# Title: Deep learning-based auto-segmentation of swallowing and chewing structures

**Authors:** Aditi Iyer[1], Maria Thor[1], Rabia Haq[1], Joseph O. Deasy[1], Aditya P. Apte[1]

*[1]Department of Medical Physics, Memorial Sloan Kettering Cancer Center, New York, USA*

1        **Abstract**

2     **Purpose**

3     Delineating the swallowing and chewing structures in Head and Neck (H&N) CT scans is

4     necessary for radiotherapy treatment (RT) planning to reduce the incidence of radiation-induced

5     dysphagia, trismus, and speech dysfunction. Automating this process would decrease the manual

6     input required and yield reproducible segmentations, but generating accurate segmentations is

7     challenging due to the complex morphology of swallowing and chewing structures and limited

8     soft tissue contrast in CT images.

9     **Methods**

10    We trained deep learning models using 194 H&N CT scans from our institution to segment the

11    masseters (left and right), medial pterygoids (left and right), larynx, and pharyngeal constrictor

12    muscle using DeepLabV3+ with the resnet-101 backbone. Models were trained in a sequential

13    manner to guide the localization of each structure group based on prior segmentations.

14    Additionally, an ensemble of models was developed using contextual information from three

15    different views (axial, coronal, and sagittal), for robustness to occasional failures of the individual

16    models. Output probability maps were averaged, and voxels were assigned labels corresponding

17    to the class with the highest combined probability.

18    **Results**

19    The median dice similarity coefficients (DSC) computed on a hold-out set of 24 CT scans were

20    $0.87\pm0.02$ for the masseters, $0.80\pm0.03$ for the medial pterygoids, $0.81\pm0.04$ for the larynx, and

21    $0.69\pm0.07$for the constrictor muscle. The corresponding 95[th] percentile Hausdorff distances were

22    0.32±0.08cm (masseters), 0.42±0.2cm (medial pterygoids), 0.53±0.3cm (larynx), and

23    0.36±0.15cm (constrictor muscle). Dose-volume histogram (DVH) metrics previously found to

24    correlate with each toxicity were extracted from manual and auto-generated contours and

25    compared between the two sets of contours to assess clinical utility. Differences in DVH metrics

26    were not found to be statistically significant (p>0.05) for any of the structures. Further, inter-

27    observer variability in contouring was studied in 10 CT scans. Automated segmentations were

28    found to agree better with each of the observers as compared to inter-observer agreement,

29    measured in terms of DSC.

30    **Conclusions**

31    We developed deep learning-based auto-segmentation models for swallowing and chewing

32    structures in CT. The resulting segmentations can be included in treatment planning to limit

33    complications following RT  for H&N cancer. The segmentation models developed in this work

34    are distributed for research use through the open-source platform CERR, accessible at

35    https://github.com/cerr/CERR.

36

37

38

39

40

41

42

43

## 1. INTRODUCTION

45 Delineating organs at risk (OAR) is central in radiotherapy (RT) treatment planning to limit normal

46 tissue complications following RT. Manual delineation is time-consuming, subjective, and prone

47 to errors due to factors such as organ complexity and level of experience [1]. Therefore, there has

48 been a great deal of interest in automating this process to generate accurate and reproducible

49 segmentations in a time-efficient manner.

50 In head and neck (H&N) cancer treatment, isolating the larynx and pharyngeal constrictor muscle

51 is of particular interest to limit speech dysfunction and dysphagia, respectively, following RT [2]

52 [3]. The masseters and medial pterygoids have also been identified as critical structures in limiting

53 radiation-induced trismus [4]. However, delineating these structures is challenging due to their

54 complex morphology and the low soft tissue contrast in CT images.

55 Few semi-automatic and automatic methods have been previously developed to segment OARs in

56 the H&N. Conventional multi-atlas-based auto-segmentation (MABAS) methods involve

57 propagating and combining manually-segmented OARs from a curated library of CT scans through

58 image registration as in [5] and [6]. Other approaches offer strategies to further refine MABAS

59 using organ-specific intensity [7] and texture [8] features or shape representation models [9] [10].

60 However, MABAS is sensitive to inter-subject anatomical variations as well as image artifacts,

61 and image registration is computationally intensive, requiring several minutes even with highly

62 efficient implementations [11].

63 Convolutional neural networks (CNNs) have recently been applied successfully to various medical

64 image segmentation tasks. Ibragimov et al. [12] trained 13 CNNs, applied in sliding-window

65  fashion to segment H&N OARs including the larynx and pharynx, which were further refined

66  using Markov Random Field (MRF)-based post-processing. In [13], Ward van Rooij et al.

67  employed the popular 3D U-net [14] architecture to segment H&N OARs including the pharyngeal

68  constrictor muscle. Zhu et all. [11] extended the U-Net model by incorporating squeeze-and

69  excitation (SE) residual blocks and a modified loss function to improve segmentation of smaller

70  structures such as the chiasm and optic nerves. In [15], Men et al. segmented nasopharyngeal tumor

71  volumes in H&N CT using a modified version of the VGG-16 [16] architecture, replacing fully-

72  connected layers with fully-convolutional layers and introducing improved decoder networks to

73  rebuild high-resolution feature maps. Tong et al [17] trained a fully convolutional neural net

74  (FCNN), incorporating prior information by training a shape representation model to regularize

75  shape characteristics of 9 H&N OARs. The FocusNet [18] developed by Gao et al. utilizes multiple

76  CNNs to segment H&N OARs including the larynx, first segmenting large structures, then

77  segmenting smaller structures with specifically designed sub-networks.

78  In this work, we present a fully automatic method to segment swallowing and chewing structures

79  and examine its suitability for clinical use. To the best of our knowledge, this is the first [35] deep

80  learning-based method for segmenting the chewing structures in CT images. We propose a novel

81  framework in which DeepLabV3+ segmentation models are trained sequentially to guide the

82  localization of each structure group based on previously-segmented structures. Model ensembles

83  are created using three orthogonal views (axial, sagittal, and coronal) in 2.5D and shown to

84  improve segmentation accuracy of H&N OARs compared to models trained on a single-view.

