of a set  $X$ ; then, the shortest distance of an arbitrary voxel  $p$  can be defined as follows:

$$d(p, \mathbf{S}_X) = \min_{s_X \in \mathbf{S}_X} \|p - s_X\|_2. \quad (8)$$

Thus, HD is defined as follows [49]:

$$HD(X, Y) = \max\left\{\max_{s_X \in \mathbf{S}_X} d(s_X, \mathbf{S}_Y) + \max_{s_Y \in \mathbf{S}_Y} d(s_Y, \mathbf{S}_X)\right\}. \quad (9)$$

The distance function is defined as:

$$D(\mathbf{S}_X, \mathbf{S}_Y) = \sum_{s_X \in \mathbf{S}_X} d(s_X, \mathbf{S}_Y), \quad (10)$$

Moreover, the ASSD can be defined as follows [49]:

$$ASSD(X, Y) = \frac{1}{|\mathbf{S}_X| + |\mathbf{S}_Y|} (D(\mathbf{S}_X, \mathbf{S}_Y) + D(\mathbf{S}_Y, \mathbf{S}_X)). \quad (11)$$

2) *Comparison*: The overall evaluation of segmentation performances is presented in Table II. We first demonstrate manual seeded (i.e., contoured) level-set-based methods [6], [7] and ToothNet [11], which are state-of-the-art methods for individual tooth segmentation. Our proposed framework,

TABLE II  
QUANTITATIVE RESULTS OF THE STATE-OF-THE-ART TOOTH SEGMENTATION METHODS

Methods	F1 Score	AJI	Precision	Sensitivity	HD [mm]	ASSD [mm]
Seed + Levelset [6]	$0.85 \pm 0.13$	N/A	$0.86 \pm 0.07$	$0.87 \pm 0.16$	$3.19 \pm 2.28$	$0.46 \pm 0.86$
Seed + Levelset [7]	$0.84 \pm 0.13$	N/A	$0.85 \pm 0.08$	$0.87 \pm 0.16$	$3.51 \pm 2.17$	$0.49 \pm 0.85$
ToothNet [11]	$0.88 \pm 0.06$	$0.66 \pm 0.07$	$0.88 \pm 0.08$	$0.89 \pm 0.11$	$2.75 \pm 1.49$	$0.32 \pm 0.14$
Mask-RCNN [17]	$0.83 \pm 0.21$	$0.47 \pm 0.18$	$0.91 \pm 0.06$	$0.82 \pm 0.26$	$3.47 \pm 3.22$	$0.56 \pm 0.86$
TRCNN+TSNet (ours w/o PA)	$0.90 \pm 0.12$	$0.84 \pm 0.00$	$0.91 \pm 0.11$	$0.91 \pm 0.14$	$2.04 \pm 1.84$	$0.29 \pm 0.51$
<b>PATRCNN+TSNet (ours)</b>	<b><math>0.93 \pm 0.03</math></b>	<b><math>0.86 \pm 0.01</math></b>	<b><math>0.93 \pm 0.04</math></b>	<b><math>0.93 \pm 0.07</math></b>	<b><math>1.59 \pm 1.22</math></b>	<b><math>0.20 \pm 0.10</math></b>

