

# 2 **Insights into Predicting Tooth Extraction from Panoramic** 3 **Dental Images: Artificial Intelligence vs. Dentists**

4 **Ila Motmaen<sup>1,†</sup>, Kunpeng Xie<sup>2,3,†</sup>, Leon Schönbrunn<sup>3,2</sup>, Jeff Berens<sup>2,3</sup>, Kim Grunert<sup>2,3</sup>,**  
5 **Anna Maria Plum<sup>2,3</sup>, Johannes Raufeisen<sup>2,3</sup>, André Ferreira<sup>4,5,3,2</sup>, Alexander Hermans<sup>6,7</sup>,**  
6 **Jan Egger<sup>5</sup>, Frank Hölzle<sup>2</sup>, Daniel Truhn<sup>7</sup>, Behrus Puladi<sup>2,3,✉</sup>**

7  
8 1 Department of Oral and Maxillofacial Surgery, University Hospital Knappschaftskrankenhaus Bochum, 44892 Bochum, Germany.

9 2 Department of Oral and Maxillofacial Surgery, University Hospital RWTH Aachen, 52074 Aachen, Germany

10 3 Institute of Medical Informatics, University Hospital RWTH Aachen, 52074 Aachen, Germany

11 4 Centre Algoritmi / LASI, University of Minho, 4710-057 Braga, Portugal

12 5 Institute for Artificial Intelligence in Medicine, Essen University Hospital, 45147 Essen, Germany

13 6 Visual Computing Institute, Computer Science and Natural Sciences, RWTH Aachen University, 52074 Aachen, Germany

14 7 Department of Diagnostic and Interventional Radiology, RWTH Aachen University, 52074 Aachen, Germany

15 † These authors have shared first authorship and contributed equally to this work.

16  
17 ✉ **Corresponding author:**

18 Dr. med. Dr. med. dent. Behrus Puladi

19 E-mail: bpuladi@ukaachen.de

20 Tel.: +49-241-80-38389

21 Fax: +49-241-80-82430

22 Department of Oral and Maxillofacial Surgery & Institute of Medical Informatics,

23 University Hospital RWTH Aachen,

24 Pauwelsstraße 30, 52074 Aachen, Germany

25  
26 **Key words:** Tooth Extraction; Surgery, Oral; Dentistry; Decision Support Techniques; Deep Learning; Artificial  
27 Intelligence.

28  
29 **Abstract:**

30 Objectives: Tooth extraction is one of the most frequently performed medical procedures. The indication is based  
31 on the combination of clinical and radiological examination and individual patient parameters and should be made  
32 with great care. However, determining whether a tooth should be extracted is not always a straightforward decision.  
33 Moreover, visual and cognitive pitfalls in the analysis of radiographs may lead to incorrect decisions. Artificial  
34 intelligence (AI) could be used as a decision support tool to provide a score of tooth extractability.

35 Material and Methods: Using 26,956 single teeth images from 1,184 panoramic radiographs (PANs), we trained a  
36 ResNet50 network to classify teeth as either extraction-worthy or preservable. For this purpose, teeth were cropped  
37 with different margins from PANs and annotated. The usefulness of the AI-based classification as well that of  
38 dentists was evaluated on a test dataset. In addition, the explainability of the best AI model was visualized via a  
39 class activation mapping using CAMERAS.

40 Results: The ROC-AUC for the best AI model to discriminate teeth worthy of preservation was 0.901 with 2%  
41 margin on dental images. In contrast, the average ROC-AUC for dentists was only 0.797. With a 19.1% tooth  
42 extractions prevalence, the AI model's PR-AUC was 0.749, while the human evaluation only reached 0.589.

43 Conclusion: AI models outperform dentists/specialists in predicting tooth extraction based solely on X-ray images,  
44 while the AI performance improves with increasing contextual information.

45 Clinical Relevance: AI could help monitor at-risk teeth and reduce errors in indications for extractions.

## 46 Introduction

47 Tooth extraction is one of the most commonly performed medical measures in the field of general dentistry/  
48 oral and maxillofacial surgery. The decision is based on the patient's records, which include medical history,  
49 clinical evaluation, and radiographs. Given its irreversible impact on the quality of life, the decision of extraction  
50 should be made with great care [1–3]. Certain X-ray signs are pivotal in determining the necessity for tooth  
51 extraction. These signs include the compromised structural integrity of the tooth, significant alveolar bone loss, or  
52 evident root fractures. In addition, massive periapical radiolucency may also suggest the extraction. Advanced  
53 internal or external resorption cases can also be identified on these radiographs, providing a clear indication for  
54 removal of the affected teeth [4].

55 Although indications are made clear in the extraction guidelines [5, 6], the decision-making process is not  
56 always easy for the practitioner in clinical practice [2, 4]. This decision may be confounded by many factors, such  
57 as the dentist's/specialist's own experience, the reliability of the clinical evidence, or even pressure from patients  
58 [5]. The interplay of these different potentially disruptive factors regarding diagnostic decision-making can lead  
59 to misdiagnosis and problematic therapy situations, especially in borderline cases. For example, incorrect tooth  
60 extraction is the third most common cause of tooth loss in periodontally damaged teeth [7].

61 However, leaving teeth that are not worthy of preservation is not an option, as they can cause massive pain [1]  
62 and can even be the starting point for life-threatening lodge abscesses in the head and neck region or cause fatal  
63 endocarditis, which ultimately affects the entire organism [8, 9]. At the same time, every tooth extraction has its  
64 risk of serious complications like persisting root fractures, dry sockets or damage to neighboring teeth. The  
65 indication is, therefore, also always a balancing of different requirements. In general, tooth extraction serves as a  
66 last resort when every other treatment option failed or is not indicated anymore [4].

67 Panoramic radiographs (PANs), commonly used due to easy access and low dosage, are crucial in evaluating  
68 a patient's dental condition, providing insights into the whole dentition and relating structures [10]. However,  
69 accurate and comprehensive interpretation of PANs requires extensive training and considerable clinical  
70 experience. This expertise may not be fully developed in young practitioners, potentially leading to variability in  
71 diagnostic decisions [11]. Furthermore, seasoned practitioners may also be susceptible to cognitive and visual  
72 pitfalls when dealing with challenging cases [12].

73 Deep learning (DL), a subfield of artificial intelligence (AI), has revolutionized the field of medical imaging  
74 by extending the capabilities of human practitioners. These models are trained on vast datasets, allowing them to  
75 recognize patterns and anomalies with superhuman precision [13]. In the context of PANs, the DL models enable  
76 the detection and segmentation of anatomical structures in seconds, with performance improvements being noted  
77 on an ongoing basis [14–18]. Moreover, DL models can identify subtle or complex pathologies that may be  
78 overlooked by the human eye, such as caries, cysts, periodontitis, and periapical lesions. These can be  
79 automatically annotated with high accuracy [19–23]. Such advancements demonstrate the potential of DL to serve  
80 as a powerful tool that enhances diagnostic accuracy and efficiency.

81 Despite these advancements, most research has focused on lesion diagnosis [24–28], with limited exploration  
82 into subsequent clinical decisions like tooth extraction. Furthermore, the model's predictions are often given with  
83 blunt probabilities without any explanation or reasoning process, which is crucial for clinical acceptance and  
84 understanding. Applying explainable DL has the potential to accelerate the decision-making process, resulting in  
85 timely and more effective interventions, ultimately leading to improved patient outcomes [29].

