

1 Automatic hyoid bone detection in fluoroscopic images using deep
2 learning

3 Zhenwei Zhang¹ James L. Coyle² Ervin Sejdić¹

4
5 **Abstract**

6 The displacement of the hyoid bone is one of the key components evaluated in the swallow
7 study, as its motion during swallowing is related to overall swallowing integrity. In daily research
8 settings, experts visually detect the hyoid bone in the video frames and manually plot hyoid
9 bone position frame by frame. This study aims to develop an automatic method to localize
10 the location of the hyoid bone in the video sequence. To automatically detect the location of
11 the hyoid bone in a frame, we proposed a single shot multibox detector, a deep convolutional
12 neural network, which is employed to detect and classify the location of the hyoid bone. We
13 also evaluated the performance of two other state-of-art detection methods for comparison. The
14 experimental results clearly showed that the single shot multibox detector can detect the hyoid
15 bone with an average precision of 89.14 % and outperform other auto-detection algorithms.
16 We conclude that this automatic hyoid bone tracking system is accurate enough to be widely
17 applied as a pre-processing step for image processing in dysphagia research, as well as a promising
18 development that may be useful in the diagnosis of dysphagia.

19 **Keywords:** hyoid bone, dysphagia, deep learning, computer vision, videofluoroscopy

20 Dysphagia, a common condition among older individuals, is defined as an impairment in swal-
21 lowing function during eating and drinking [1]. Dysphagia causes subjective discomfort and objec-
22 tive difficulty in the formation or transportation of a bolus from mouth to stomach, and prevention

¹Zhenwei Zhang and Ervin Sejdić are with the Department of Electrical and Computer Engineering, Swanson School of Engineering, University of Pittsburgh, Pittsburgh, PA, 15261, USA. E-mail: esejdic@ieee.org. Ervin Sejdić is the corresponding author.

²James L. Coyle are with Department of Communication Science and Disorders, University of Pittsburgh, Pittsburgh, PA, 15260, USA.

23 of errant entry of swallowed material into the airway. Dysphagia is a frequent clinical sign in
24 patients with stroke, head and neck cancer and a variety of other medical conditions [2–4]. The
25 prevalence of dysphagia is very high: stroke is the most commonly reported etiology with over 50
26 % of patients exhibiting dysphagia in the immediate post-onset stage of recovery, diminishing to
27 a lower prevalence of around 11 % within 6 months of onset [5]. Additionally, chronic dysphagia
28 affects 7.2% of people with other neurological diseases, and 4.9% of patients treated for head and
29 neck cancer [6]. Up to 40% of people over 65 years old and more than 60% of adults in nursing
30 home [7] suffer from dysphagia. It is estimated that 25% – 50% of Americans over 60 [2] and
31 17% of citizens over 65 in Europe [8] will suffer from dysphagia, leading to increased risk of poor
32 nutrition or dehydration. The variation in estimation may be due to different definitions of dys-
33 phagia, the method of swallowing assessment and the number of patients investigated. As a more
34 immediate clinical consequence, dysphagia may lead to misdirection of food and colonized saliva
35 into the airway, possibly causing pneumonia and chronic lung disease. In many cases aspiration
36 occurs without any obvious clinical signs of dysphagia (silent aspiration), postponing early iden-
37 tification and preventive treatment therefore lowering patient survival [9]. Efforts to accurately
38 evaluate swallowing function early after the onset of conditions leading to dysphagia can mitigate
39 many of these health risks [10].

40 The videofluoroscopic swallowing studies (VFSS), also known as modified barium swallow study,
41 is the gold standard test for dysphagia evaluation [11–14]. VFSS, unlike bedside clinical examina-
42 tion, enables the examiner to visualize oral, pharyngeal and upper esophageal structure and function
43 during patient swallowing. VFSS also evaluate errors of biomechanical coordination that lead to
44 bolus misdirection. Patients with dysphagia may not exhibit overt signs of swallowing problems
45 at the bedside. VFSS excels at allowing clinicians to identify occult disorders in airway protection
46 and biomechanical errors leading to impaired airway protection and transfer of food to the diges-
47 tive system. Airway closure and upper esophageal sphincter opening are largely influenced by the
48 timing and displacement of the hyolaryngeal complex during the pharyngeal stage of swallowing.
49 During VFSS, the hyoid bone is the most salient anatomic structure for detecting hyolaryngeal mo-
50 tion [15]. Hyolaryngeal excursion is an important feature considered by clinicians and researchers
51 because disordered motion may signify dysphagia. Clinicians make subjective judgments about
52 the completeness of hyoid displacement by gross visual inspection of VFSS images. In dysphagia
53 research labs, expert judges annotate hyoid position and its key components in each image frame.
54 However, the subjective clinical process is prone to judgment error, and frame-by-frame annotation
55 done by researchers is time consuming and is prone to inter- and intra-rater variation [16].

56 Efforts by researchers to develop hyoid tracking methods that combine human judgment with
57 automated image processing and machine learning are still quite limited. Patrick *et al.* proposed a
58 method to define the hyoid bone in a calibration frame by identifying a region of interest manually
59 and using Sobel edge detection to track the hyoid bone in subsequent frames [17]. Hoaasin *et*
60 *al.* proposed a semi-automatic hyoid bone tracking system that can match the hyoid bone by
61 Haar classifier matching. However, their method still requires manual identification of regions that
62 clearly contain the hyoid bone [18]. Lee *et al.* developed a software platform that extracted the
63 trajectory of the moving hyoid bone by calculating local binary patterns and multi-scale local binary
64 patterns [19]. Kim *et al.* developed software which can track, smooth and segment the hyoid bone
65 motion from VFSS [20].