85  **2. MATERIALS AND METHODS**

86  **2.1 Dataset**

87  CT scans of 243 H&N cancer patients from our institution were retrospectively collected under

88  IRB 16-142 and accessed under IRB 16-1488 to develop the auto-segmentation methods in this

89  work. The dataset was randomly partitioned into training (80%), validation (10%), and testing

90  (10%) sets. The validation set (10%) was used to estimate model performance while tuning

91  hyperparameters, and the testing set (10%) was used for unbiased evaluation of the final model.

92  Representative CT images showing considerable variation in head pose, shape, and appearance

93  around the structures of interest, including slices with dental artifacts due to dental implants, were

94  included in the dataset. Additional data characteristics are listed in Table 1. Manual segmentations

95  of the masseters (left, right), medial pterygoids (left, right),  larynx, and constrictor muscle,

96  generated using MSKCC's in-house treatment planning system, served as the reference standard.

97  Reference  contours of the masseters and medial pterygoids were available in 60% of the scans,

98  and the constrictor muscle in 97% of the scans (larynx was provided in all scans).

99  **Table 1**. Summary of data characteristics for axial images

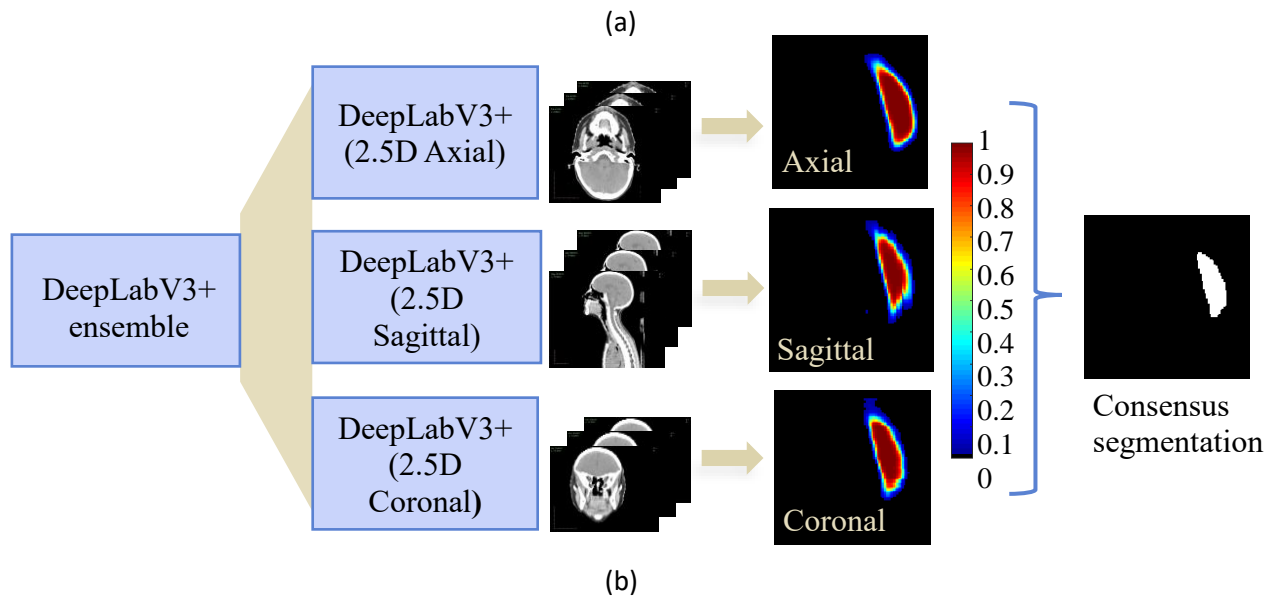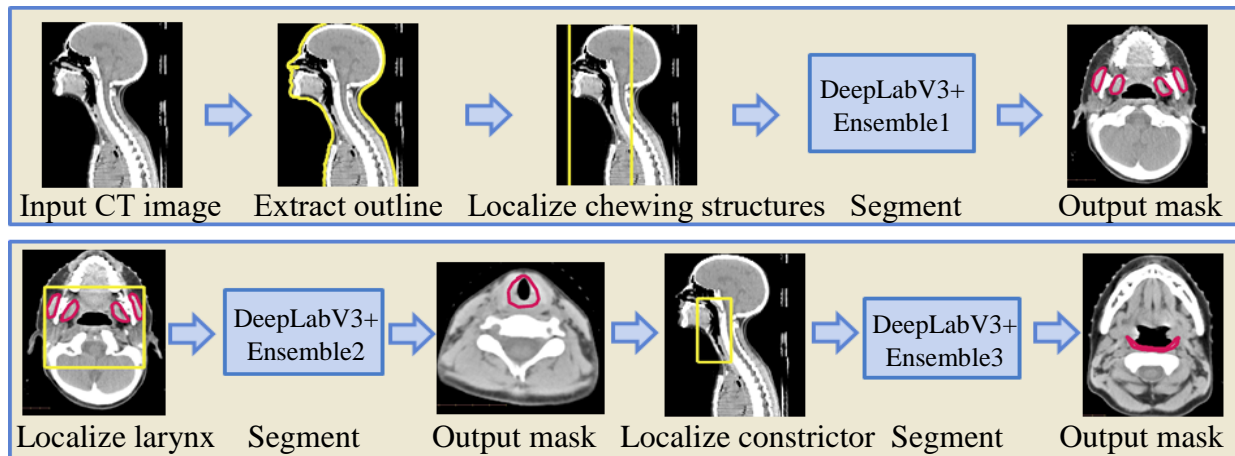| Attribute | Median | Min | Max | 5th percentile | 95th percentile |
|---|---|---|---|---|---|
| Number of slices | 80 | 48 | 279 | 65 | 103 |
| In-plane spacing (mm) | 0.78 | 0.51 | 1.37 | 0.61 | 1.17 |
| Slice thickness (mm) | 3.00 | 1.25 | 6.00 | 2.50 | 3.27 |

100

101  **2.2 CNN segmentation model**

102  The DeepLabV3+ CNN [19] was selected due to its impressive performance on the PASCAL VOC

103  2012 and Cityscapes datasets. Moreover, it has previously been applied successfully to medical

104  image segmentation tasks, e.g., by Elguindi et al. [26] to segment the prostate in MR images and

105    by Haq et al. [27] to segment cardio-pulmonary substructures in CT images. This architecture

106    offers the advantage of an encoder network that can capture contextual information at different

107    scales using multiple dilated convolution layers applied in parallel at different rates and a decoder

108    network capable of effectively recovering object boundaries.