86 The study's main objective is to develop and internally validate a model that can predict the need for tooth  
87 extraction from PANs and compare its performance to dentists/specialists. Furthermore, the effect of contextual  
88 knowledge of teeth on the model's performance and its possible explainability will be visualized.

## 89 Material and Methods

### 90 *Study Design and Patients*

91 The study used retrospective PANs from 2011 to 2021 from patients who underwent tooth extraction at the  
92 Department of Oral and Maxillofacial Surgery of the University Hospital RWTH Aachen. Patients with edentulous  
93 conditions, or without available panoramic radiographs taken within six months post-treatment were excluded.  
94 Additionally, patients with significant artifacts in their preoperative panoramic radiographs that affected the teeth  
95 were also removed from the study cohort.

96 The study was approved by the Ethics Committee of the University Hospital RWTH Aachen (approval number  
97 EK 068/21, chairs: Prof. Dr. G. Schmalzing and PD Dr. R. Hausmann, approval date 25.02.2021) and followed  
98 the MI-CLAIM reporting guideline for the development of AI models [30].

99 *Dataset Preparation*

100 For the study, all PANs were exported in DICOM format from the hospital's picture archiving and  
101 communication system. If a patient had received more than one PAN within six months post-treatment, the last  
102 PAN would be taken as the postoperative image. After the cohort's statistical summary, all PANs were stratified  
103 by patients and converted to PNG format for anonymization purposes.

104 Annotations and labeling of teeth in the preoperative PANs were performed by four investigators (I.M., J.B.,  
105 K.G. and B.P.) using LabelMe [31]. For this purpose, all teeth were marked with a bounding box on the  
106 preoperative image and divided into a preserved and extracted class according to their presence in the postoperative  
107 image (Figure 1). Implants or residual roots were marked in the same way as teeth. For quality control, the  
108 annotated images and labels were then reviewed by two investigators (I.M. and B.P.) for a second round.

109 The bounding boxes were then used to export single tooth images with different margins, as well as their class  
110 (preserved or extracted tooth). Since the distances (in mm) in PANs are not uniform and the teeth themselves have  
111 different sizes, we defined the margins in % of the PAN image height and width. Images were then exported with  
112 margins ranging from -0.5% to 10%, with 0% being the bounding box itself, resulting in 8 datasets. Figure 1  
113 describes the pipeline of the dataset preparation.

114 *Model Development and Validation*

115 The dataset was stratified by patient and randomly divided into a training set (17,874), validation set (4,784),  
116 and test set (4,298) in a 4:1:1 ratio. During training, we apply a random crop to the image, then resize it to 224x224  
117 pixels and perform horizontal flip augmentations to enhance model generalization. Validation and test sets images  
118 are resized to 256x256 pixels and the 224x224 center-crop is extracted.

119 The training was conducted on a high-performance cluster at RWTH Aachen University. We adopted a  
120 ResNet50 model pre-trained on ImageNet. The binary cross-entropy loss was used for our binary classification  
121 tasks. Training spans 50 epochs. The model employs the SGD optimizer with a learning rate of 0.01 and  
122 momentum of 0.9. A learning rate scheduler reduces the learning rate by a factor of 10 every 7 epochs, aiding in  
123 precise model tuning as training progresses (reduce by  $< 1 =$  increase). Model performance was evaluated based  
124 on accuracy and ROC-AUC metrics, with periodic checks to save the best-performing model based on the highest  
125 ROC-AUC achieved. Predictions were made on the test set using these best models, and the predictions were  
126 evaluated and saved. The corresponding code can be found on GitHub (<https://github.com/OMFSdigital/PAN-AI-X>).  
127

128 *Performance of Dentists*

129 In addition, the test images were evaluated by 5 dentists/specialists (A.P., J.B., I.M., K.X., B.P.) with different  
130 levels of experience (dentist in first year to specialist in oral and maxillofacial surgery) to evaluate human  
131 performance. For this purpose, the 4,298 test images (2% margin) were randomly distributed among the  
132 investigators. Each dental image was then given a score between 0 (preserved) and 10 (extracted) to determine the  
133 the likelihood with which a human investigator would recommend a removal of the to.. The 2% margin was chosen  
134 to compare human performance to the DL model with the best performance. To avoid a learning effect between  
135 the annotation in the PANs and the scoring of the individual tooth images by the investigators, there was a 6-month  
136 time delay between initial annotation and scoring.

137 *Model Explainability*

138 To explain the basis of the prediction of the AI models, CAMERAS [32] was used. It uses class activation  
139 mapping to help visualize the regions of the input image that are important for the model's decision-making process  
140 (Figure 4, 5). In our case of binary classification where outcomes are extraction or preservation, CAMERAS  
141 highlights features based on the binary outcome. If the model predicts extraction, it highlights features leading to  
142 this decision; conversely, a prediction of preservation highlights or lacks features, indicating why the preservation  
143 is predicted. The intensity and frequency of these highlights can aid in interpreting model outputs, where more  
144 frequent or intense highlights correlates with a prediction with a higher probability.

145 *Statistical analysis*

146 The statistical analysis was performed in Python (version 3.11.0) using the scikit-learn package (version 1.4.0).  
147 The performance of the AI classifiers and dentists was assessed by using the area under the curve of the receiver  
148 operating characteristic curve (ROC-AUC) and the precision-recall curve (PR-AUC). We then calculated the  
149 maximum Youden's index for each ROC curve and acquired the optimal threshold for the corresponding model.

150 Metrics of accuracy, specificity, precision (syn. positive predictive value), and sensitivity (syn. recall) were  
151 calculated with the thresholds above. The F1 score was calculated from precision and sensitivity. We used a set of  
152 thresholds of 0.3 and 0.7 to plot the confusion matrices with clinically relevant decisions, namely extraction,  
153 monitoring, and preservation.

## 154 **Results**

### 155 *Patients*

156 1,184 patients who met the criteria were selected in this study. The average age of patients was 50.0 years  
157 (range 11 – 99 years), with a standard deviation of 20.3 years. The gender ratio of the cohort was 61:39, with 722  
158 males and 462 females. A total of 26,956 teeth were annotated in 1,184 PANs with bounding boxes and classified  
159 into preservation (21,797) and extraction (5,159). The prevalence of tooth extraction in our dataset was 19.1%,  
160 compared to the majority of 80.9% of preserved teeth. The demographic and clinical characteristics of patients are  
161 described in Table 1.

### 162 *Performance of AI models*

163 Eight different ResNet-50 models were trained on single tooth images with margin settings from -0.5% to 10%.  
164 The performance of models is summarized in Table 2 and Figure 2 based on the thresholds at the maximum  
165 Youden's index. The model with 2% margin setting yielded the best results in both ROC-AUC (0.901) and PR-  
166 AUC (0.749). It also exhibited the best performance in all other metrics except for sensitivity. Shrinking of the  
167 bounding boxes (margin -0.5%) produced worse results in ROC-AUC and PR-AUC than the baseline (margin 0%).  
168 A general increase can be observed in both ROC-AUC and PR-AUC as the margin increases from -0.5% to 2%.  
169 Models with a 5% margin setting have achieved the highest sensitivity (0.835). However, increasing the margin  
170 further to 10% reduced both ROC-AUC and PR-AUC. In confusion matrices, with thresholds of 0.3 and 0.7 for  
171 monitoring, the 2% margin model had the least cases of false positive (53). The model with 3% margin had the  
172 highest accuracy (3455/4298).