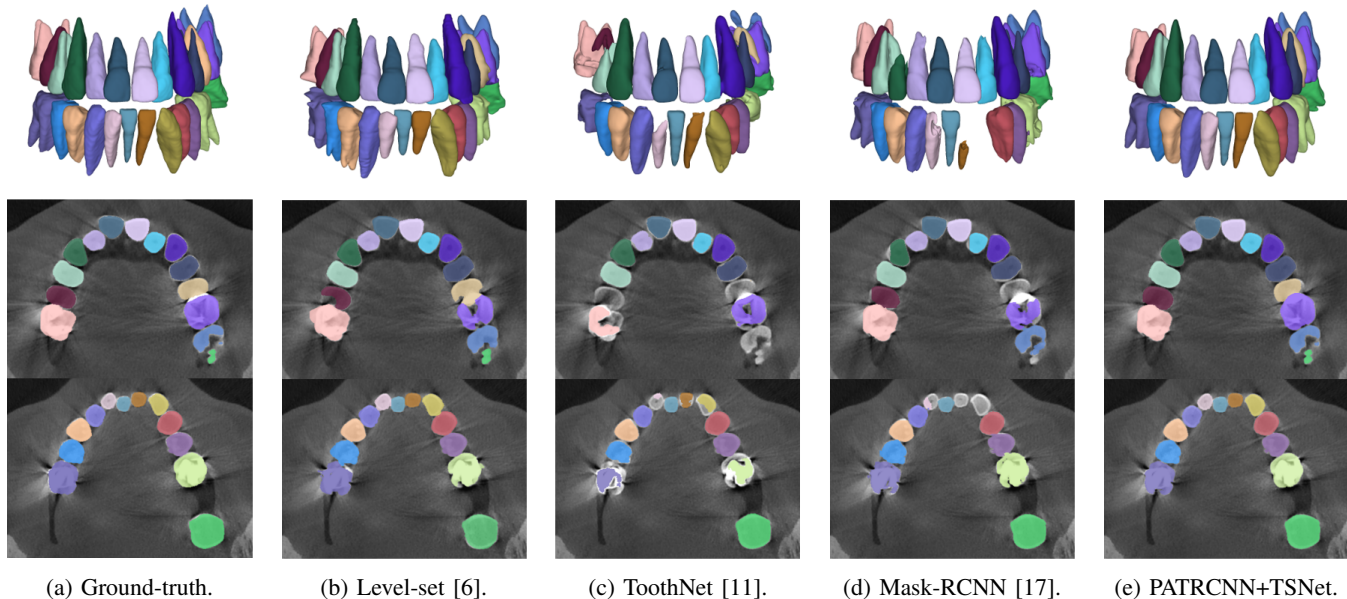


Fig. 7: Visualization of the segmentation results on the test image: (a) First column shows the ground-truth label annotated by the experts in the field. The columns from the second indicate the following: (b) level-set method [6], (c) ToothNet [11], (d) Mask-RCNN [17], and (e) our method, respectively. Note that the level-set method (i.e., the second column) is performed by manually specifying initial contours for each tooth.

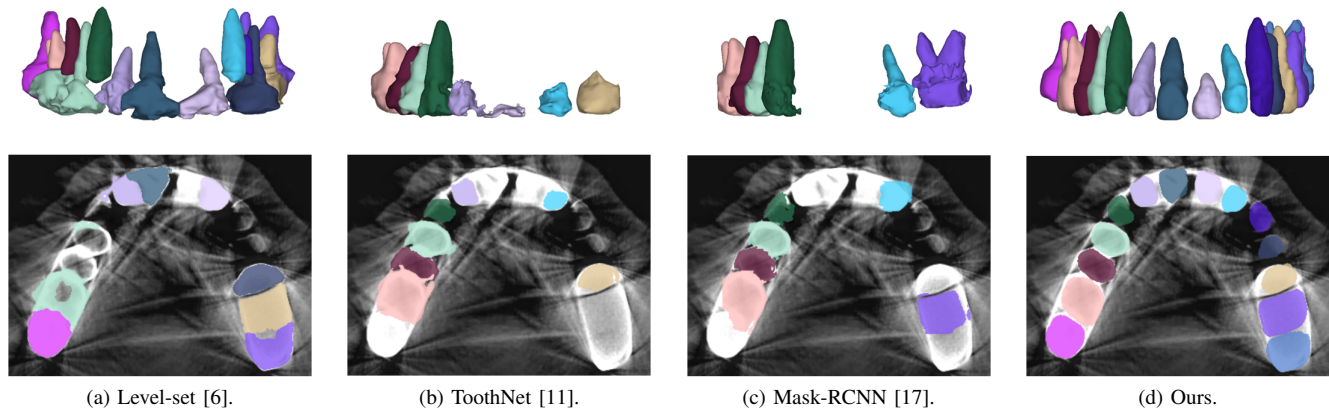


Fig. 8: Visualization of the segmentation results on the test image with severe metal artifacts. The example case shows one of the relatively more severe metal artifacts observed in the clinics. Note that the ground-truth label is not obtained, even by the experts, owing to the difficulty of tooth delineation on metal artifacts. (a) The level-set method [6] suffers from accurate propagation of tooth boundaries. The results of (b) ToothNet [11] and (c) Mask-RCNN [17] demonstrate that many teeth with metal artifacts are difficult to detect. (d) Our proposed PATRCNN+TSNet shows the best detection and segmentation results.