66 Remarkable progress has been made in medical imaging techniques due to the large number of
67 databases and deep convolutional neural networks (CNNs) [21, 22]. Currently, the ideas of CNNs
68 are mainly employed in various medical imaging modalities such as conventional X-ray fluoroscopy,
69 MRI and CT for classification and segmentation [23–26]. The medical applications of CNNs tech-
70 niques are to help clinicians diagnose and classify diseases more quickly, including segmentation of
71 various tissues such as brain and organs; classification of cancer, fractures, neurological diseases
72 and biomedical image retrieval systems. Research based on segmentation and object detection has
73 closely followed the development of CNNs in the last few years. Almost all recent works for the
74 object detectors and segmentation are based on CNNs, a deep architecture using pretraining on
75 ImageNet which is trainable end-to-end. Girshick *et al.* first described Region-based Convolutional
76 Neural Networks (RCNN) that dramatically increased the performance of object detection com-
77 pared to traditional features based classifiers [27]. Traditional methods usually use sliding windows
78 for region proposal, histograms of gradient orientation (HoG) or scale-invariant feature transform
79 (SIFT) as feature extraction [28, 29], and support vector machine (SVM) and Boosting methods
80 as classifiers [30, 31]. Fast-RCNN extended the idea of RCNN and improved system performance
81 by sharing the computation across the proposed image regions [32]. Then, Faster-RCNN improved
82 the region proposer method and sped up the overall process [33]. In this method, only one CNN is
83 trained and the region proposal reused the results of the same CNN instead of running a separate
84 searching algorithm in the previous work. You Only Look Once (YOLO) [34] and Single Shot
85 MultiBox Detector (SSD) [35] are existing methods that focus on better computation speed and
86 performance. These two methods classify and regress a set of anchor boxes without using the idea
87 of Regions of Interests. YOLO applies a simpler network structure, predicting bounding boxes
88 and class probabilities directly from the last convolutional feature maps. SSD uses features from

89 different layers progressively to predict the various size of bounding boxes. Features from the early
90 layers were applied to predict the small-sized boxes while features from the latter layers are applied
91 for larger boxes.

92 In previous research related to the hyoid bone motion, users manually marked a region of
93 interest in the first frame after which their algorithm tracked or detected the motion of hyoid bone.
94 The number of images used in these studies was not representative of a patient population. The
95 hyoid bone motion analysis provides meaningful solutions in clinical research settings. However,
96 the manual tracking is time consuming and impractical in real-life cases. Improved hyoid bone
97 localization and an automatic hyoid bone tracking system can help clinicians provide a quicker
98 assessment of the patient. Therefore, we sought to develop a software platform that can localize
99 the region of interest containing the hyoid bone in subsequent video frames. The proposed method
100 relies on the CNN based object detection method. We hypothesized that our detection algorithms
101 would accurately detect the location of the hyoid bone in each video frame with high accuracy when
102 compared to the gold-standard manual detection method (visual inspection with frame-by-frame
103 plotting).

104 The paper is organized as follows. Section 2 reports the background and the current state-of-
105 the-art object detection methods; section 3 proposes the methodology, followed by the analysis of
106 the experimental results and discussion; and section 4 concludes the paper.

107 **1 Material and Methods**

108 **1.1 Data Collection**

109 In this investigation, 265 patients with swallowing difficulty underwent videofluoroscopic exam-
110 ination at the Presbyterian University Hospital of the University of Pittsburgh Medical Center
111 (Pittsburgh, Pennsylvania). The protocol for this study was approved by the Institutional Review
112 Board at the University of Pittsburgh and all participants provided informed consent. The age
113 range of these subjects was from 19 to 94, and the average age was 64.833 ± 13.56 years old. The
114 distribution of ages is illustrated in Fig 1. There were no significant differences in hyoid bones be-
115 tween younger and older patients in the detection task. The main difference in the anatomy of the
116 hyoid bone across the lifespan is density and when the greater cornua fuses to the body of the hyoid.
117 Hyoid bone tracking with VFSS relies on identification of landmarks on the body of the hyoid bone
118 without regard to cornua. Patients swallowed radiopaque liquid boluses of different consistencies

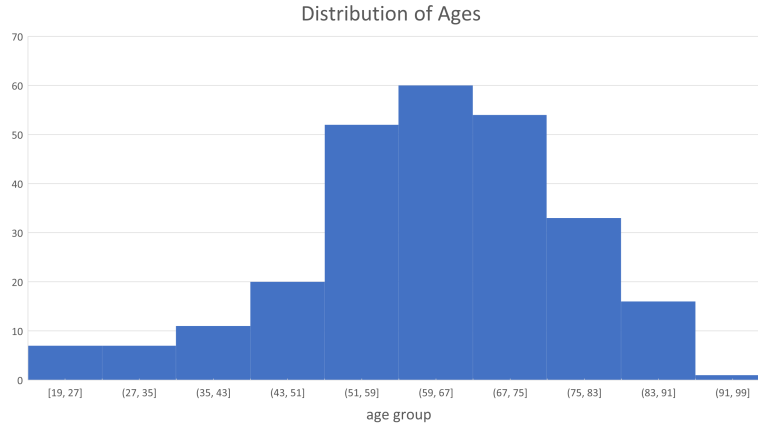


Figure 1: The age range of participants are from 19 to 94. Most of subjects are in the age range 43-83 years old.

119 and volumes as well as pureed food and cookies during their VFSS examination. The volumes
 120 and viscosity of material administered to patients were determined during the examinations in real
 121 time by clinicians based on factors such as the patient’s history and clinical indications. These
 122 liquids included thin liquid (Varibar Thin Liquid with < 5 cPs viscosity), and nectar-thick liquid
 123 (Varibar Nectar with about 300 cPs viscosity). The position of patients during swallowing was pri-
 124 marily neutral head position though some swallows were performed in a head-neck flexion position.
 125 Patients swallowed liquid boluses from a spoon containing 3-5mL volumes, or self-administered
 126 boluses from a cup containing patient self-selected, comfortable volumes between 10-20mL.

127 Videofluoroscopy was set at 30 pulses per second (full motion) and video images were acquired
 128 at 60 frames per second by a video card (AccuStream Express HD, Foresight Imaging, Chelmsford,
 129 MA) and collected into a hard drive with a LabVIEW program. The videos were two-dimensional
 130 digital movie clips of 720 x 1080 resolution, and in this investigation, we down-sampled the video
 131 clips to 30 frames/second to eliminate duplicated frames.

132 1.2 Methods

133 In this investigation, our solution is to build a detection system based on the single shot multibox
 134 detector, which is one of the most popular detection algorithms in recent years. The SSD algorithm
 135 can generate high detection performance at the cost of high computational complexity. Thus, we
 136 also evaluated the performance of several other state-of-the-art detection methods, i.e., Faster-
 137 RCNN and YOLOv2, for comparison. The following paragraphs describe the SSD approach, the

138 data set ground truth creation and the training and testing details.

139 1.2.1 Network Architecture

140 Machine learning has been widely used in medical imaging and videos to help users better un-
141 derstand the properties of these data [36]. Neural networks are one of the most popular types of
142 machine learning models. The basic idea of neural networks is to multiply the input data with layers
143 of weighted connections. Deep neural networks consist of a typical architecture of neural networks,
144 constructed by multiple layers. Each layer implements a series of convolution operators on input,
145 followed by a non-linear activation function, such as a logistic function or a rectified linear unit
146 (Relu). Then a pooling layer is applied to reduce the size of features to the following layers [37].
147 Popular convolutional neural networks for image tasks include AlexNet [38], GoogleNet [39], VGG
148 net [40] and Residual Net [41].