### 173 *Performance of Dentists*

174 In contrast, the human assessment (average of 5 dentists/specialists) had a lower performance based on the 2%  
175 dental images compared to the AI models. The ROC-AUC was only 0.797 or PR-AUC of 0.589. This is also  
176 reflected by the confusion matrices where human have the most false positives (131) and lowest accuracy  
177 (3085/4298).

### 178 *Explainability*

179 Figure 4 and 5 shows the activation map of the *extracted and preserved* predictions generated by CAMERAS  
180 with a 2% margin setting. In extraction cases, the model focused on the areas where roots are exposed in low  
181 density regions and crowns are buried in bone. In preservation cases, on the other hand, alveolar ridge and  
182 periapical regions were the most relevant.

## 183 **Discussion**

184 In this study, to our knowledge, we present the first clinical prediction model using DL to make a  
185 recommendation about teeth extractions. The main results of the study are, 1) the best model achieved a ROC-  
186 AUC of 0.901 with a PR-AUC of 0.749; 2) outperforming dentists/specialists, who on average achieved a ROC-  
187 AUC of 0.797 with a PR-AUC of 0.589; 3) additional contextual information through wide margins around the  
188 tooth led to a better prediction; 5) the visual explainability of the prediction for tooth extraction or preservation  
189 was comprehensible.

190 Decision aids are a useful tool, for example in healthcare, to reduce the dentists' workload, as suggestions  
191 calculated by algorithms can contribute to the final decision-making or diagnosis and significantly speed up this  
192 process [33]. Similarly, decision aids can be used as an objective perspective, especially in borderline cases where  
193 otherwise subjective approaches are applied by the clinicians alone [33, 34]. In this regard, work in the medical  
194 field has already been done on identifying pathologies in medical imaging like X-ray scans. One of the first  
195 applications used for detection was in 1995 to detect nodules in X-rays of the lungs [35]. Another object detection  
196 algorithm was developed to detect and classify several entities in chest X-rays like cardiomegaly, calcified

197 granulomas, catheters, surgical instruments or thoracic vertebrae [36]. The emergence of convolutional neural  
198 networks / DL more than a decade ago opened up completely new possibilities [37].

199 One recent application is described by Yoo et al. who proposed a DL model (VGG16 pre-trained on ImageNet)  
200 to predict the difficulty of extracting a mandibular third molar from PANs [38]. The model was trained to predict  
201 the difficulty of mandibular third molar extraction in terms of depth, ramal relationship, and angulation. The  
202 accuracies of the model for different difficulty parameters (depth, ramal relationship, angulation) were found to  
203 be 78.9%, 82.0%, and 90.2%, respectively. Yet the model was made to predict the difficulty rather than the  
204 necessity of the extraction.

205 In our study, we used a residual neural network (ResNet-50) pretrained on ImageNet for the development of  
206 our clinical prediction model. Compared to other convolutional neural networks, a ResNet is characterized by so-  
207 called residual skip connections, which add inputs to outputs of small blocks of layers in the network. These skip  
208 connections improve the gradient flow during training and significantly improve the performance of very deep  
209 networks [39]. An outstanding strength of our model was its ability to classify teeth not worthy of preservation  
210 across multiple indications, such as extractions for orthodontic space, misplaced wisdom teeth, caries-destroyed  
211 teeth, periodontally compromised teeth or teeth from mixed dentition. Equally noteworthy was the reliable  
212 classification even in radiographs with more difficult classification conditions, such as anatomical superimposition  
213 effects.

214 Yet, evidence-based medicine encourages decisions based on patient-specific clinical evidence. However, DL  
215 models often provide blunt predictions without any explanation [40]. This results in a low acceptance among  
216 practitioners of these predictions due to the lack of visible evidence [29]. To address this problem, class activation  
217 map offers a solution to visualize and highlight the critical area of the image where the predictions are made [41,  
218 32]. In the case of the caries classification task in the study of Vinayahalingam et al., areas that leads to the  
219 classification by DL model were be highlighted [42]. Such visual prompts can then correlate with established  
220 dental knowledge of the practitioners, which in turn explains the classification or recommendations.

221 We used CAMERAS, which, in contrast to methods such as GCAM or NormGrad, provides high-resolution  
222 mapping for ResNet and, thus, new insights into the explainability of DL methods [32]. The explainability can be  
223 illustrated using the examples of extracted teeth (Figure 4) and preserved teeth (Figure 5), including their prediction  
224 probability. In the case of healthy teeth, for example, this leads to activation of the bone, whereas in the case of  
225 root remnants this leads directly to the root itself. In addition to the recommendation, this activation map could  
226 also be offered directly to the dentist.

227 Interestingly, however, it can also be seen that due to the additional context information provided by the  
228 extended margin (2%) in Figures 4 and 5, neighboring root residues are also included in the classification and may  
229 possibly lead to a misclassification. This could be remedied in the future by more modern architectures that  
230 consider the entire PAN instead of individual image sections with a tooth and the adjacent bone.

231 Besides these technical aspects, the question arises as to how such a model could be translated into practice.  
232 An important challenge is that DL models fall under regulatory requirements such as FDA / Medical Device  
233 Regulation (MDR) as medical software. This means that the models developed in research cannot simply be  
234 applied in clinical encounters [43]. An important step here would be the external validation of the developed model  
235 [44]. At our department, the prevalence of tooth extraction was 19.1% (Table 1). This is influenced by the present  
236 population with its socioeconomic status, but certainly to some extent also to the treating specialty (conservative  
237 dentistry, prosthodontics, orthodontics, oral and maxillofacial surgery) has an impact that cannot be dismissed out  
238 of hand, as well as the pre-selection of cases. This could represent a bias if the model is applied elsewhere. On the  
239 other hand, it could be argued that the reasons for tooth extraction are universal worldwide [3, 45]. Periapical  
240 radiolucency or deep caries are not treated much differently around the world.

241 Clinical prediction models such as ours usually divide cases into two treatment recommendations based on a  
242 single threshold (perceive / extract). For an actual application scenario, however, the question of design is  
243 particularly crucial for optimal clinical usefulness [46]. This could involve dividing teeth into three groups based  
244 on two thresholds. Using a low threshold (with a high negative predictive value) to distinguish teeth that are  
245 definitely worth preserving from suspect teeth. Another higher threshold (with a high positive predictive value)  
246 could separate suspect teeth from definitely not preservable ones. The suspect teeth could then be monitored  
247 closely, while the healthy teeth would be ignored, and the decayed teeth would be extracted. An example for this  
248 approach is shown in Figure 3.

249 However, a major limitation of our results is that our model does not include clinical information (pain, tooth  
250 vitality, course of disease, diagnosis). On the one hand, this is impressive because a high level of accuracy has  
251 been achieved despite the lack of any clinical information surpassing humans. Nevertheless, in a real clinical  
252 setting this information would be available and should be used. In the future, multimodal AI models could be used  
253 to process additional clinical information and improve prediction.

254 Another limitation is that there was a maximum period of 6 months between pre- and postoperative PAN.  
255 Usually, significant changes are visible during this period, but the causes for the extraction may not have been  
256 visible on the preoperative image used in some cases, but only shortly before the extraction itself (such as the  
257 involvement of teeth in a mandibular fracture).

## 258 **Conclusion**

259 In summary, our study presented the first AI model to our knowledge to assist dentists/specialists in making  
260 tooth extraction decisions based on radiographs alone. The developed AI models outperform humans, with AI  
261 performance improving as contextual information increases. Future models may integrate clinical data. This study  
262 provides a good foundation for further research in this area. In the future, AI could help monitor at-risk teeth and  
263 reduce errors in indications for extraction. By providing a class activation map, clinicians could be able to  
264 understand and verify the AI decision.