TABLE III  
ASSESSMENT OF SEGMENTATION NETWORKS BASED ON THE GROUND-TRUTH BOXES

Methods	F1 Score	Precision	Sensitivity	HD [mm]	ASSD [mm]
GT+3D U-net [26]	0.89 ± 0.03	0.92 ± 0.05	0.87 ± 0.07	2.71 ± 1.26	0.30 ± 0.11
GT+VoxResNet [27]	0.93 ± 0.02	0.91 ± 0.05	0.96 ± 0.03	1.06 ± 0.60	0.18 ± 0.04
GT+DenseVNet [35]	0.92 ± 0.03	0.97 ± 0.03	0.88 ± 0.07	1.70 ± 1.23	0.21 ± 0.08
GT+TSNet (ours w/o PA)	<b>0.93 ± 0.03</b>	<b>0.97 ± 0.02</b>	<b>0.89 ± 0.06</b>	<b>1.08 ± 0.67</b>	<b>0.18 ± 0.07</b>
PAGT+3D U-net [26]	0.92 ± 0.02	0.92 ± 0.05	0.91 ± 0.05	2.10 ± 1.13	0.23 ± 0.08
PAGT+VoxResNet [27]	0.94 ± 0.02	0.91 ± 0.05	0.97 ± 0.03	0.93 ± 0.54	0.16 ± 0.04
PAGT+DenseVNet [35]	0.93 ± 0.03	0.97 ± 0.02	0.89 ± 0.06	1.25 ± 0.81	0.19 ± 0.07
PAGT+TSNet (ours)	<b>0.96 ± 0.01</b>	<b>0.96 ± 0.03</b>	<b>0.96 ± 0.02</b>	<b>0.86 ± 0.44</b>	<b>0.15 ± 0.04</b>

TABLE IV  
PERFORMANCE OF PA DETECTORS

Methods	AP <sub>50</sub>	OIR
PA-Faster R-CNN [18]	72.20	0.86 ± 0.03
PATRCNN	<b>89.75</b>	<b>0.90 ± 0.02</b>
PATRCNN w/o NMS sampling	76.72	0.89 ± 0.02
PATRCNN w/o grouping	84.59	0.88 ± 0.03

TABLE V  
PERFORMANCE OF TSNet AND ITS VARIANTS

Methods	F1 Score	HD [mm]	ASSD [mm]
PAGT+TSNet	<b>0.96 ± 0.01</b>	<b>0.86 ± 0.44</b>	<b>0.15 ± 0.04</b>
PAGT+TSNet-C	0.94 ± 0.02	0.90 ± 0.48	0.16 ± 0.03
PAGT+TSNet-D	0.91 ± 0.04	0.96 ± 0.51	0.24 ± 0.07
PAGT+TSNet-CD	0.92 ± 0.04	0.96 ± 0.50	0.22 ± 0.07

together with the Mask-RCNN [17], are presented to show the effectiveness of the PA detection and the subsequent tooth segmentation.

The manual seeded level-set methods [6], [7] demonstrate inferior performances while segmenting teeth with predominant metal artifacts (Table II). The energy-based (i.e., level-set function) algorithmic approach failed while defining an optimum stopping criterion (Fig. 7b). The AJI measures for level-set methods are omitted because the box detectors do not exist in the level-set-based methods. In particular, the low precision and sensitivity values indicate that level-set methods have over-segmented or under-segmented the teeth in many cases (Figs. 7b and 8a). The accuracies of the Mask-RCNN [17] and ToothNet [11] showed comparable performance to the level-set-based methods. The instance segmentation framework successfully detected and segmented the teeth automatically. However, it failed to detect all the teeth regions accurately, which resulted in a degradation of the segmentation performance. Moreover, the performance of the segmentation itself also showed low accuracy owing to the metal artifacts (Figs. 7c and 7d). Figure 8 illustrates the more