149 The SSD is a feed-forward convolutional neural network built on image classification neural net-
150 work, called base network, such as VGGNet, ZFNet or ResNet [35]. Eight additional convolutional
151 feature layers are added after these base networks to replace the last few layers of the base networks.
152 The size of these layers decreased progressively and were used as output layers for the prediction
153 of detections at multiple resolutions. SSD integrated both higher and lower feature layers, as the
154 lower layers contain better location information and the higher layers have more image details [42].
155 The images are divided into different grid sizes which are associated to default bounding boxes.
156 The correspondence between the position of the default box and the feature cell are fixed. SSD
157 predicts the objects based on default boxes instead of predicting the bounding boxes directly. The
158 default boxes are assigned with different scales and aspect ratios, which provides information on
159 different object scales. The scale of each feature map is manually designed as:

$$s_k = s_{min} + \frac{s_{max} - s_{min}}{m - 1}(k - 1), \quad k \in [1, m]$$

160 where m is the number of feature maps used for prediction. s_{min} is 0.2 and s_{max} is 0.9.

161 Each feature map cell is correspondent to 6 default boxes, which are assigned with different
162 aspect ratios, denoted as $\alpha_\gamma = \{1, 2, 3, \frac{1}{2}, \frac{1}{3}\}$. The width and height of the default box is computed
163 as $w_k^\alpha = s_k \sqrt{\alpha_\gamma}$ and $h_k^\alpha = s_k / \sqrt{\alpha_\gamma}$. For the aspect ratio of 1, another scale $s'_k = \sqrt{s_k s_{k+1}}$ is added
164 for the default box as well. The center of each default box is set at $(\frac{i+0.5}{|f_k|}, \frac{j+0.5}{|f_k|})$, and $|f_k|$ is the
165 size of k-th feature map. By using these default boxes with various scales and aspect ratios from
166 all locations of added feature maps, SSD predictions can cover different input sizes and shapes. Fig

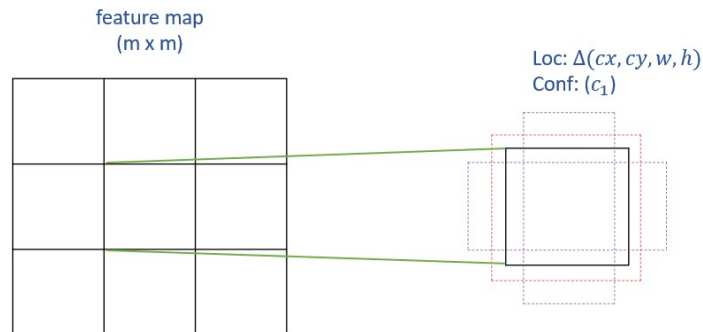


Figure 2: The idea of default boxes applied in SSD. For each default box, the offsets and confidence for categories are predicted.

167 2 illustrates the idea of default boxes.

168 A set of convolutional filters are applied to the added features layers to perform the bounding
 169 box regression and category classification. For each feature layer of size $m \times n$ with p channels,
 170 a $3 \times 3 \times p$ small kernel filter is applied to produce one value at each feature map cell, where the
 171 outputs are classification scores as well as the offsets relative to the bounding box shape.

172 The label of SSD includes the class and the offsets from the default boxes. The default boxes
 173 are matched with ground truth if their intersection over union (IOU) is over 0.5. IOU is defined
 174 as *Area of Overlap/Area of Union*. The loss function of SSD combines a softmax loss for the
 175 confidence loss and a Smooth L1 loss for localization loss. The overall objective loss function is

$$L_{tot} = \frac{1}{N}(L_{conf} + \alpha L_{loc})$$

176 where N is the number of matched default boxes and α is set to 1 by cross-validation. The SSD
 177 framework is shown in Fig 3. For more details of the SSD network and loss function please refer
 178 to [35].

179 1.2.2 Training and Testing

180 Expert judges in VFSS image measurements manually annotated the hyoid bone location (coordi-
 181 nate of left corner, height and width) in each frame of the videos. To evaluate the reliability of the
 182 swallowing analysis, 10 swallow cases were utilized. Three experts analyzed the same 10 swallows.
 183 Inter-rater reliability was tested between raters and experts analyzed the same cases one month