## 265 **Declaration**

266 **Author Contributions:** Conceptualization, B.P., I.M., and K.X.; methodology, I.M., L.S., J.R., A.H., B.P., and  
267 J.E.; software, L.S., I.M. and J.R.; validation, K.X., B.P., I.M., J.B., K.G. and A.P.; formal analysis, K.X., B.P.,  
268 A.F., A.H., F.H. and D.T.; investigation, I.M., L.S., K.X. and B.P.; resources, B.P., F.H. and D.T.; data curation,  
269 I.M., K.X. and L.S.; writing—original draft preparation, K.X., B.P. and I.M.; writing—review and editing, B.P.,  
270 K.X., I.M., L.S., J.B., K.G., A.P., J.R., A.F., A.H., J.E., F.H. and D.T.; visualization, B.P., L.S. and K.X.;  
271 supervision, B.P.; project administration, B.P.; All authors have read and agreed to the published version of the  
272 manuscript.

273  
274 **Funding:** André Ferreira was funded by the Advanced Research Opportunities Program (AROP) of RWTH  
275 Aachen University. Behrus Puladi was funded by the Medical Faculty of RWTH Aachen University as part of the  
276 Clinician Scientist Program.

277  
278 **Institutional Review Board Statement:** The study approved by the Institutional Review Board (or Ethics  
279 Committee) of University Hospital RWTH Aachen (approval number EK 068/21, chairs: Prof. Dr. G. Schmalzing  
280 and PD Dr. R. Hausmann, approval date 25.02.2021).

281  
282 **Informed Consent Statement:** Informed consent was obtained from all subjects involved in the study.

283  
284 **Code Availability Statement:** All code was implemented in Python. The source code, including the model  
285 weights, is available on GitHub (<https://github.com/OMFSdigital/PAN-AI-X>).

286  
287 **Data Availability Statement:** The data presented in this study are available upon reasonable request from the  
288 corresponding author.

289  
290 **Acknowledgments:** Computations were performed with computing resources granted by RWTH Aachen  
291 University under project rwth1410.

292  
293 **Conflicts of Interest:** The authors declare no conflict of interest.

294

## Reference

- 296 1. Gilbert GH, Meng X, Duncan RP et al. (2004) Incidence of tooth loss and prosthodontic dental care: effect  
297 on chewing difficulty onset, a component of oral health-related quality of life. *J Am Geriatr Soc* 52:880–  
298 885. <https://doi.org/10.1111/j.1532-5415.2004.52253.x>
- 299 2. Avila G, Galindo-Moreno P, Soehren S et al. (2009) A novel decision-making process for tooth retention  
300 or extraction. *Journal of Periodontology* 80:476–491. <https://doi.org/10.1902/jop.2009.080454>
- 301 3. Broers DLM, Dubois L, Lange J de et al. (2022) Reasons for Tooth Removal in Adults: A Systematic  
302 Review. *Int Dent J* 72:52–57. <https://doi.org/10.1016/j.identj.2021.01.011>
- 303 4. Sambrook PJ, Goss AN (2018) Contemporary exodontia. *Australian Dental Journal* 63 Suppl 1:S11-S18.  
304 <https://doi.org/10.1111/adj.12586>
- 305 5. Broers DLM, Brands WG, Welie JVM et al. (2010) Deciding about patients' requests for extraction:  
306 ethical and legal guidelines. *J Am Dent Assoc* 141:195–203.  
307 <https://doi.org/10.14219/jada.archive.2010.0139>
- 308 6. Alkhalifah S, Alkandari H, Sharma PN et al. (2017) Treatment of Cracked Teeth. *J Endod* 43:1579–1586.  
309 <https://doi.org/10.1016/j.joen.2017.03.029>
- 310 7. Lundgren D, Rylander H, Laurell L (2008) To save or to extract, that is the question. Natural teeth or  
311 dental implants in periodontitis-susceptible patients: clinical decision-making and treatment strategies  
312 exemplified with patient case presentations. *Periodontol* 2000 47:27–50. <https://doi.org/10.1111/j.1600-0757.2007.00239.x>
- 313
- 314 8. Hansen BW, Ryndin S, Mullen KM (2020) Infections of Deep Neck Spaces. *Semin Ultrasound CT MR*  
315 41:74–84. <https://doi.org/10.1053/j.sult.2019.10.001>
- 316 9. Nomura R, Matayoshi S, Otsugu M et al. (2020) Contribution of Severe Dental Caries Induced by  
317 *Streptococcus mutans* to the Pathogenicity of Infective Endocarditis. *Infect Immun* 88.  
318 <https://doi.org/10.1128/IAI.00897-19>
- 319 10. Perschbacher S (2012) Interpretation of panoramic radiographs. *Australian Dental Journal* 57 Suppl 1:40–  
320 45. <https://doi.org/10.1111/j.1834-7819.2011.01655.x>
- 321 11. Geibel M-A, Carstens S, Braisch U et al. (2017) Radiographic diagnosis of proximal caries-influence of  
322 experience and gender of the dental staff. *Clin Oral Invest* 21:2761–2770. <https://doi.org/10.1007/s00784-017-2078-2>
- 323
- 324 12. Aeffner F, Wilson K, Martin NT et al. (2017) The Gold Standard Paradox in Digital Image Analysis:  
325 Manual Versus Automated Scoring as Ground Truth. *Arch Pathol Lab Med* 141:1267–1275.  
326 <https://doi.org/10.5858/arpa.2016-0386-RA>
- 327 13. Çallı E, Sogancioglu E, van Ginneken B et al. (2021) Deep learning for chest X-ray analysis: A survey.  
328 *Med Image Anal* 72:102125. <https://doi.org/10.1016/j.media.2021.102125>
- 329 14. Lee J-H, Han S-S, Kim YH et al. (2020) Application of a fully deep convolutional neural network to the  
330 automation of tooth segmentation on panoramic radiographs. *Oral Surg Oral Med Oral Pathol Oral Radiol*  
331 129:635–642. <https://doi.org/10.1016/j.oooo.2019.11.007>
- 332 15. Bilgir E, Bayrakdar İŞ, Çelik Ö et al. (2021) An artificial intelligence approach to automatic tooth  
333 detection and numbering in panoramic radiographs. *BMC Med Imaging* 21:124.  
334 <https://doi.org/10.1186/s12880-021-00656-7>
- 335 16. Cha J-Y, Yoon H-I, Yeo I-S et al. (2021) Panoptic Segmentation on Panoramic Radiographs: Deep  
336 Learning-Based Segmentation of Various Structures Including Maxillary Sinus and Mandibular Canal. *J*  
337 *Clin Med* 10. <https://doi.org/10.3390/jcm10122577>
- 338 17. Vinayahalingam S, Goey R-S, Kempers S et al. (2021) Automated chart filing on panoramic radiographs  
339 using deep learning. *J Dent* 115:103864. <https://doi.org/10.1016/j.jdent.2021.103864>
- 340 18. Jeon KJ, Choi H, Lee C et al. (2023) Automatic diagnosis of true proximity between the mandibular canal  
341 and the third molar on panoramic radiographs using deep learning. *Sci Rep* 13:22022.  
342 <https://doi.org/10.1038/s41598-023-49512-4>
- 343 19. Yang H, Jo E, Kim HJ et al. (2020) Deep Learning for Automated Detection of Cyst and Tumors of the  
344 Jaw in Panoramic Radiographs. *J Clin Med* 9. <https://doi.org/10.3390/jcm9061839>
- 345 20. Lian L, Zhu T, Zhu F et al. (2021) Deep Learning for Caries Detection and Classification. *Diagnostics*  
346 (Basel) 11. <https://doi.org/10.3390/diagnostics11091672>
- 347 21. Watanabe H, Arijji Y, Fukuda M et al. (2021) Deep learning object detection of maxillary cyst-like lesions  
348 on panoramic radiographs: preliminary study. *Oral Radiol* 37:487–493. <https://doi.org/10.1007/s11282-020-00485-4>
- 349