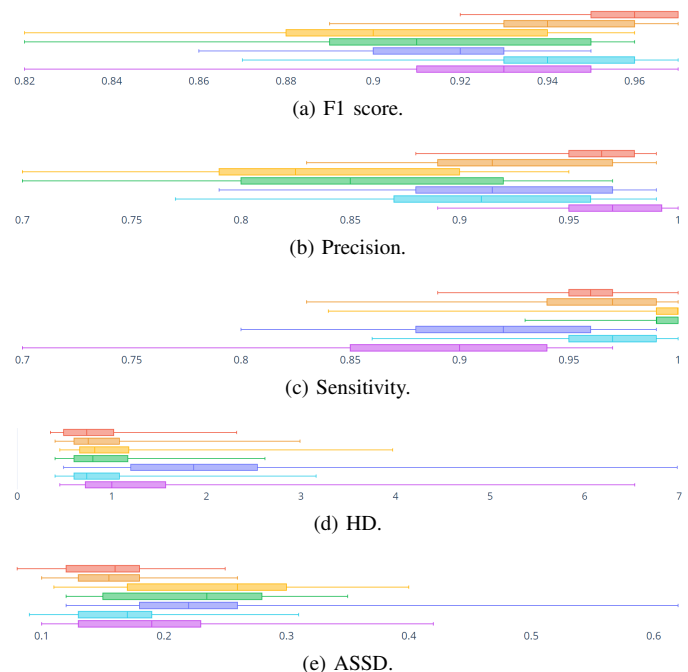


Fig. 9: Box plots of PAGT-based segmentation methods.

severe conditions of metal artifacts. The performance of the detection and segmentation processes significantly degraded in ToothNet [11] (Fig. 8b) and Mask-RCNN [17] (Fig. 8c). Conversely, our proposed method outperformed the other state-of-the-art methods (Figs. 7e and 8d); further, the comparison between the PA- and non-PA-based results also demonstrated that the employment of a PA detector significantly improved the proposed architecture (Table II). The superior AJI value clearly shows that successful detection improved the overall performance.

Table III shows that our proposed TSNet is superior to other networks in a stand-alone segmentation performance (i.e., using the ground-truth boxes). Moreover, the overall performance of the networks that utilized PA-based ground-truth (PAGT) showed better accuracy than the ones that used the original image axes. The original 3D U-net [26] and others failed to segment the teeth with metal artifacts (Figs. 9 and



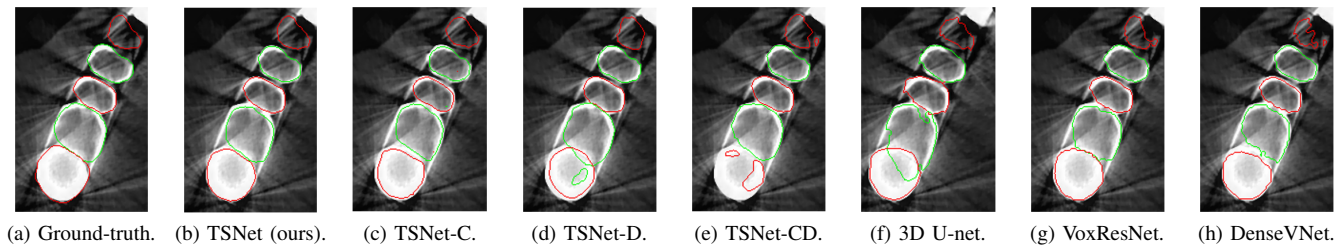


Fig. 10: Visualization of the PAGT-based networks.

10). The DenseVNet [35] showed relatively higher precision than the 3D U-net [26] and VoxResNet [27] as DenseVNet employed a shape-prior based on trainable parameters [35]. The trained shape-prior based on well-bounded images suppressed the false positive responses. However, DenseVNet demonstrated inaccurate segmentation (i.e., sensitivity) due to the large shape variance of teeth.

3) *Ablation Studies*: The performance of our proposed detector (TRCNN) and its variants are listed in Table IV. The PATRCNN without NMS sampling indicates that while training the anchor classifier, we applied a top-k sampling metric with random negative sampling, which is common for an instance segmentation framework. In the case of PATRCNN without grouping, a single-class classification framework was used for all the classifiers in the network. The result demonstrated that the most significant improvement in the accuracy was obtained through true example mining based on the NMS sampling. The multiclass classification metric by grouping also aided accurate detection. The PA-Faster R-CNN method showed a marginal improvement when compared to TRCNN (in Table I), which indicates that the PA method indeed enhances the performance of the detector and neither NMS sampling nor the multiclass classification method can improve the detection without the PA mechanism.