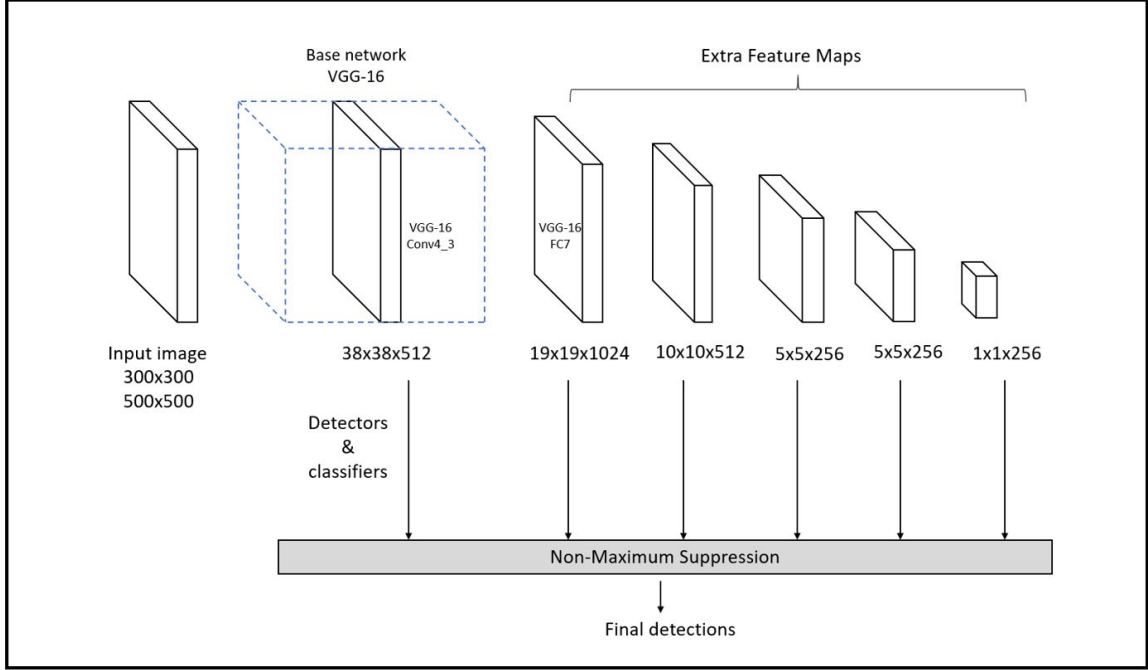


Figure 3: Architecture of Single shot multibox detector

184 later for intra-rater reliability. ICC score were over 0.9 for all measures of reliability. The swallow
 185 data were split and distributed to each of the experts. Their annotations were considered as ground
 186 truth (gold standard). The data were randomly separated by patients: 70 % of the patients were
 187 split into training data which contained around 30,000 frames with annotations, while 30 % of the
 188 patients were split into test data which contained around 18,000 frames. We chose both VGG-16
 189 and ResNet-101 as base networks, and considered two image resolutions for inputs: 300×300 and
 190 500×500 . We compared models trained on both base networks and both resolutions inputs. The
 191 input with size 500×500 should provide better performance as more details can be detected in
 192 higher resolution images. However, a larger image size increases the computation complexity. Fur-
 193 thermore, we compared the results with YOLO and Faster-RCNN and used a training procedure
 194 similar to the original papers. We chose 0.0005 as our learning rate with multi-steps, dividing by
 195 10 for iteration 4000 and 8000. The momentum is 0.9 and gamma is 0.1 for the optimizers.

196 1.2.3 Evaluation of Accuracy

197 The performance of the detection module is measured by mean average precision (mAP), which
 198 is the most commonly used evaluation method for object detection. Average precision estimated

Table 1: Comparison of mAP with different models

| Model | mean average precision |
|------------------|------------------------|
| YOLOv2 | 33.10% |
| Faster-RCNN + ZF | 69.01% |
| SSD300-VGG | 84.37% |
| SSD300-ResNet | 81.92% |
| SSD500-VGG | 89.14% |
| SSD500-ResNet | 89.03% |

whether detected bounding boxes match the corresponding ground truth. Mean average precision is the area below the precision-recall curve, which integrates precision and recall while varying from 0 to 1. As we have just one class to classify, mean average precision is the average precision for the hyoid bone class. The bounding box is labeled as true positive if IOU is greater than 0.5. Precision evaluates the fraction of true positive bounding box over all predictions and recall evaluates the fraction of the true positive detected bounding boxes among all ground truths.

2 Results

Table 1 shows results of the state-of-the-art published methods on our VFSS image dataset. Overall, SSD method outperforms the results produced by YOLOv2 and Faster-RCNN. Among SSD method, VGGNet with input size of 500×500 produced the best result compared to ResNet and input size of 300×300 . The mAP of SSD500-VGGNet is 89.14%, which is 0.11% better than using ResNet-101 as base network and 2.45% better than using the smaller image input size. Figure 3 shows the example results by manual segmentation, SSD500-VGGNet, Faster-RCNN and YOLOv2. We selected two different cases as examples: patient swallowing the bolus in neutral head position or in chin down position. In comparing automated hyoid detection to the ground truth, we used the bounding box to locate the hyoid bone. Most of the object detection methods use the bounding box to locate and classify the content inside. In the example case, all three tested methods revealed a positive result, detecting the hyoid bone location successfully. However, the Faster-RCNN method produced two regions of interest that it considered as the hyoid bone with a close confidence score.

Figure 5 illustrates results using the SSD500-VGGNet method with different hyoid bone locations (under the mandible and behind the mandible), and the results are shown with different image

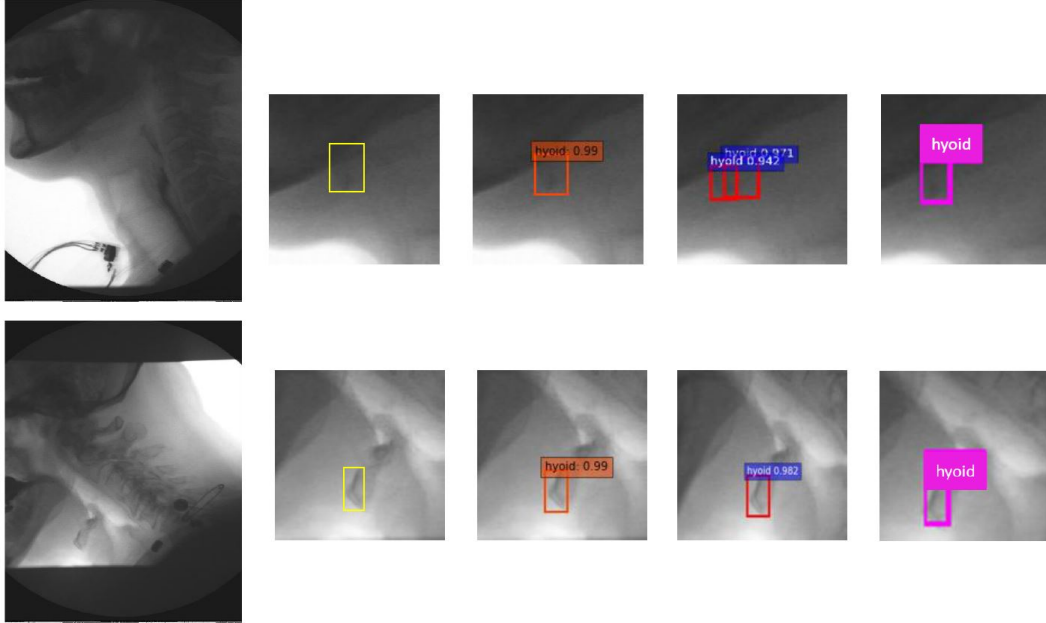


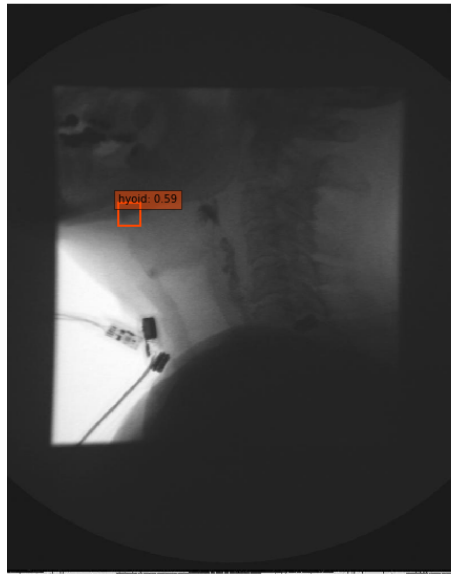
Figure 4: The identification of hyoid bone using different method: ground truth (yellow), SSD500-VGG (orange), Faster-RCNN (red), and YOLOv2 (pink)

220 qualities. From these results, SSD500-VGGNet showed stable detection results, clearly finding the
 221 hyoid bone. When the hyoid bone is hidden behind the mandible in case (a) and (b), the algorithm
 222 detected the hyoid bone with a relatively low confidence score. It performed well in case (c) and
 223 (d) where the hyoid bone is present under the mandible.