- 350 22. Endres MG, Hillen F, Salloumis M et al. (2020) Development of a Deep Learning Algorithm for Periapical  
351 Disease Detection in Dental Radiographs. *Diagnostics (Basel)* 10.  
352 <https://doi.org/10.3390/diagnostics10060430>
- 353 23. Guler Ayyildiz B, Karakis R, Terzioglu B et al. (2024) Comparison of deep learning methods for the  
354 radiographic detection of patients with different periodontitis stages. *Dentomaxillofac Radiol* 53:32–42.  
355 <https://doi.org/10.1093/dmfr/twad003>
- 356 24. Liu Z, Liu J, Zhou Z et al. (2021) Differential diagnosis of ameloblastoma and odontogenic keratocyst by  
357 machine learning of panoramic radiographs. *Int J Comput Assist Radiol Surg* 16:415–422.  
358 <https://doi.org/10.1007/s11548-021-02309-0>
- 359 25. Kwon O, Yong T-H, Kang S-R et al. (2020) Automatic diagnosis for cysts and tumors of both jaws on  
360 panoramic radiographs using a deep convolution neural network. *Dentomaxillofac Radiol* 49:20200185.  
361 <https://doi.org/10.1259/dmfr.20200185>
- 362 26. Ekert T, Krois J, Meinhold L et al. (2019) Deep Learning for the Radiographic Detection of Apical  
363 Lesions. *J Endod* 45:917-922.e5. <https://doi.org/10.1016/j.joen.2019.03.016>
- 364 27. Sukegawa S, Fujimura A, Taguchi A et al. (2022) Identification of osteoporosis using ensemble deep  
365 learning model with panoramic radiographs and clinical covariates. *Sci Rep* 12:6088.  
366 <https://doi.org/10.1038/s41598-022-10150-x>
- 367 28. Arijji Y, Yanashita Y, Kutsuna S et al. (2019) Automatic detection and classification of radiolucent lesions  
368 in the mandible on panoramic radiographs using a deep learning object detection technique. *Oral Surg Oral*  
369 *Med Oral Pathol Oral Radiol* 128:424–430. <https://doi.org/10.1016/j.oooo.2019.05.014>
- 370 29. Loh HW, Ooi CP, Seoni S et al. (2022) Application of explainable artificial intelligence for healthcare: A  
371 systematic review of the last decade (2011-2022). *Computer Methods and Programs in Biomedicine*  
372 226:107161. <https://doi.org/10.1016/j.cmpb.2022.107161>
- 373 30. Norgeot B, Quer G, Beaulieu-Jones BK et al. (2020) Minimum information about clinical artificial  
374 intelligence modeling: the MI-CLAIM checklist. *Nat Med* 26:1320–1324. <https://doi.org/10.1038/s41591-020-1041-y>
- 375 31. Kentaro Wada, mpitid, Martijn Buijs et al. (2021) `wkentaro/labelme: v4.6.0`.  
376 <https://doi.org/10.5281/zenodo.5711226>
- 377 32. Jalwana MAAK, Akhtar N, Bennamoun M et al. (2021) CAMERAS: Enhanced Resolution And Sanity  
378 preserving Class Activation Mapping for image saliency. In: 2021 IEEE/CVF Conference on Computer  
379 Vision and Pattern Recognition (CVPR). IEEE, pp 16322–16331
- 380 33. Razzak MI, Naz S, Zaib A (2018) Deep Learning for Medical Image Processing: Overview, Challenges  
381 and the Future. In: Dey N, Ashour AS, Borra S (eds) *Classification in BioApps: Automation of Decision*  
382 *Making*, vol 26. Springer International Publishing, Cham, pp 323–350
- 383 34. Bini SA (2018) Artificial Intelligence, Machine Learning, Deep Learning, and Cognitive Computing: What  
384 Do These Terms Mean and How Will They Impact Health Care? *The Journal of Arthroplasty* 33:2358–  
385 2361. <https://doi.org/10.1016/j.arth.2018.02.067>
- 386 35. Lo SB, Lou SA, Lin JS et al. (1995) Artificial convolution neural network techniques and applications for  
387 lung nodule detection. *IEEE Trans Med Imaging* 14:711–718. <https://doi.org/10.1109/42.476112>
- 388 36. Shin H-C, Roberts K, Lu L et al. (2016) Learning to Read Chest X-Rays: Recurrent Neural Cascade Model  
389 for Automated Image Annotation. In: 2016 IEEE Conference on Computer Vision and Pattern Recognition  
390 (CVPR). IEEE, pp 2497–2506
- 391 37. Corbella S, Srinivas S, Cabitza F (2021) Applications of deep learning in dentistry. *Oral Surg Oral Med*  
392 *Oral Pathol Oral Radiol* 132:225–238. <https://doi.org/10.1016/j.oooo.2020.11.003>
- 393 38. Yoo J-H, Yeom H-G, Shin W et al. (2021) Deep learning based prediction of extraction difficulty for  
394 mandibular third molars. *Sci Rep* 11:1954. <https://doi.org/10.1038/s41598-021-81449-4>
- 395 39. He K, Zhang X, Ren S et al. Deep Residual Learning for Image Recognition.  
396 <https://doi.org/10.48550/arXiv.1512.03385>
- 397 40. Taylor J, Fenner J (2019) The challenge of clinical adoption-the insurmountable obstacle that will stop  
398 machine learning? *BJR Open* 1:20180017. <https://doi.org/10.1259/bjro.20180017>
- 399 41. Viton F, Elbattah M, Guerin J-L et al. (2020) Heatmaps for Visual Explainability of CNN-Based  
400 Predictions for Multivariate Time Series with Application to Healthcare. In: 2020 IEEE International  
401 Conference on Healthcare Informatics (ICHI). IEEE, pp 1–8
- 402 42. Vinayahalingam S, Kempers S, Limon L et al. (2021) Classification of caries in third molars on panoramic  
403 radiographs using deep learning. *Sci Rep* 11:12609. <https://doi.org/10.1038/s41598-021-92121-2>
- 404 43. Karnik K (2014) FDA regulation of clinical decision support software. *J Law Biosci* 1:202–208.  
405 <https://doi.org/10.1093/jlb/lisu004>
- 406



- 407 44. Beckers R, Kwade Z, Zanca F (2021) The EU medical device regulation: Implications for artificial  
408 intelligence-based medical device software in medical physics. *Phys Med* 83:1–8.  
409 <https://doi.org/10.1016/j.ejmp.2021.02.011>
- 410 45. Passarelli PC, Pagnoni S, Piccirillo GB et al. (2020) Reasons for Tooth Extractions and Related Risk  
411 Factors in Adult Patients: A Cohort Study. *International Journal of Environmental Research and Public*  
412 *Health* 17:2575. <https://doi.org/10.3390/ijerph17072575>
- 413 46. Steyerberg EW (2019) *Clinical Prediction Models: A practical approach to development, validation, and*  
414 *updating*, Second edition. Springer eBooks Mathematics and Statistics. Springer International Publishing,  
415 Cham  
416