109    Here, training was performed using ResNet-101 [20] as the backbone of the encoder network. A

110    publicly-distributed implementation [21] of DeepLabV3+ using the Pytorch [33] framework was

111    utilized, and a soft-max layer was appended to obtain voxel-wise probabilities. On training a single

112    multi-class model to segment all the structures of interest, it was observed that the imbalance in

113    class labels due to large differences in OAR sizes led to poor performance on smaller structures.

114    Augmenting the loss function based on class frequencies did not produce significant improvement,

115    as previously observed in [18]. Consequently, three separate models were developed- one to

116    segment the chewing structures and one each for the larynx and the constrictor muscle. A

117    sequential segmentation strategy (figure 1) was utilized in which each segmented OAR was used

118    to constrain the location of subsequently segmented OARs. Additionally, an ensemble of three

119    models was developed per OAR group using 2D axial, sagittal, and coronal slices. As compared

120    to 3D convolutional neural networks (CNNs) which are highly memory-intensive, training in 2D

121    enabled us to employ a more complex CNN while still providing contextual information and

122    increased redundancy from the three different orientations.

123    A high-performance cluster with four NVIDIA GeForce GTX 1080Ti GPUs with 11GB of

124    memory each was used in training. Batch-normalization was applied using mini-batches of 8

125    images, resized to 320 x 320 voxels. Input images were standardized by scaling the intensities to

126    [0,1] and normalizing to zero-mean and unit-variance. Data augmentation was performed through

127    image scaling, cropping, and rotation. Input channels were populated using three consecutive slices

128    (axial, sagittal, and coronal, respectively, for the three models), and the cross-entropy loss was

129    employed in training. Data preprocessing and export to HDF5 [28] format was performed using

130    the Computational Environment for Radiological Research (CERR) [23]. Details of the individual

131    models are presented in the following sections, and Table 2 summarizes the hyperparameters used.
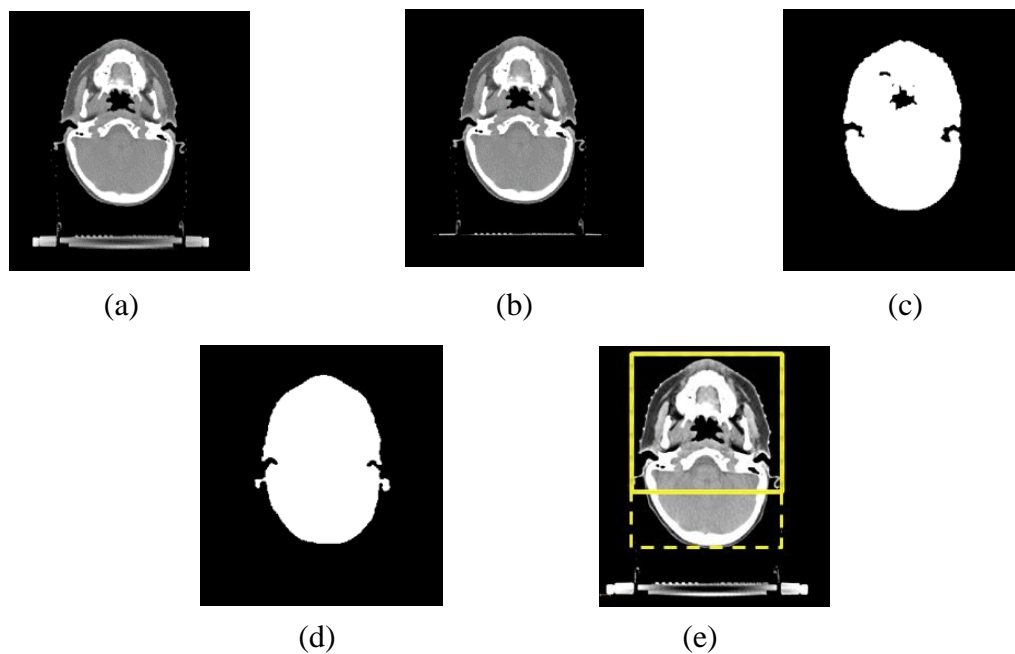


(a)



(b)

132    **Figure 1.** (a) Sequential framework for segmenting chewing and swallowing structures using deep

133    learning, in which each segmented OAR group is used to improve localization of subsequently

134    segmented OARs. (b) Example showing consensus segmentation of left masseter using ensemble

135    model trained on 3 orthogonal views (axial, sagittal, and coronal).

136    **2.3 Chewing structures**

137    A multi-label model was trained to segment the left and right masseters and medial pterygoids.

138    H&N CT scans were automatically cropped around the patient's outline prior to training. The

139    couch was first detected using the Hough transform and masked out. This was followed by

140    intensity-based thresholding and morphological post-processing to extract the patient's outline. A

141    bounding box was generated around the outline and its posterior extent was limited by 25%. 2D

142    axial, sagittal, and coronal images and corresponding masks of the chewing structures were

143    extracted within these extents for training.



(a)                                        (b)                                        (c)

(d)                                        (e)

144

145    **Figure 2.** Illustration of method to localize chewing structures in axial H&N CT scans. (a) Input

146    CT image (b) Hough transform-based couch segmentation (c) Intensity thresholding to extract

147    patient outline (d) Morphological post-processing (e) Bounding box with reduced posterior extent.

148    Model weights were optimized using Stochastic Gradient Descent (SGD) with hyperparameters

149    listed in Table 2 and learning rate was decayed following the polynomial scheduling policy. An

150 early stopping strategy was employed to avoid overfitting if there was no improvement in the

151 validation loss.