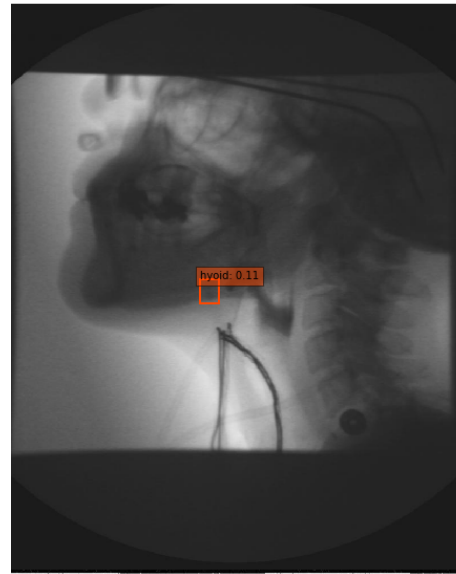
224 Figure 6 shows the change of training loss function and the performance on test data during the
 225 training of SSD models. These figures illustrate how the performance of the model changes during
 226 training. The loss function dramatically decreased in the first 1000 iterations and the loss function
 227 only slightly decreased in the following training iterations. The training errors of SSD300-VGG
 228 were always higher than those of SSD500-VGG. SSD300 with different pre-trained models showed
 229 a similar training loss trend and test accuracy.

230 3 Discussion

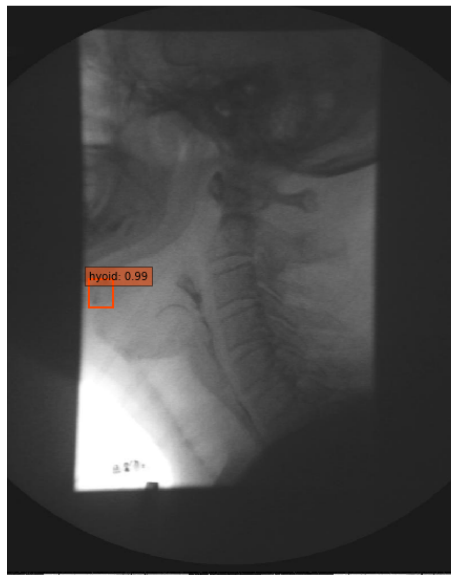
231 In this investigation, we aimed to detect the location of the hyoid bone in the videofluoroscopic
 232 images without human intervention. The hyoid bone is an important structure considered in
 233 dysphagia assessment. Its motion can be related to the severity of dysphagia and is used to



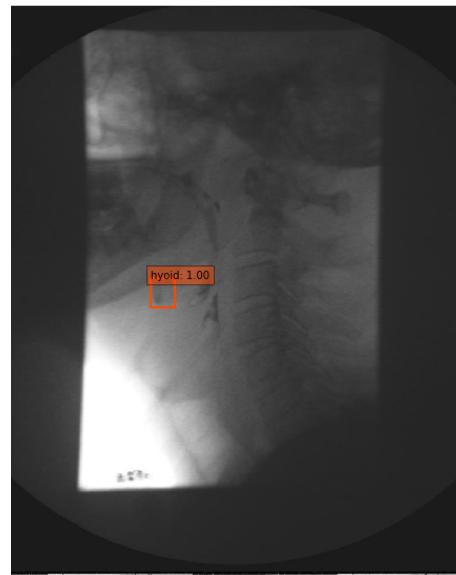
(a)



(b)



(c)



(d)

Figure 5: Results on different image conditions using SSD500-VGGNet: (a)(b) hyoid bone hides behind mandible (c)(d) hyoid bone is slightly blurred during motion

234 assess treatment effectiveness. Manual tracking of hyoid bone data from VFSS is the gold standard
235 accepted by experts and clinicians. Manually segmenting and annotating is time-consuming and
236 prone to judgment error. The hyoid bone motion data presented in this paper can be applied in
237 further investigations such as statistical methods and classification based on machine learning. A
238 quantitative and qualified computer-aided system is highly desirable in clinical work in which the
239 availability of an expert clinician to judge VFSS is not ubiquitous. Currently in dysphagia research,
240 human judgment is necessary to annotate hyoid position in initial video frames. Elimination or
241 mitigation of human judgment regarding hyoid motion could speed up image processing without
242 compromising accuracy. The following sections discuss the performance of each method and possible
243 factors that may have influenced the results.

244 We examined the performance of different object detection methods (Faster-RCNN, YOLOv2,
245 and SSD) to locate hyoid bone in our VFSS image dataset. For the deep architecture, we employed
246 the medium-size network VGGNet, the relatively larger-size network ResNet 101 for the SSD and
247 a small network ZFNet for Faster-RCNN. YOLOv2 is from the original Darknet model [34]. The
248 SSD500-VGGNet achieved better results than other CNN based models, indicating that it is the
249 most suitable method for hyoid bone detection in VFSS images. It is not surprising that YOLO
250 achieve the worst performance on VFSS data. The hyoid bone is a small object in the VFSS images.
251 YOLOv2 is a fast object detection method but is weak for small object detection as it applies global
252 features which doesn't obtain enough details for small objects. SSD500 is better than SSD 300 in
253 all settings by using ResNet-101 or VGGNet-16. The reasons might be as follows. SSD resizes the
254 input images to a fixed size: SSD300 resizes the images into 300×300 while SSD500 resizes images
255 into 500×500 . The training errors of SSD300 model is higher than those in SSD500. Resizing the
256 already small hyoid bone in images into a smaller size may result in a loss information. SSD300
257 cannot learn the details of the hyoid bone, which leads to worse performance. Furthermore, ResNet
258 reached a similar mAP to VGGNet in SSD500 but it has worse performance in SSD300. ResNet-
259 101 is a neural network with 101 layers, while VGG-16 has 16 layers. The similar performance
260 in SSD500 may indicate that both networks provide detailed information for the added features
261 layers. In the case of SSD300, the models with VGG networks had slightly smaller training loss after
262 iteration 8000, which might explain why VGG performed better on test data. The SSD method is
263 a powerful tool to detect the hyoid bone location, however, training SSD models with ResNet-101
264 and VGGNet with larger input size is time-consuming. We implemented our algorithms on the
265 NVIDIA Tesla M40 GPU. It took over one week to train both the SSD500-VGG16 models and
266 SSD500 with ResNet-101. The Faster-RCNN took only one day because ZFNet is a small neural