417 **Tables**

418 *Table 1*

Parameters		Training	Validation	Testing	Total
Patients	Total	787	206	191	1,184
	(%)	(66.5%)	(17.4%)	(16.1%)	(100%)
Age	Mean	50.5	48.6	49.5	50.0
	(SD)	(20.3)	(21.6)	(20.7)	(20.6)
	Range	12 – 92	27 – 66	12 – 99	11 – 99
Gender	Female	310	78	74	462
	(%)	(39.4%)	(37.9%)	(38.7%)	(39.0%)
	Male	477	128	117	722
	(%)	(60.6%)	(62.1%)	(61.3%)	(61.0%)
Teeth	Extracted	3,410	876	873	5,159
	(%)	(66.1%)	(17.0%)	(16.9%)	(19.1%)
	Preserved	14,464	3,908	3,425	21,797
	(%)	(66.4%)	(17.9%)	(15.7%)	(80.9%)
	Total	17,874	4,784	4,298	26,956
	(%)	(66.3%)	(17.7%)	(15.9%)	(100%)

419 **Table 1.** Demographic and dental characteristics of the patients and distribution across training, validation, and  
420 testing datasets.

421 *Table 2*

Predictor	Margin	Accuracy	Sensitivity	Specificity	Precision	F1_Score	ROC-AUC	PR-AUC
AI (ResNet 50)	-0.5%	0.775	0.700	0.794	0.464	0.558	0.815	0.590
	0%	0.768	0.781	0.765	0.459	0.578	0.852	0.657
	0.5%	0.778	0.772	0.780	0.472	0.586	0.852	0.650
	1%	0.813	0.716	0.837	0.529	0.608	0.864	0.693
	<b>2%</b>	<b>0.834</b>	0.797	<b>0.843</b>	<b>0.564</b>	<b>0.661</b>	<b>0.901</b>	<b>0.749</b>
	3%	0.812	0.805	0.814	0.525	0.635	0.895	0.743
	5%	0.805	<b>0.835</b>	0.797	0.512	0.634	0.899	0.741
	10%	0.817	0.816	0.818	0.533	0.645	0.890	0.727
Human (Avg.)	2%	0.775	0.674	0.800	0.462	0.548	0.797	0.589

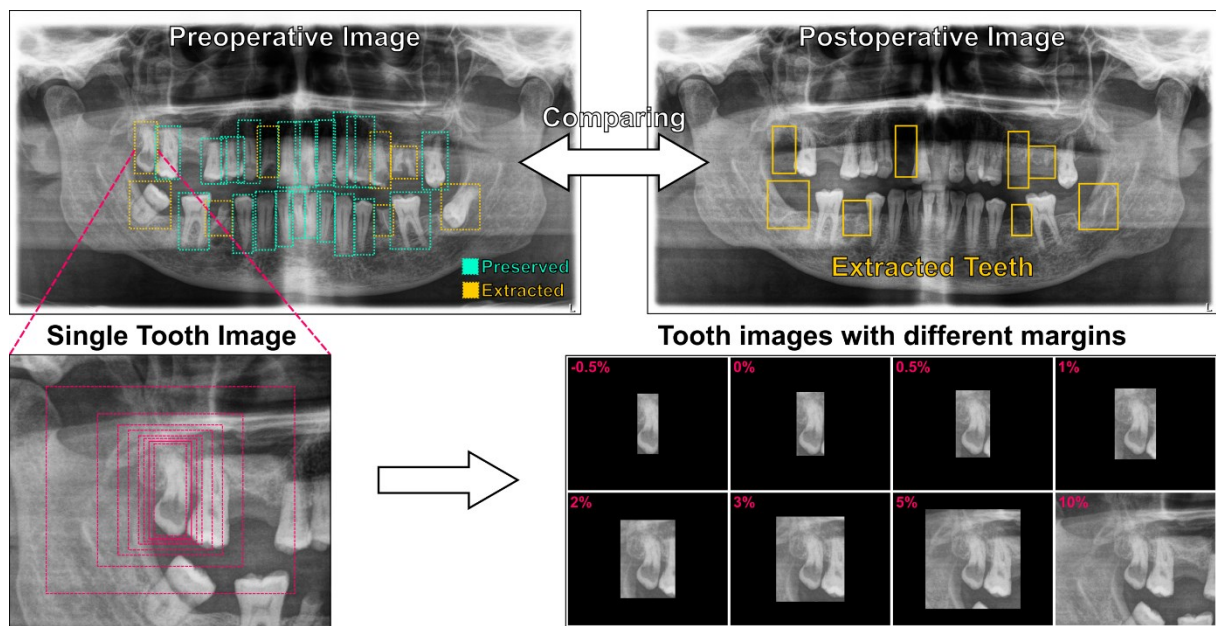
422 **Table 2.** Performance at Youden's index of AI models with different margin settings as well as human  
423 performance..