152 **2.4 Swallowing structures**

153 **2.4.1 Larynx**

154 CT scans were cropped further, with the anterior, left, right, and superior limits defined by the

155 corresponding extents of the previously segmented chewing structures, as shown in figure 3. 2D

156 axial, sagittal and coronal images and corresponding masks of the larynx were extracted within the

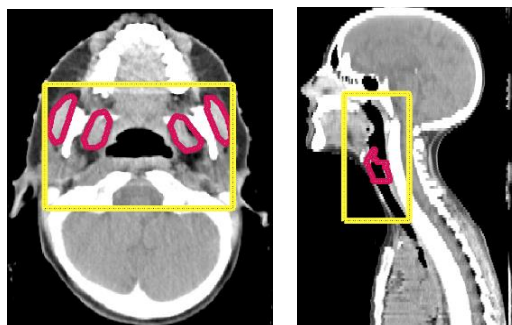157 resulting bounding box and optimization was performed using SGD.



158

159 **Figure 3.** Bounding box (yellow) for localization of larynx using previously segmented chewing

160 structures (red).

161 **2.4.2 Constrictor muscle**

162 To localize the constrictor muscle, CT scans were cropped with the anterior, left, right, and

163 superior limits defined by the corresponding extents of the chewing structures. The posterior and

164 inferior limits were defined by corresponding extents of the larynx, with sufficient padding. 2D

165 axial, sagittal,  and coronal images and corresponding masks of the constrictor muscle were

166 extracted within the resulting bounding box, shown in figure 4. Optimization was performed using

167 adaptive moment estimation (Adam) [22].



168

169 **Figure 4.** Bounding box (yellow) for localization of constrictor muscle using previously

170 segmented chewing structures and larynx (red).

171 For each of the above OARs, probability maps returned by the three models (axial, sagittal, and

172 coronal) were averaged and voxels were assigned to the class with the highest resulting probability

173 to produce stable consensus segmentations, robust to occasional failures of the individual models.

174 The segmentation masks were further post-processed to remove isolated voxels by discarding all

175 but the largest connected component. Morphological processing was performed to fill holes and

176 obtain smooth contours.

177 **Table 2.** Summary of hyperparameters for training.

| Model | Structure(s) | Learning rate | Optimizer | Momentum | Weight Decay |
|-------|-------------|---------------|-----------|----------|--------------|
| 1 | MML[a], MMR[b], PML[c], PMR[d] | Axial: 0.002<br>Sagittal: 0.003<br>Coronal: 0.002 | SGD | 0.9 | 0.0001 |

| 2 | Larynx | Axial: 0.0003 | SGD | 0.9 | 0.0002 |
| | | Sagittal: 0.0003 | | | |
| | | Coronal: 0.0003 | | | |
| 3 | CM$^e$ | Axial: 1x10$^{-6}$ | Adam | 0.9 | 0.0001 |
| | | Sagittal: 8x10$^{-7}$ | | | |
| | | Coronal : 2x10$^{-6}$ | | | |

$^a$ Masseter (left); $^b$ masseter (right); $^c$ medial pterygoid (left); $^d$ medial pterygoid (right); $^e$ constrictor muscle.
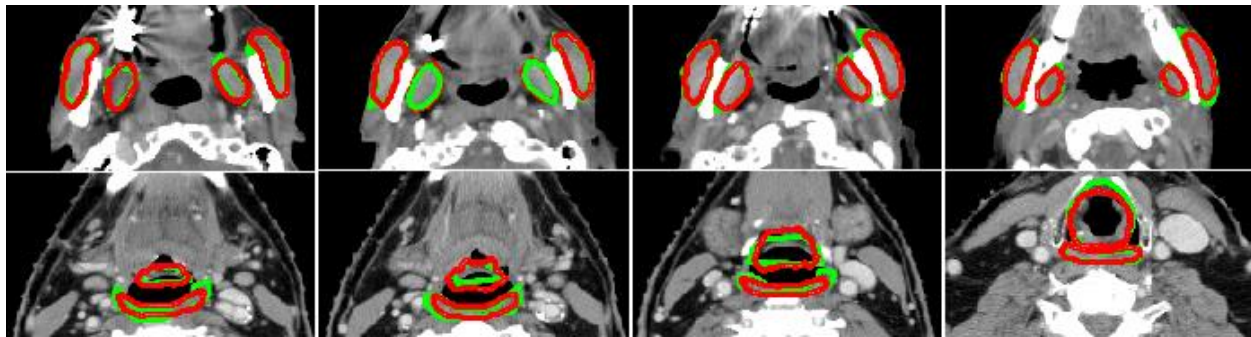
## 3. RESULTS

### 3.1 Comparison to reference standard

The performance of the proposed method was evaluated on the test set of 24 H&N cancer patients by comparing the results against existing manually delineated segmentations. The Dice Similarity Coefficient (DSC) was used to assess the degree of overlap between manual (A) and automated (B) segmentations, computed as:

$$DSC = 2\ \frac{|A \cap B|}{|A| + |B|}$$

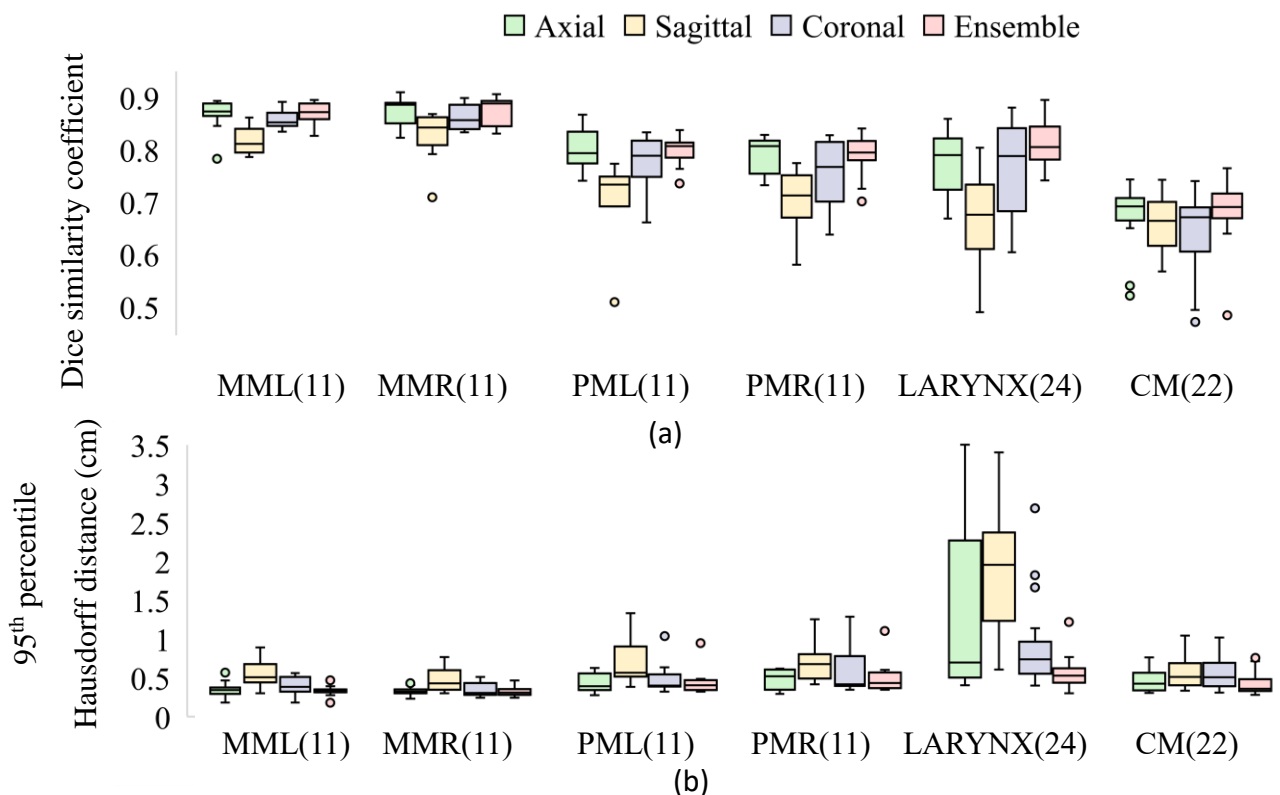Additionally, the 95$^{th}$ percentile of the Hausdorff distance (HD$_{95}$), i.e., maximum distance between boundary points of A and B, was computed to capture the impact of a few sizeable segmentation errors on the overall segmentation quality. Both DSC and HD$_{95}$ were measured for the individual

188    models (axial, sagittal, coronal) as well as the ensemble to investigate the benefits of a multi-view

189    ensemble.

190    Examples of segmentations resulting from our algorithm are presented in figure 5 for qualitative

191    assessment. Box plots of DSC and HD$_{95}$ are presented in figure 6.



192

193    **Figure 5.** Auto-segmentation results for chewing structures (row-1) and swallowing structures

194    (row-2), shown in four axial cross-sections. Manual reference segmentations are depicted in green

195    and deep-learning-based auto-segmentations in red.

196 **Figure 6.** Performance of deep learning models (axial, sagittal, coronal, and ensemble) compared

197 to manual reference segmentations in terms of (a) DSC and (b) $HD_{95}$. Of the 24 test patients, the

198 number with manual contours available for comparison is noted in parentheses. *MML, MMR:*

199 *masseters (left and right), PML, PMR: medial pterygoids (left and right), CM : constrictor muscle.*

200 Differences in mean dose, previously identified [29-30] as a possible factor in radiation-induced

201 complications were also computed between deep learning-based and manual contours, and the

202 Wilcoxon signed-rank test was applied to investigate potential statistical disparities (Table 3).

203 Differences in mean dose were not found to be statistically significant for all tested structures at

204 significance level 5%.

205 **Table 3**. Comparison of mean doses extracted using manual and auto-generated contours. Median,

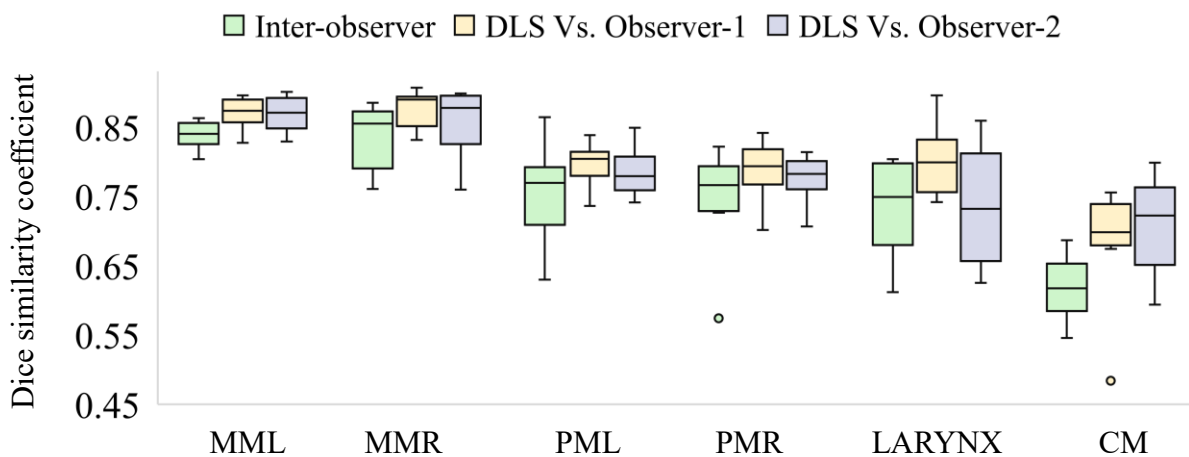206 first and third quartiles of percentage differences are presented.

| Structure | Metric | Difference (%) | p-value | No. test patients |
|-----------|--------|----------------|---------|-------------------|
| MML[b] U MMR[c] | Ipsilateral[a] mean dose | -0.01 (-1.20, 1.32) | 1.00 | 11 |
| PML[d] U PMR[e] | Ipsilateral[a] mean dose | 0.53 (-0.77, 1.11) | 0.70 | 11 |
| Larynx | Mean dose | 0.38 (-2.02, 6.59) | 0.33 | 24 |
| CM[f] | Mean dose | 0.15 (-0.49, 1.28) | 0.29 | 22 |

[a]Ipsilaterality for the paired structures was decided based upon the side with the highest dose. [b]

Masseter (left); [c] masseter (right); [d] medial pterygoid (left); [e] medial pterygoid (right); [f]

constrictor muscle.