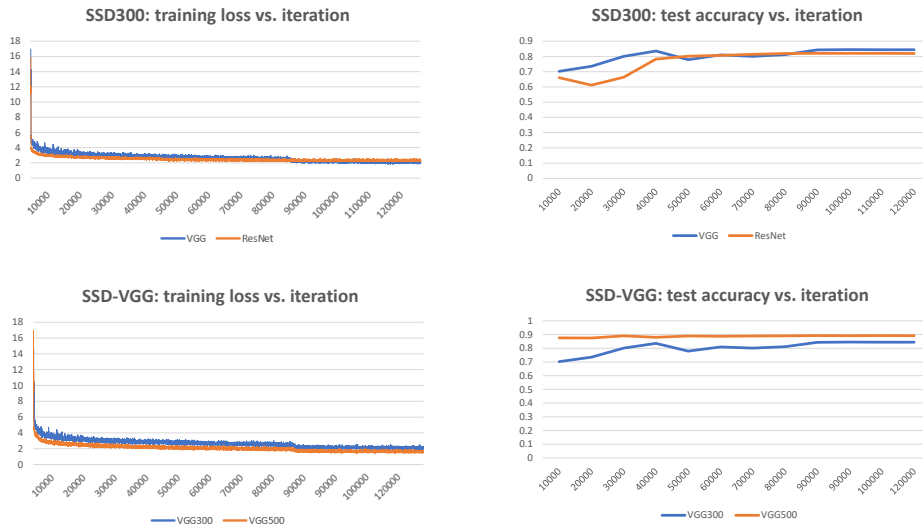


Figure 6: The influence of training loss and model performance of SSD models with different input sizes and pre-trained models.

267 network.

268 The hyoid bone moves upward and forward during a patient’s swallow. It will sometimes rise
 269 into the radiographic shadow of the mandible, obscuring its visibility by the judge/examiner. The
 270 judges must compare adjacent frames to infer the hyoid’s actual location when it is obscured by
 271 the mandible. Figure 5 (a) and (b) show the detection of the hyoid bone. Although the confidence
 272 score is low, our algorithm can be considered successful because experts may not be able to locate
 273 the hyoid bone. Figure 5 (c) and (d) are examples of blurred hyoid bone. The hyoid bone may be
 274 blurred when it moves quickly between two frames, but the algorithm can detect the moving hyoid
 275 bone with a high confidence score.

276 X-ray images vary in quality because clinicians control dosage to patients to the least amount
 277 of radiation as possible. Thus, as shown in the Figure 5, the brightness and, contrast of each x-ray
 278 image is different, altering the amount of useful information in each image. As shown in Figure 7,
 279 the SSD method detects the obscured hyoid bone location with a low confidence score or does not
 280 detect the hyoid bone location, similar to a guess when humans attempt to locate these cases. We
 281 know the location of the hyoid bone as the pre-knowledge, and seek to find a target around the
 282 predicted location while eliminating impossible regions one by one. The object detection algorithm

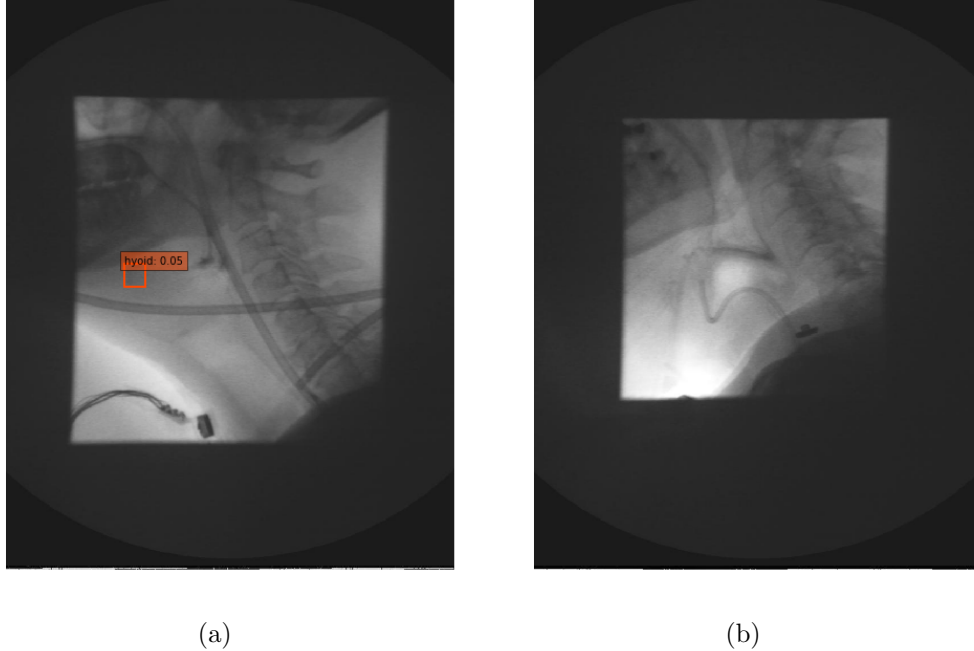


Figure 7: The cases which algorithm didn't detect the hyoid bone (a) the case with low confidence score (b) the case totally not detected

283 classifies the regions based on the default boxes, which is a direct way to make the decision and
 284 can't fully make use of outside information.

285 We investigated the performance of our model in the hyoid bone location task, however, our
 286 research had several limitations. X-rays images are often low quality, and the quality may vary
 287 from machine to machine. Whether the model can achieve similar performances across varied
 288 image quality requires further investigation. Furthermore, our investigation included data from
 289 265 patients from the same hospital, which may provide limited diagnostic variability in patients.
 290 Additional data should be collected to improve the performance and stability of our model. Prior
 291 research [43] indicated that Faster-RCNN with inception ResNet v2 has the best object detection
 292 results when compared to other modern object detection methods. Furthermore, several studies
 293 focused on small object detection, such as feature pyramid network [44], which may be a direc-
 294 tion for further research to increase the detection performance of the hyoid bone. For clinical
 295 relevance, future work should investigate automatic segmentation of hyoid bone areas, examine
 296 data to determine whether or not hyoid displacement is disordered, and determine if hyoid motion
 297 is the biomechanical etiology of impaired airway closure or upper esophageal sphincter opening.
 298 Moreover, since SSD detection methods detected the hyoid bone, future investigations will explore

299 detecting other key components in videofluoroscopy images. Given the millions of VFSS studies
300 implemented, high-accuracy component detection can save experts considerable time during their
301 diagnosis.