424

425

426 **Figures**

427 *Figure 1*

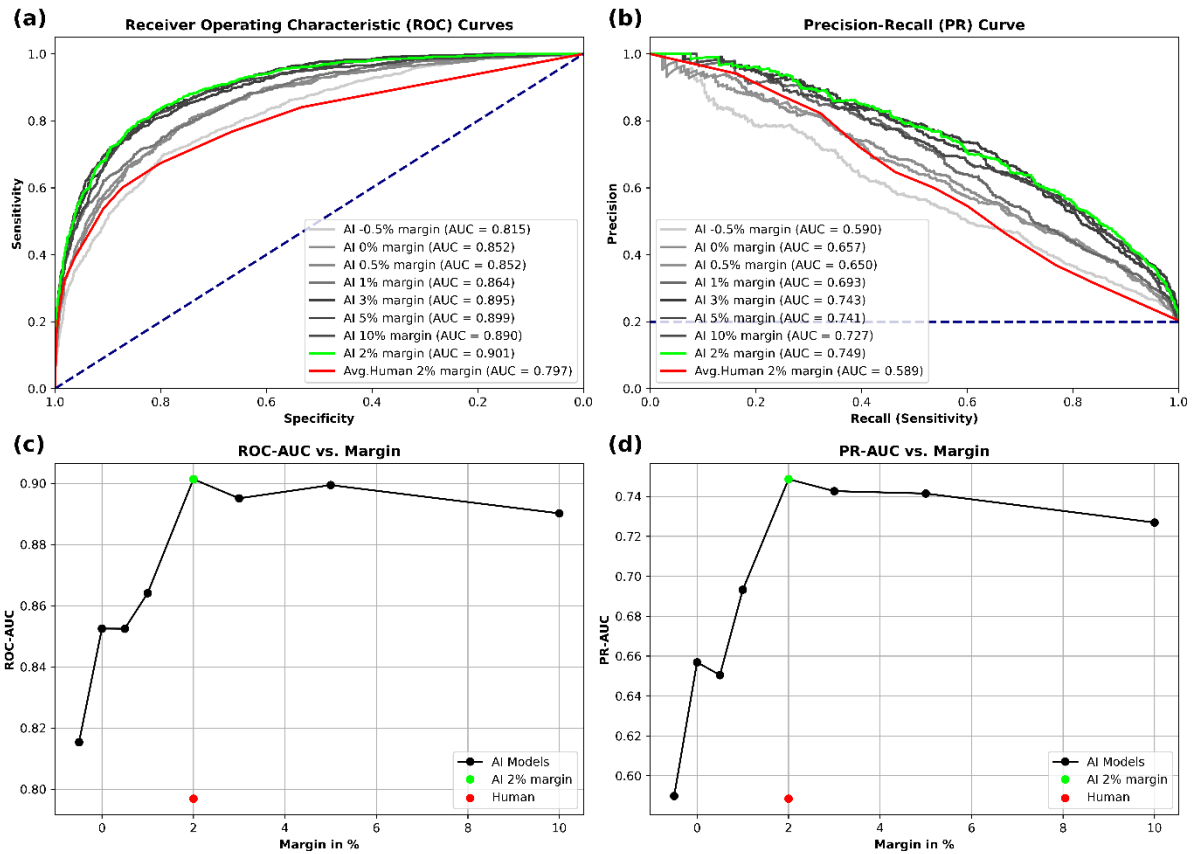


428

429 **Figure 1.** Pipeline to prepare the dataset. Panoramic radiographs from the same patient were compared and  
430 annotations of teeth were made on the preoperative image with bounding boxes and labeled as preserved (green)  
431 or extracted (yellow). Different margin factors were used to resize the bounding boxes (red) in width and height.  
432 Teeth images were then cropped from the original image with margins (-0.5% to 10%).