207 **3.2 Variability of reference standards**

208    The manual reference segmentations came from various sources: the larynx was delineated for

209    treatment planning by multiple observers with variation in level of expertise; the constrictors and

210    chewing structures were each delineated post-treatment by one of two radiation oncology residents

211    for studying dysphagia and trismus, respectively [29] [30]. Regardless of origin, these contours

212    are referred to as observer-1. A randomly selected subset of 10 CT scans from the testing dataset

213    were then re-segmented by a medical physicist (observer-2) to provide a second reference

214    standard. The agreement between the independent observers was measured in terms of DSC and

215    compared to agreement of each observer with the deep learning-based contours (figure 7).

216    Automated segmentations showed better agreement with each observer as compared to the inter-

217    observer agreement.



218

**Figure 7.** Comparing agreement of deep learning-based segmentations (DLS) with independent

220    observers vs. inter-observer agreement. *MML, MMR: Masseters (left and right), PML, PMR:*

221    *medial pterygoids (left and right), CM (constrictor muscle)*

## 3.3 Comparison to previously reported methods

223    We compared the performance of our method with H&N OARs segmentation models from the

224    published literature. These included state-of-the-art CNN models such as the 3D Unet [13],

225  FocusNet [18], a H&N OAR-focused CNN [12], commercial atlas-based segmentation tool SPICE

226  [5] and other atlas-based methods [24,25]. The mean DSCs for the different methods are

227  summarized in Table 4. The performance of our segmentation models matched or exceeded

228  previously presented methods. However, it should be noted that these results were reported on

229  different datasets and do not represent a direct comparison performed on our testing dataset.

**3.4 Distribution of trained models**

231  The segmentation models developed in this work are publicly distributed as a part of CERR's

232  library of model implementations [31]. Model dependencies were encapsulated using Singularity

233  [34] containers, enabling deployment on a variety of scientific computing architectures and aiding

234  in portability and reproducibility. This also facilitates further analysis through integration with

235  CERR's radiomics toolbox [32] and dosimetric models [31]. CERR can be downloaded from

236  https://github.com/cerr/CERR. The list of currently available segmentation containers and

237  corresponding links for download are available at https://github.com/cerr/CERR/wiki/Auto-

238  Segmentation-models. It should be noted that the segmentation models are distributed strictly for

239  research use. Clinical or commercial use is prohibited. CERR and containerized model

240  implementations have not been approved by the U.S. Food and Drug Administration (FDA).

241  **Table 4.** DSC (mean ± std. deviation) for swallowing and chewing structures using the proposed

242  method (column-1) and results from previously published methods.

243

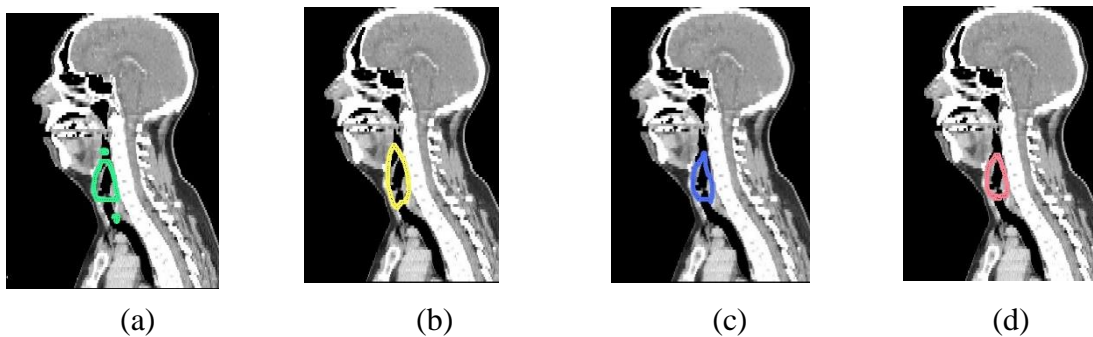| Structure | Proposed | van Rooij et al. [13] (2019)[1] | Gao et al. [18] (2019)[1] | Ibragimov et al. [12] (2017)[1] | Tao et al. [5] (2015)[2] | Thomson et al. [24] (2015)[2] | Han et al. [25] (2008)[2] |
|---|---|---|---|---|---|---|---|
| Masseters | $0.87 \pm 0.02$ | - | - | - | - | - | 0.83 |
| Pterygoids | $0.80 \pm 0.03$ | - | - | - | - | - | 0.83 |
| Larynx | $0.81 \pm 0.04$ | $0.78 \pm 0.05$ | $0.66 \pm 0.29$ | $0.86 \pm 0.04$ | $0.73 \pm 0.04$ | 0.58 | - |
| Constrictor | $0.68 \pm 0.07$ | $0.68 \pm 0.09$ | - | $0.69 \pm 0.06$ | $0.65 \pm 0.06$ | 0.50 | - |

244 [1]Deep learning based  [2] atlas based.

## 4. DISCUSSION

246  We trained three distinct model ensembles to segment structures of different sizes and constrained

247  the location of each structure group based on the extents of previously identified structures. This

248  sequential localization framework was able to handle imbalance in class labels due to differences

249  in the sizes of the OARs. Additionally, we used an ensemble of models trained on three different

250  image orientations to capture important contextual information and provide redundancy in case of

251  failures of the individual models. The performance of the of the multi-view ensemble in terms of

252  DSC either matched or out-performed the individual models and showed lower variance across all

253  structures, as evidenced by tighter bounds in the box plots (figure 6). Further, the multi-view

254  consensus reduced the worst-case segmentation errors (as captured by $HD_{95}$) across all the

255  structures. This suggests that employing a multi-view ensemble may help produce more stable

256  segmentations than single-view models

257  Of the structures considered, the pharyngeal constrictor muscle was most challenging to segment

258  due to its morphological complexity, high anatomical variability, and low soft tissue contrast. This

259  was reflected in our analysis of manual segmentations by different observers. The lowest inter-

260  observer agreement was observed in delineating the constrictor. In some cases, segmentation was

261  further complicated by tumor infiltration around the periphery. For the larynx, standardizing the

262  reference segmentations (generated by multiple observers with varying levels of experience) prior

263  to training could potentially yield better results. We also observed a few spurious and

264  discontinuous detections of the larynx when using models trained only on axial or sagittal images.