302 **4 Conclusion**

303 In this paper, we investigated hyoid bone detection in videofluoroscopy images using a deep learning
304 approach. We used 1434 swallows on VFSS videos as our dataset. The hyoid bone location was
305 manually annotated in each frame of the videos. We considered each frame as the single sample
306 and trained 70% of the frames using state-of-the-art object detection methods. The SSD-500
307 model tracked the location of the hyoid bone on each frame accurately. Ideally, hyoid bone motion
308 information can be used for physiological analysis. We believe that this proposed model has the
309 potential to improve the diagnosis assessment of dysphagia.

310 **Acknowledgment**

311 Research reported in this publication was supported by the Eunice Kennedy Shriver National
312 Institute of Child Health & Human Development of the National Institutes of Health under Award
313 Number R01HD092239, while the data was collected under Award Number R01HD074819. The
314 content is solely the responsibility of the authors and does not necessarily represent the official
315 views of the National Institutes of Health.

316 **Additional Information**

317 The authors declare that there is no conflict of interest.

318 **References**

- 319 [1] L. Sura, A. Madhavan, G. Carnaby, and M. A. Crary, “Dysphagia in the elderly: management
320 and nutritional considerations,” *Clinical Interventions in Aging*, vol. 7, p. 287, 2012.
- 321 [2] G. Mann, G. J. Hankey, and D. Cameron, “Swallowing disorders following acute stroke: preva-
322 lence and diagnostic accuracy,” *Cerebrovascular Diseases*, vol. 10, no. 5, pp. 380–386, 2000.

- 323 [3] N. P. Nguyen, C. Frank, C. C. Moltz, P. Vos, H. J. Smith, P. V. Bhamidipati, U. Karlsson,
324 P. D. Nguyen, A. Alfieri, L. M. Nguyen *et al.*, “Aspiration rate following chemoradiation for
325 head and neck cancer: an underreported occurrence,” *Radiotherapy and Oncology*, vol. 80,
326 no. 3, pp. 302–306, 2006.
- 327 [4] J. M. Dudik, I. Jestrović, B. Luan, J. L. Coyle, and E. Sejdić, “A comparative analysis of
328 swallowing accelerometry and sounds during saliva swallows,” *Biomedical Engineering online*,
329 vol. 14, no. 1, p. 3, 2015.
- 330 [5] D. G. Smithard, P. A. O’Neill, R. E. England, C. L. Park, R. Wyatt, D. F. Martin, and
331 J. Morris, “The natural history of dysphagia following a stroke,” *Dysphagia*, vol. 12, no. 4, pp.
332 188–193, 1997.
- 333 [6] N. Bhattacharyya, “The prevalence of dysphagia among adults in the united states,”
334 *Otolaryngology–Head and Neck Surgery*, vol. 151, no. 5, pp. 765–769, 2014.
- 335 [7] P. Clavé, R. Terré, M. De Kraa, and M. Serra, “Approaching oropharyngeal dysphagia,”
336 *Revista Espanola de Enfermedades Digestivas*, vol. 96, no. 2, pp. 119–131, 2004.
- 337 [8] L. Rofes, V. Arreola, J. Almirall, M. Cabré, L. Campins, P. García-Peris, R. Speyer, and
338 P. Clavé, “Diagnosis and management of oropharyngeal dysphagia and its nutritional and
339 respiratory complications in the elderly,” *Gastroenterology Research and Practice*, vol. 2011,
340 2010.
- 341 [9] O. B. Harrington, J. K. Duckworth, C. L. Starnes, P. White, L. Fleming, S. B. Kritchevsky,
342 and R. Pickering, “Silent aspiration after coronary artery bypass grafting,” *The Annals of*
343 *Thoracic Durgery*, vol. 65, no. 6, pp. 1599–1603, 1998.
- 344 [10] J. A. Hinchey, T. Shephard, K. Furie, D. Smith, D. Wang, S. Tonn *et al.*, “Formal dysphagia
345 screening protocols prevent pneumonia,” *Stroke*, vol. 36, no. 9, pp. 1972–1976, 2005.
- 346 [11] M. M. B. Costa, “Videofluoroscopy: the gold standard exam for studying swallowing and its
347 dysfunction,” *Arquivos de Gastroenterologia*, vol. 47, no. 4, pp. 327–328, 2010.
- 348 [12] S. ODonoghue and A. Bagnall, “Videofluoroscopic evaluation in the assessment of swallowing
349 disorders in paediatric and adult populations,” *Folia Phoniatica et Logopaedica*, vol. 51, no.
350 4-5, pp. 158–171, 1999.

- 351 [13] B. Martin-Harris, J. A. Logemann, S. McMahon, M. Schleicher, and J. Sandidge, “Clinical
352 utility of the modified barium swallow,” *Dysphagia*, vol. 15, no. 3, pp. 136–141, 2000.
- 353 [14] R. J. Hazelwood, K. E. Armeson, E. G. Hill, H. S. Bonilha, and B. Martin-Harris, “Identifica-
354 tion of swallowing tasks from a modified barium swallow study that optimize the detection of
355 physiological impairment,” *Journal of Speech, Language, and Hearing Research*, vol. 60, no. 7,
356 pp. 1855–1863, 2017.
- 357 [15] J. A. Logemann and J. A. Logemann, “Evaluation and treatment of swallowing disorders,”
358 1983.
- 359 [16] G. H. McCullough, R. T. Wertz, J. C. Rosenbek, R. H. Mills, W. G. Webb, and K. B. Ross,
360 “Inter-and intrajudge reliability for videofluoroscopic swallowing evaluation measures,” *Dys-*
361 *phagia*, vol. 16, no. 2, pp. 110–118, 2001.
- 362 [17] P. M. Kellen, D. L. Becker, J. M. Reinhardt, and D. J. Van Daele, “Computer-assisted assess-
363 ment of hyoid bone motion from videofluoroscopic swallow studies,” *Dysphagia*, vol. 25, no. 4,
364 pp. 298–306, 2010.
- 365 [18] I. Hossain, A. Roberts-South, M. Jog, and M. R. El-Sakka, “Semi-automatic assessment of
366 hyoid bone motion in digital videofluoroscopic images,” *Computer Methods in Biomechanics*
367 *and Biomedical Engineering: Imaging & Visualization*, vol. 2, no. 1, pp. 25–37, 2014.
- 368 [19] J. C. Lee, K. W. Nam, D. P. Jang, N. J. Paik, J. S. Ryu, and I. Y. Kim, “A supporting plat-
369 form for semi-automatic hyoid bone tracking and parameter extraction from videofluoroscopic
370 images for the diagnosis of dysphagia patients,” *Dysphagia*, vol. 32, no. 2, pp. 315–326, 2017.
- 371 [20] W.-S. Kim, P. Zeng, J. Q. Shi, Y. Lee, and N.-J. Paik, “Semi-automatic tracking, smoothing
372 and segmentation of hyoid bone motion from videofluoroscopic swallowing study,” *PloS one*,
373 vol. 12, no. 11, p. e0188684, 2017.
- 374 [21] S. Wang and R. M. Summers, “Machine learning and radiology,” *Medical Image Analysis*,
375 vol. 16, no. 5, pp. 933–951, 2012.
- 376 [22] Y. Le Cun, L. Jackel, B. Boser, J. Denker, H. Graf, I. Guyon, D. Henderson, R. Howard,
377 and W. Hubbard, “Handwritten digit recognition: Applications of neural network chips and
378 automatic learning,” *IEEE Communications Magazine*, vol. 27, no. 11, pp. 41–46, 1989.