433

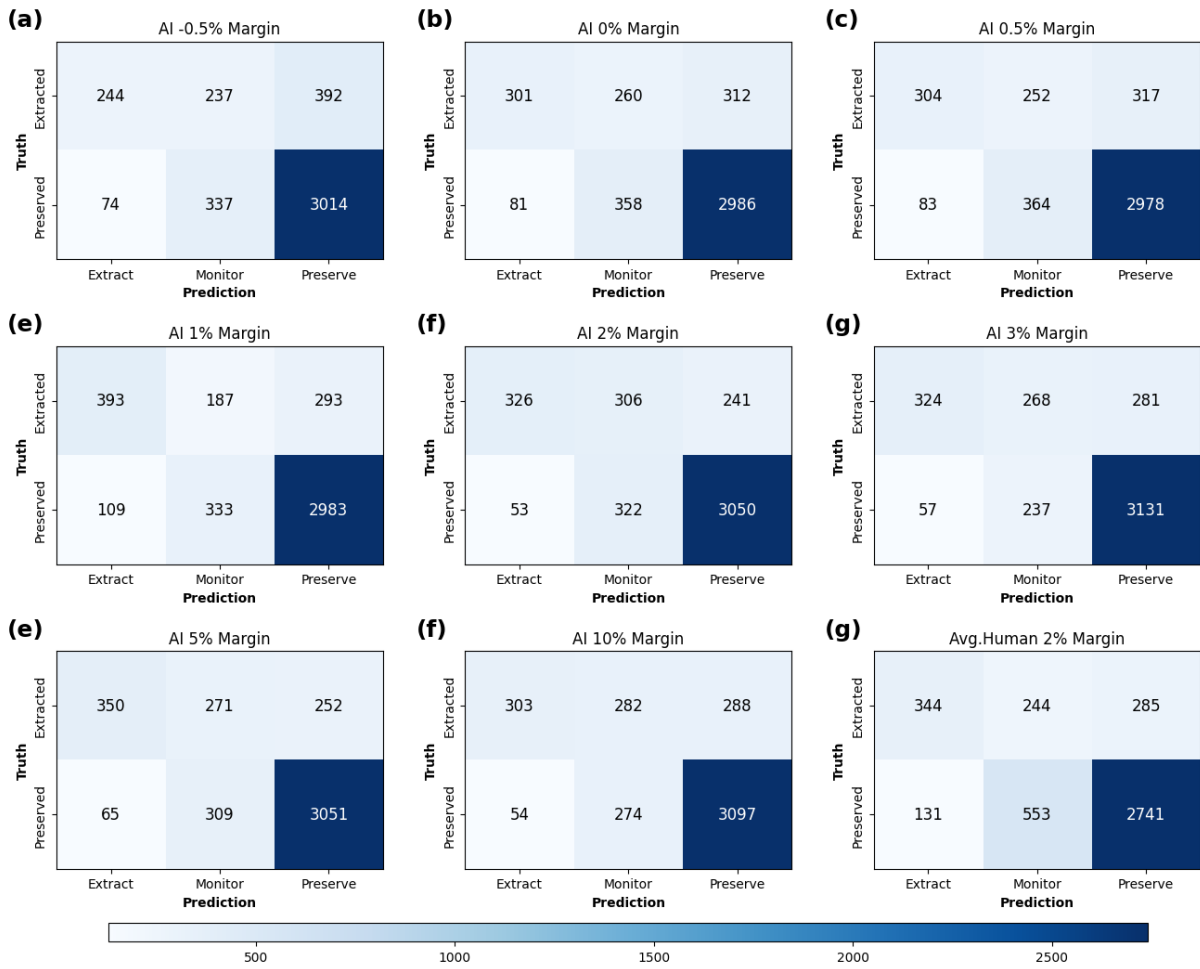
434



436

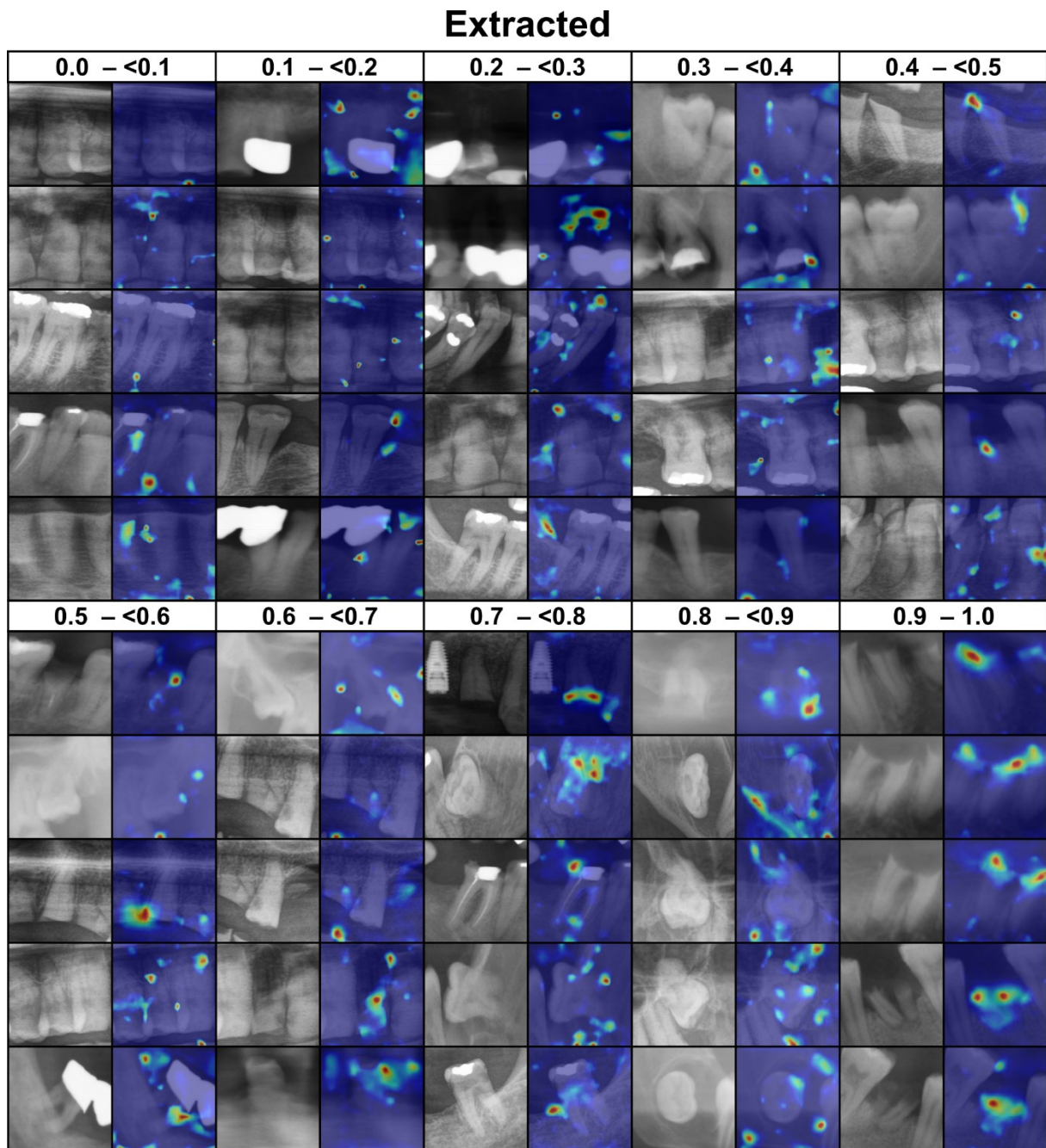
437 **Figure 2.** (a) ROC curves and (b) PPV-Sensitivity curves of models with different margin settings. The 2%  
 438 margin model performed best in both ROC-AUC (0.901) and PR-AUC (0.749), the average human performance  
 439 was ROC-AUC (0.797) and PR-AUC (0.589). Relationship between ROC-AUC and margins is displayed in (c).  
 440 Relationship between PR-AP and margins is displayed in (d). A steep increase observed for both metrics from -  
 441 0.5% to 2% margin and slightly drop from 5% to 10% margin.

442



444

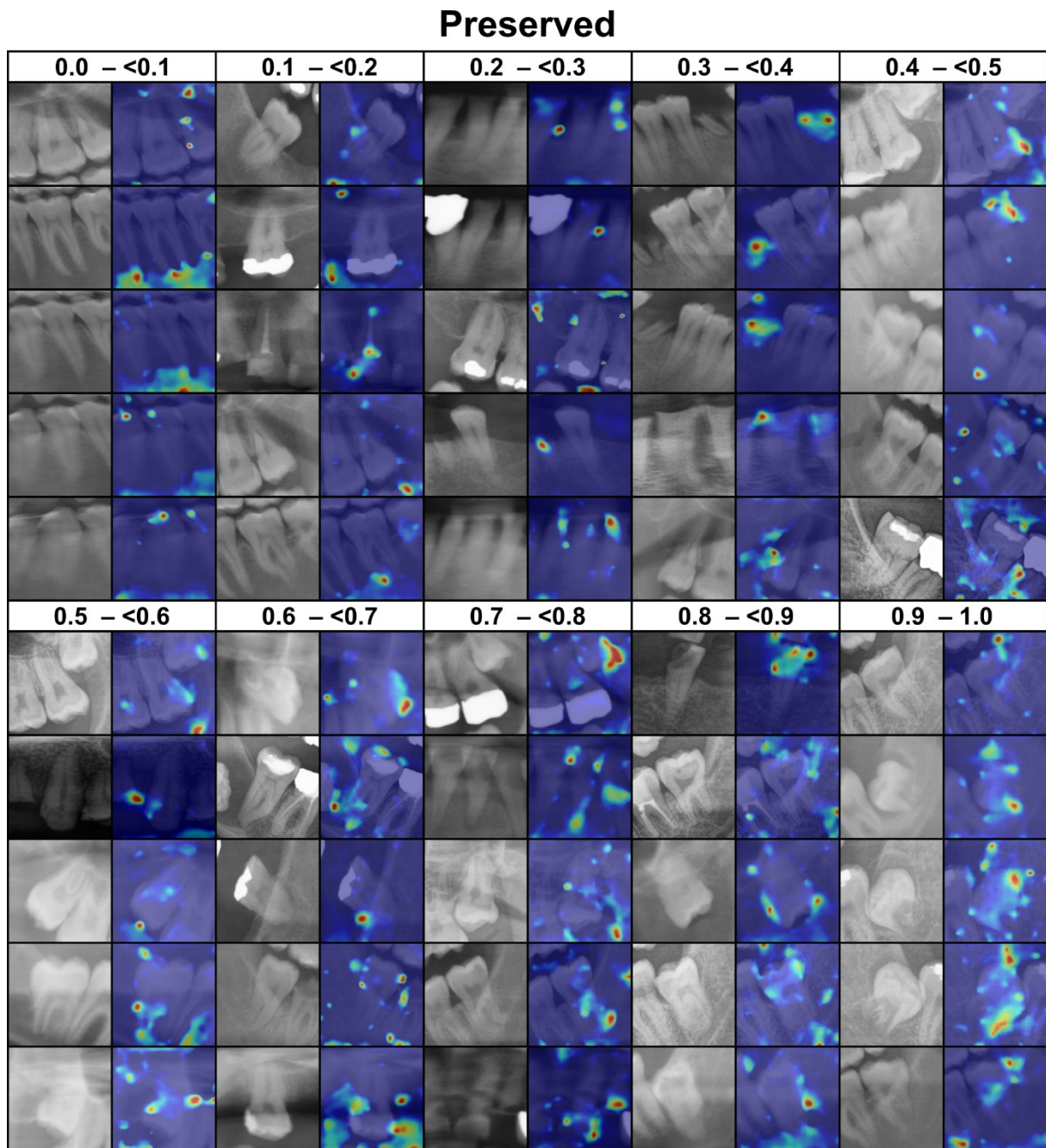
445 **Figure 3.** Confusion matrices showing prediction results. The results from AI models (a)–(f) and dentists (g) with  
 446 different margins were split into 3 decisions, namely extraction, monitoring, and preservation. Teeth with  
 447 prediction probabilities from 0.3 to 0.7 were recommended to “Monitor”. Teeth with prediction probabilities below  
 448 0.3 were recommended to “Extract” while above 0.7 to “Preserve”. True labels were marked in y-axis.



450

451 **Figure 4.** Activation gradient heatmap generated by CAMERAS for extracted teeth with a margin of 2%. The  
 452 probability (0 to 1, where 0 indicates preservation and 1 indicates extraction) of the prediction is shown in the first  
 453 row. The left image in each column is the tooth image used for the prediction, the right image is the class activation  
 454 mapping with CAMERAS. Blue indicates no activation and red indicates strong activation. Green and yellow are  
 455 in between.

456



458

459 **Figure 5.** Activation gradient heatmap generated by CAMERAS for preserved teeth with a margin of 2%. The  
 460 probability (0 to 1, where 0 indicates preservation and 1 indicates extraction) of the prediction is shown in the first  
 461 row. The left image in each column is the tooth image used for the prediction, the right image is the class activation  
 462 mapping with CAMERAS. Blue indicates no activation and red indicates strong activation. Green and yellow are  
 463 in between.

464

465

466