265  This was mitigated by generating a consensus segmentation using information from all 3 views

266  (e.g. figure 8). Finally, the deep learning models for all structures were found to generalize well as

267  the auto-generated results showed good agreement with delineations by a new (unseen) observer.

268  We further investigated the potential for clinical application of the trained models by comparing

269  mean doses extracted from manual and automated segmentations. No statistical differences were

270  observed at the 5% significance level.



(a)　　　　　　　(b)　　　　　　　(c)　　　　　　　(d)

271  **Figure 8.** Sagittal cross sections showing auto-segmented larynx using (a) axial only (b) sagittal

272  only (c) coronal only and (d) ensemble models.

273  **5. CONCLUSIONS**

274  We developed a fully-automatic, accurate, and time-efficient method to segment swallowing and

275  chewing structures in CT images and demonstrated its potential for clinical use. The proposed

276  method introduces a novel framework for sequential localization and segmentation to handle

277  imbalances in OAR sizes. Additionally, the potential advantage of a multi-view ensemble over a

278  single-view model was investigated. As hypothesized, the ensemble models were found to yield

279  more stable segmentations across all structures. This ensemble approach could be applied to

280  improve segmentation quality in other sites as well. The trained models are publicly distributed

281  for research use through the open-source platform CERR using Singularity containers.

282  **REFERENCES**

283  1.  Harari PM, Song S, Tomé WA. Emphasizing conformal avoidance versus target definition

284      for IMRT planning in head-and-neck cancer. International journal of radiation oncology,

285      biology, physics. https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2905233/.

286  2.  Choi M, Refaat T, Lester MS, Bacchus I, Rademaker AW, Mittal BB. Development of a

287      standardized method for contouring the larynx and its substructures. Radiation oncology

288      (London, England). https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4266916/.

289  3.  Levendag PC, Teguh DN, Voet P, et al. Dysphagia disorders in patients with cancer of the

290      oropharynx are significantly affected by the radiation therapy dose to the superior and middle

291      constrictor muscle: A dose-effect relationship. Radiotherapy and Oncology. 2007;85(1):64-

292      73. doi:10.1016/j.radonc.2007.07.009

293  4.  Kraaijenga SA, Hamming-Vrieze O, Verheijen S, et al. Radiation dose to the masseter and

294      medial pterygoid muscle in relation to trismus after chemoradiotherapy for advanced head

295      and neck cancer. Head & Neck. 2019;41(5):1387-1394. doi:10.1002/hed.25573

296  5.  Tao C-J, Yi J-L, Chen N-Y, et al. Multi-subject atlas-based auto-segmentation reduces

297      interobserver variation and improves dosimetric parameter consistency for organs at risk in

298    nasopharyngeal carcinoma: A multi-institution clinical study. Radiotherapy and Oncology.

299    2015;115(3):407-411. doi:10.1016/j.radonc.2015.05.012

300    6.  Levendag P, Hoogeman M, Teguh D, et al. Atlas Based Auto-segmentation of CT Images:

301    Clinical Evaluation of using Auto-contouring in High-dose, High-precision Radiotherapy of

302    Cancer in the Head and Neck. International Journal of Radiation Oncology*Biology*Physics.

303    2008;72(1). doi:10.1016/j.ijrobp.2008.06.1285

304    7.  Fortunati V, Verhaart RF, Lijn FVD, et al. Tissue segmentation of head and neck CT images

305    for treatment planning: A multiatlas approach combined with intensity modeling. Medical

306    Physics. 2013;40(7):071905. doi:10.1118/1.4810971

307    8.  Qazi AA, Pekar V, Kim J, Xie J, Breen SL, Jaffray DA. Auto-segmentation of normal and

308    target structures in head and neck CT images: A feature-driven model-based approach.

309    Medical Physics. 2011;38(11):6160-6170. doi:10.1118/1.3654160

310    9.  Wang Z, Wei L, Wang L, Gao Y, Chen W, Shen D. Hierarchical Vertex Regression-Based

311    Segmentation of Head and Neck CT Images for Radiotherapy Planning. IEEE Transactions

312    on Image Processing. 2018;27(2):923-937. doi:10.1109/tip.2017.2768621

313    10. Fritscher KD, Peroni M, Zaffino P, Spadea MF, Schubert R, Sharp G. Automatic

314    segmentation of head and neck CT images for radiotherapy treatment planning using multiple

315    atlases, statistical appearance models, and geodesic active contours. Medical Physics.

316    2014;41(5):051910. doi:10.1118/1.4871623

317    11. Zhu W, Huang Y, Zeng L, et al. AnatomyNet: Deep learning for fast and fully automated

318    whole-volume segmentation of head and neck anatomy. Medical Physics. 2018;46(2):576-

319    589. doi:10.1002/mp.13300

320    12.  Ibragimov B, Xing L. Segmentation of organs-at-risks in head and neck CT images using

321         convolutional neural networks. Medical Physics. 2017;44(2):547-557. doi:10.1002/mp.12045

322    13. Rooij WV, Dahele M, Brandao HR, Delaney AR, Slotman BJ, Verbakel WF. Deep Learning-

323         Based Delineation of Head and Neck Organs at Risk: Geometric and Dosimetric Evaluation.

324         International Journal of Radiation Oncology*Biology*Physics. 2019;104(3):677-684.

325         doi:10.1016/j.ijrobp.2019.02.040

326    14. Ronneberger O, Fischer P, Brox T. U-Net: Convolutional Networks for Biomedical Image