We extended the evaluation of segmentation network with the following variants (Table V): TSNet-C, TSNet-D, and TSNet-CD, which represent without cutout augmentation, without distance loss (3), and without both the cutout and the distance loss (i.e., 3D U-net [26] with SkipBlocks), respectively. We used the dice loss in TSNet-D and TSNet-CD. The result shows that the distance loss metric was the primary factor to achieve success in the individual tooth segmentation. Figure 9 illustrates the box plots of the PAGT-based results and the ablations (Tables III and V) of the proposed network; further, Fig. 10 illustrates that the proposed TSNet is superior to other methods.

## V. DISCUSSION AND CONCLUSION

The presence of metal artifacts in CBCT images, which is prevalent condition in clinical practices, hinders the accurate detection and segmentation of teeth. The proposed instance segmentation architecture overcame the challenges by realigning the VOI, improving the detector, and reformulating the segmentation task (i.e., pixel-wise classification) into a distance map regression. The proposed segmentation network focused on identifying metal artifacts and proximate objects.

Moreover, the fully automated framework neither required manual inputs nor priors. The instance segmentation framework also demonstrated the advantage of avoiding the difficult pose-based procedure of separating individual teeth.

The pose regression-based VOI extraction and realignment method aided the performance of the tooth detector, and also achieved accurate individual tooth segmentation results. The performance of the tooth detector improved by a significant margin owing to two major factors: 1) reduction in the overlapping ratio (4) between ground-truth boxes and 2) formation of anatomical groups for multiclass classification. A reduced overlapping ratio significantly improved the NMS-based sampling in the true example mining stage, which resulted in detector improvement. By employing NMS-based true example sampling, we acquired accurate region proposals without the difficulty of hard example mining. The multiclass classification framework boosted the accuracy of metal-teeth classification, which also resulted in the improvement of tooth detection. Additionally, the VOI realignment enhanced the tooth-to-volume ratio (TVR) of the input image. The enhancement of TVR demonstrated a huge benefit for instance segmentation because the anchor-based RPN module demonstrates a high probability of ignoring small objects. It is more critical to the volumetric input images owing to the limitation of GPU memories. Unlike [11], our proposed network did not require a tedious procedure of cropping the input images while training. The proposed pose regression framework suggests that obtaining a simple cue (i.e., pose) can ease the complexity of the original task in many aspects.

## REFERENCES

- [1] H. Akhoondali, R. Zoroofi, and G. Shirani, "Rapid automatic segmentation and visualization of teeth in ct-scan data," *Journal of Applied Sciences*, vol. 9, no. 11, pp. 2031–2044, 2009.
- [2] S. Kakehbaraei, H. Seyedarabi, and A. T. Zenouz, "Dental segmentation in cone-beam computed tomography images using watershed and morphology operators," *Journal of medical signals and sensors*, vol. 8, no. 2, p. 119, 2018.
- [3] Y. Fan, R. Beare, H. Matthews, P. Schneider, N. Kilpatrick, J. Clement, P. Claes, A. Penington, and C. Adamson, "Marker-based watershed transform method for fully automatic mandibular segmentation from cbct images," *Dentomaxillofacial Radiology*, vol. 48, no. 2, p. 20180261, 2019.
- [4] H. Gao and O. Chae, "Individual tooth segmentation from ct images using level set method with shape and intensity prior," *Pattern Recognition*, vol. 43, no. 7, pp. 2406–2417, 2010.
- [5] D. X. Ji, S. H. Ong, and K. W. C. Foong, "A level-set based approach for anterior teeth segmentation in cone beam computed tomography images," *Computers in biology and medicine*, vol. 50, pp. 116–128, 2014.