- 379 [23] W. Sun, T.-L. B. Tseng, J. Zhang, and W. Qian, “Enhancing deep convolutional neural network
380 scheme for breast cancer diagnosis with unlabeled data,” *Computerized Medical Imaging and*
381 *Graphics*, vol. 57, pp. 4–9, 2017.
- 382 [24] M. H. Le, J. Chen, L. Wang, Z. Wang, W. Liu, K.-T. T. Cheng, and X. Yang, “Automated
383 diagnosis of prostate cancer in multi-parametric mri based on multimodal convolutional neural
384 networks,” *Physics in Medicine & Biology*, vol. 62, no. 16, p. 6497, 2017.
- 385 [25] K. Chockley and E. Emanuel, “The end of radiology? three threats to the future practice of
386 radiology,” *Journal of the American College of Radiology*, vol. 13, no. 12, pp. 1415–1420, 2016.
- 387 [26] Y. Dong, Y. Pan, J. Zhang, and W. Xu, “Learning to read chest x-ray images from 16000+
388 examples using CNN,” in *2017 IEEE/ACM International Conference on Connected Health:*
389 *Applications, Systems and Engineering Technologies (CHASE)*. IEEE, 2017, pp. 51–57.
- 390 [27] R. Girshick, J. Donahue, T. Darrell, and J. Malik, “Rich feature hierarchies for accurate object
391 detection and semantic segmentation,” in *Proceedings of the IEEE Conference on Computer*
392 *Vision and Pattern Recognition*, 2014, pp. 580–587.
- 393 [28] A. Sotiras, C. Davatzikos, and N. Paragios, “Deformable medical image registration: A survey,”
394 *IEEE Transactions on Medical Imaging*, vol. 32, no. 7, pp. 1153–1190, 2013.
- 395 [29] K. Keraudren, V. Kyriakopoulou, M. Rutherford, J. V. Hajnal, and D. Rueckert, “Localisation
396 of the brain in fetal mri using bundled sift features,” in *International Conference on Medical*
397 *Image Computing and Computer-Assisted Intervention*. Springer, 2013, pp. 582–589.
- 398 [30] J. Tang, R. M. Rangayyan, J. Xu, I. El Naqa, and Y. Yang, “Computer-aided detection
399 and diagnosis of breast cancer with mammography: recent advances,” *IEEE Transactions on*
400 *Information Technology in Biomedicine*, vol. 13, no. 2, pp. 236–251, 2009.
- 401 [31] T. Acharya and A. K. Ray, *Image processing: principles and applications*. John Wiley &
402 Sons, 2005.
- 403 [32] R. Girshick, “Fast R-CNN,” in *Proceedings of the IEEE International Conference on Computer*
404 *Vision*, 2015, pp. 1440–1448.
- 405 [33] S. Ren, K. He, R. Girshick, and J. Sun, “Faster r-cnn: Towards real-time object detection with
406 region proposal networks,” in *Advances in Neural Information Processing Systems*, 2015, pp.
407 91–99.

- 408 [34] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, “You only look once: Unified, real-time
409 object detection,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern
410 Recognition*, 2016, pp. 779–788.
- 411 [35] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, “Ssd: Single
412 shot multibox detector,” in *European Conference on Computer Vision*. Springer, 2016, pp.
413 21–37.
- 414 [36] G. Wang, M. Kalra, and C. G. Orton, “Machine learning will transform radiology significantly
415 within the next 5 years,” *Medical Physics*, 2017.
- 416 [37] Y. LeCun, Y. Bengio, and G. Hinton, “Deep learning,” *Nature*, vol. 521, no. 7553, pp. 436–444,
417 2015.
- 418 [38] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional
419 neural networks,” in *Advances in Neural Information Processing Systems*, 2012, pp. 1097–1105.
- 420 [39] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and
421 A. Rabinovich, “Going deeper with convolutions,” in *Proceedings of the IEEE Conference on
422 Computer Vision and Pattern Recognition*, 2015, pp. 1–9.
- 423 [40] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recog-
424 nition,” *arXiv preprint arXiv:1409.1556*, 2014.
- 425 [41] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in
426 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp.
427 770–778.
- 428 [42] J. Long, E. Shelhamer, and T. Darrell, “Fully convolutional networks for semantic segmenta-
429 tion,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*,
430 2015, pp. 3431–3440.
- 431 [43] J. Huang, V. Rathod, C. Sun, M. Zhu, A. Korattikara, A. Fathi, I. Fischer, Z. Wojna, Y. Song,
432 S. Guadarrama *et al.*, “Speed/accuracy trade-offs for modern convolutional object detectors,”
433 *arXiv preprint arXiv:1611.10012*, 2016.
- 434 [44] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, “Feature pyramid
435 networks for object detection,” *arXiv preprint arXiv:1612.03144*, 2016.

436 [45] E. Fisher, D. Austin, H. M. Werner, Y. J. Chuang, E. Bersu, and H. K. Vorperian, “Hyoid
437 bone fusion and bone density across the lifespan: prediction of age and sex,” *Forensic science,*
438 *medicine, and pathology*, vol. 12, no. 2, pp. 146–157, 2016.

439 **Author Contribution**

440 Zhenwei Zhang performed the experiments and wrote the manuscript with support from James
441 L. Coyle and Ervin Sejdić. All authors provided critical feedback and helped shape the research,
442 analysis and manuscript.