327         Segmentation. Lecture Notes in Computer Science Medical Image Computing and

328         Computer-Assisted Intervention – MICCAI 2015. 2015:234-241. doi:10.1007/978-3-319-

329         24574-4_28

330    15. Men K, Chen X, Zhang Y, et al. Deep Deconvolutional Neural Network for Target

331         Segmentation of Nasopharyngeal Cancer in Planning Computed Tomography Images.

332         Frontiers in Oncology. 2017;7. doi:10.3389/fonc.2017.00315

333    16. Simonyan K, Zisserman A. Very Deep Convolutional Networks for Large-Scale Image

334         Recognition. https://arxiv.org/abs/1409.1556.

335    17. Tong N, Gou S, Yang S, Ruan D, Sheng K. Fully automatic multi-organ segmentation for

336         head and neck cancer radiotherapy using shape representation model constrained fully

337         convolutional neural networks. Medical Physics. 2018;45(10):4558-4567.

338         doi:10.1002/mp.13147

339    18.  Gao Y, Huang R, Chen M, et al. FocusNet: Imbalanced Large and Small Organ

340         Segmentation with an End-to-End Deep Neural Network for Head and Neck CT Images.

341         Lecture Notes in Computer Science Medical Image Computing and Computer Assisted

342         Intervention – MICCAI 2019. 2019:829-838. doi:10.1007/978-3-030-32248-9_92

343    19. Chen L-C, Zhu Y, Papandreou G, Schroff F, Adam H. Encoder-Decoder with Atrous

344       Separable Convolution for Semantic Image Segmentation. Computer Vision – ECCV 2018

345       Lecture Notes in Computer Science. 2018:833-851. doi:10.1007/978-3-030-01234-2_49

346    20. He K, Zhang X, Ren S, Sun J. Deep Residual Learning for Image Recognition. 2016 IEEE

347       Conference on Computer Vision and Pattern Recognition (CVPR). 2016.

348       doi:10.1109/cvpr.2016.90

349    21. Zhang J. jfzhang95/pytorch-deeplab-xception. GitHub. https://github.com/jfzhang95/pytorch-

350       deeplab-xception.git.

351    22. Kingma D, Ba J. Adam: A Method for Stochastic Optimization. arXiv preprint.

352       https://arxiv.org/pdf/1412.6980.

353    23. Deasy JO, Blanco AI, Clark VH. CERR: A computational environment for radiotherapy

354       research. Medical Physics. 2003;30(5):979-985. doi:10.1118/1.1568978

355    24. Thomson D, Boylan C, Liptrot T, et al. Evaluation of an automatic segmentation algorithm

356       for definition of head and neck organs at risk. Radiation Oncology. 2014;9(1):173.

357       doi:10.1186/1748-717x-9-173.

358    25. Han X, Hoogeman MS, Levendag PC, et al. Atlas-Based Auto-segmentation of Head and

359       Neck CT Images. Medical Image Computing and Computer-Assisted Intervention – MICCAI

360       2008. 2008:434-441. doi:10.1007/978-3-540-85990-1_52

361    26. Elguindi S, Zelefsky MJ, Jiang J, et al. Deep learning-based auto-segmentation of targets and

362       organs-at-risk for magnetic resonance imaging only planning of prostate radiotherapy.

363       Physics and Imaging in Radiation Oncology. 2019;12:80-86. doi:10.1016/j.phro.2019.11.006

364    27. Haq R, Hotca A, Apte A, Rimner A, Deasy JO, Thor M. Cardio-pulmonary substructure

365       segmentation of radiotherapy computed tomography images using convolutional neural

366    networks for clinical outcomes analysis. Physics and Imaging in Radiation Oncology.

367    https://www.sciencedirect.com/science/article/abs/pii/S2405631620300221

368    28. Folk, M., Cheng, A., & Yates, K. (1999, November). HDF5: A file format and I/O library for

369    high performance computing applications. In Proceedings of supercomputing (Vol. 99, pp. 5-

370    33).

371    29.  Tsai CJ, Jackson A, Setton J, et al. Modeling Dose Response for Late Dysphagia in Patients

372    With Head and Neck Cancer in the Modern Era of Definitive Chemoradiation. JCO Clinical

373    Cancer Informatics. 2017;(1):1-7. doi:10.1200/cci.17.00070

374    30.  Rao SD, Saleh ZH, Setton J, et al. Dose-volume factors correlating with trismus following

375    chemoradiation for head and neck cancer. Acta Oncologica. 2015;55(1):99-104.

376    doi:10.3109/0284186x.2015.1037864

377    31. Apte AP, Iyer A, Thor M, et al. Library of deep-learning image segmentation and outcomes

378    model-implementations. Physica Medica.

379    https://www.sciencedirect.com/science/article/pii/S1120179720300922. Published May 1,

380    2020. Accessed July 27, 2020.

381    32. Apte AP, Iyer A, Crispin-Ortuzar M, et al. Technical Note: Extension of CERR for

382    computational radiomics: A comprehensive MATLAB platform for reproducible radiomics

383    research. Medical Physics. 2018;45(8):3713-3720. doi:10.1002/mp.13046

384    33. Paszke, A.; Gross, S.; Massa, F.; Lerer, A.; Bradbury, J.; Chanan, G.; Killeen, T.; Lin, Z.;

385    Gimelshein, N.; Antiga, L.; et al. PyTorch: An imperative style, high-performance deep

386    learning library. In Proceedings of the Advances in Neural Information Processing Systems,

387    Vancouver, BC, Canada, 8–14 December 2019; pp. 8024–8035.

388   34. Kurtzer GM, Sochat V, Bauer MW. Singularity: Scientific containers for mobility of

389      compute. Plos One. 2017;12(5). doi:10.1371/journal.pone.0177459

390   35. Iyer A, Thor M, Haq R, Deasy JO, Apte AP. Deep learning-based auto-segmentation of

391      swallowing and chewing structures in CT. bioRxiv.

392      https://www.biorxiv.org/content/10.1101/772178v1. Published January 1, 2019. Accessed

393      July 27, 2020.