- [6] Y. Gan, Z. Xia, J. Xiong, G. Li, and Q. Zhao, "Tooth and alveolar bone segmentation from dental computed tomography images," *IEEE journal of biomedical and health informatics*, vol. 22, no. 1, pp. 196–204, 2017.
- [7] Y. Wang, S. Liu, G. Wang, and Y. Liu, "Accurate tooth segmentation with improved hybrid active contour model," *Physics in Medicine & Biology*, vol. 64, no. 1, p. 015012, 2018.
- [8] S. Barone, A. Paoli, and A. V. Razonale, "Ct segmentation of dental shapes by anatomy-driven reformation imaging and b-spline modelling," *International journal for numerical methods in biomedical engineering*, vol. 32, no. 6, p. e02747, 2016.
- [9] Y. Pei, X. Ai, H. Zha, T. Xu, and G. Ma, "3d exemplar-based random walks for tooth segmentation from cone-beam computed tomography images," *Medical physics*, vol. 43, no. 9, pp. 5040–5050, 2016.
- [10] Y. Miki, C. Muramatsu, T. Hayashi, X. Zhou, T. Hara, A. Katsumata, and H. Fujita, "Classification of teeth in cone-beam ct using deep convolutional neural network," *Computers in biology and medicine*, vol. 80, pp. 24–29, 2017.
- [11] Z. Cui, C. Li, and W. Wang, "Toothnet: Automatic tooth instance segmentation and identification from cone beam ct images," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 6368–6377.
- [12] L. Hiew, S. Ong, and K. W. Foong, "Tooth segmentation from cone-beam ct using graph cut," in *Proceedings of the Second APSIPA Annual Summit and Conference*, 2010, pp. 272–275.
- [13] L. Wang, Y. Gao, F. Shi, G. Li, K.-C. Chen, Z. Tang, J. J. Xia, and D. Shen, "Automated segmentation of dental cbct image with prior-guided sequential random forests," *Medical physics*, vol. 43, no. 1, pp. 336–346, 2016.
- [14] Z. Zou, S. Luo, J.-S. Pan, S. Liu, and S. Liao, "A semi-automatic segmentation for tooth on cone beam ct volume following the anatomic guidance," *Journal of Information Hiding and Multimedia Signal Processing*, vol. 8, pp. 774–783, 2017.
- [15] M. Zhao, L. Ma, W. Tan, and D. Nie, "Interactive tooth segmentation of dental models," in *2005 IEEE Engineering in Medicine and Biology 27th Annual Conference*. IEEE, 2006, pp. 654–657.
- [16] Z. Xia, Y. Gan, L. Chang, J. Xiong, and Q. Zhao, "Individual tooth segmentation from ct images scanned with contacts of maxillary and mandible teeth," *Computer methods and programs in biomedicine*, vol. 138, pp. 1–12, 2017.
- [17] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2961–2969.
- [18] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *Advances in neural information processing systems*, 2015, pp. 91–99.
- [19] R. Girshick, "Fast r-cnn," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1440–1448.
- [20] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 779–788.
- [21] L.-C. Chen, A. Hermans, G. Papandreou, F. Schroff, P. Wang, and H. Adam, "Masklab: Instance segmentation by refining object detection with semantic and direction features," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4013–4022.
- [22] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "Ssd: Single shot multibox detector," in *European conference on computer vision*. Springer, 2016, pp. 21–37.
- [23] J. Redmon and A. Farhadi, "Yolov3: An incremental improvement," *arXiv preprint arXiv:1804.02767*, 2018.
- [24] J. Dai, K. He, and J. Sun, "Instance-aware semantic segmentation via multi-task network cascades," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 3150–3158.
- [25] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical image computing and computer-assisted intervention*. Springer, 2015, pp. 234–241.
- [26] Ö. Çiçek, A. Abdulkadir, S. S. Lienkamp, T. Brox, and O. Ronneberger, "3d u-net: learning dense volumetric segmentation from sparse annotation," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2016, pp. 424–432.
- [27] H. Chen, Q. Dou, L. Yu, J. Qin, and P.-A. Heng, "Voxresnet: Deep voxelwise residual networks for brain segmentation from 3d mr images," *NeuroImage*, 2017.
- [28] F. Milletari, N. Navab, and S.-A. Ahmadi, "V-net: Fully convolutional neural networks for volumetric medical image segmentation," in *3D Vision (3DV), 2016 Fourth International Conference on*. IEEE, 2016, pp. 565–571.
- [29] K. Kamnitsas, C. Ledig, V. F. Newcombe, J. P. Simpson, A. D. Kane, D. K. Menon, D. Rueckert, and B. Glocker, "Efficient multi-scale 3d cnn with fully connected crf for accurate brain lesion segmentation," *Medical image analysis*, vol. 36, pp. 61–78, 2017.
- [30] M. Havaei, A. Davy, D. Warde-Farley, A. Biard, A. Courville, Y. Bengio, C. Pal, P.-M. Jodoin, and H. Larochelle, "Brain tumor segmentation with deep neural networks," *Medical image analysis*, vol. 35, pp. 18–31, 2017.
- [31] Q. Dou, L. Yu, H. Chen, Y. Jin, X. Yang, J. Qin, and P.-A. Heng, "3d deeply supervised network for automated segmentation of volumetric medical images," *Medical image analysis*, vol. 41, pp. 40–54, 2017.
- [32] S. Jégou, M. Drozdal, D. Vazquez, A. Romero, and Y. Bengio, "The one hundred layers tiramisu: Fully convolutional densenets for semantic segmentation," in *Computer Vision and Pattern Recognition Workshops (CVPRW), 2017 IEEE Conference on*. IEEE, 2017, pp. 1175–1183.
- [33] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "DeepLab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs," *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 4, pp. 834–848, 2018.
- [34] O. Oktay, E. Ferrante, K. Kamnitsas, M. Heinrich, W. Bai, J. Caballero, S. A. Cook, A. de Marvao, T. Dawes, D. P. O'Regan *et al.*, "Anatomically constrained neural networks (acnns): application to cardiac image enhancement and segmentation," *IEEE transactions on medical imaging*, vol. 37, no. 2, pp. 384–395, 2018.
- [35] E. Gibson, F. Giganti, Y. Hu, E. Bonmati, S. Bandula, K. Gurusamy, B. Davidson, S. P. Pereira, M. J. Clarkson, and D. C. Barratt, "Automatic multi-organ segmentation on abdominal ct with dense v-networks," *IEEE Transactions on Medical Imaging*, 2018.
- [36] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [37] M. Chung, J. Lee, W. Song, Y. Song, I.-H. Yang, J. Lee, and Y.-G. Shin, "Automatic registration between cone-beam ct and scanned surface via deep-pose regression neural networks and clustered similarities," *arXiv preprint arXiv:1907.12250*, 2019.
- [38] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [39] X. Glorot and Y. Bengio, "Understanding the difficulty of training deep feedforward neural networks," in *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, 2010, pp. 249–256.
- [40] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [41] S. J. Nelson, *Wheeler's Dental Anatomy, Physiology and Occlusion-E-Book*. Elsevier Health Sciences, 2014.
- [42] T. DeVries and G. W. Taylor, "Improved regularization of convolutional neural networks with cutout," *arXiv preprint arXiv:1708.04552*, 2017.
- [43] Z. Zhong, L. Zheng, G. Kang, S. Li, and Y. Yang, "Random erasing data augmentation," *arXiv preprint arXiv:1708.04896*, 2017.
- [44] F. Chollet, "Xception: Deep learning with depthwise separable convolutions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1251–1258.
- [45] G. Borgefors, "Distance transformations in digital images," *Computer vision, graphics, and image processing*, vol. 34, no. 3, pp. 344–371, 1986.
- [46] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *European conference on computer vision*. Springer, 2014, pp. 740–755.
- [47] N. Kumar, R. Verma, S. Sharma, S. Bhargava, A. Vahadane, and A. Sethi, "A dataset and a technique for generalized nuclear segmentation for computational pathology," *IEEE transactions on medical imaging*, vol. 36, no. 7, pp. 1550–1560, 2017.
- [48] P. Naylor, M. Laé, F. Reyat, and T. Walter, "Segmentation of nuclei in histopathology images by deep regression of the distance map," *IEEE transactions on medical imaging*, vol. 38, no. 2, pp. 448–459, 2018.
- [49] T. Heimann, B. Van Ginneken, M. A. Styner, Y. Arzhaeva, V. Aurich, C. Bauer, A. Beck, C. Becker, R. Beichel, G. Bekes *et al.*, "Comparison and evaluation of methods for liver segmentation from ct datasets," *IEEE transactions on medical imaging*, vol. 28, no. 8, pp. 1251–1265, 2009.