

1 **AIBx, artificial intelligence model to risk stratify thyroid nodules**

2

3 Johnson Thomas MD<sup>1</sup>, Tracy Haertling<sup>2</sup>

4

5

6 <sup>1</sup>Department of Endocrinology, Mercy Hospital, Springfield, Missouri. <sup>2</sup>Mercy Research

7 Springfield, Missouri

8

9 Johnson Thomas MD, FACE, 3231 S National Ave, Suite 440, Springfield, Missouri,

10 USA 65807 [johnson.thomas@mercy.net](mailto:johnson.thomas@mercy.net)

11

12 Tracy Haertling, DO, MS, Manager, Investigator-Initiated Trials - Mercy Research

13 Mercy, 3231 S. National Ave., Suite 210 | Springfield, Missouri 65807

14 [Tracy.Haertling@Mercy.net](mailto:Tracy.Haertling@Mercy.net)

15

16 **\*\*Note:** A portion of the data in this manuscript was reported as an oral abstract at the

17 89<sup>th</sup> annual meeting of the American Thyroid Association, Chicago, USA, October 30 –

18 Nov 3, 2019.

19 Author contributions: Concept, coding, image collection, statistical analysis performed

20 by JT. TH helped in creating the study protocol and edited the article.

21

22 Running Title: AIBx, AI thyroid nodule risk stratification

23 Keywords: Artificial intelligence, Image similarity, thyroid nodule, thyroid cancer

24 **Peer reviewed article can be viewed at official Thyroid Journal website -**

25 **<https://www.liebertpub.com/doi/full/10.1089/thy.2019.0752>**

26 **Abstract**

27 Background

28 Current classification systems for thyroid nodules are very subjective. Artificial  
29 intelligence (AI) algorithms have been used to decrease subjectivity in medical image  
30 interpretation. 1 out of 2 women over the age of 50 may have a thyroid nodule and at  
31 present the only way to exclude malignancy is through invasive procedures. Hence,  
32 there exists a need for noninvasive objective classification of thyroid nodules. Some  
33 cancers have benign appearance on ultrasonogram. Hence, we decided to create an  
34 image similarity algorithm rather than image classification algorithm.

35

36 Methods

37 Ultrasound images of thyroid nodules from patients who underwent either biopsy  
38 or thyroid surgery from February of 2012 through February of 2017 in our institution  
39 were used to create AI models. Nodules were excluded if there was no definitive  
40 diagnosis of benignity or malignancy. 482 nodules met the inclusion criteria and all  
41 available images from these nodules were used to create the AI models. Later, these AI  
42 models were used to test 103 thyroid nodules which underwent biopsy or surgery from  
43 March of 2017 through July of 2018.

44

45 Results

46 Negative predictive value of the image similarity model was 93.2%. Sensitivity,  
47 specificity, positive predictive value and accuracy of the model was 87.8%, 78.5%,  
48 65.9% and 81.5% respectively.

49

50 Conclusion

51 When compared to published results of ACR TIRADS and ATA classification  
52 system, our image similarity model had comparable negative predictive value with better  
53 sensitivity specificity and positive predictive value. By using image similarity AI models,  
54 we can eliminate subjectivity and decrease the number of unnecessary biopsies. Using  
55 image similarity AI model, we were able to create an explainable AI model which  
56 increases physician's confidence in the predictions.

57

58

59

60

61

62

63

64

65

66

67

68

69

70

## 71 **Introduction**

72 Ubiquitous use of imaging modalities for evaluation of various medical conditions  
73 leads to the discovery of incidentalomas. Being present in more than half of women  
74 over age 50, thyroid nodules are common incidentalomas. Analysis of Medicare data  
75 (1) showed that thyroid ultrasound as the initial imaging modality in the cohort has risen  
76 by 20.9% year over year.

77 Current classification systems for thyroid nodules are labor intensive and are  
78 subjective (2). The most common systems used to classify thyroid nodules are TIRADS  
79 and American Thyroid Association (ATA) classification system (3, 4). These systems  
80 are fraught with problems. Varying results can be seen when different classification  
81 systems are used to assess the same thyroid nodule. The ability to make a useful  
82 distinction especially by less experienced users is limited by the inherent subjectivity  
83 and the inter and intra reader variability of these visual classification systems. Using the  
84 above systems, Follicular carcinomas, Hurthle cell cancer and follicular variant of  
85 papillary thyroid cancer may end up being classified as benign (5). Not all nodules can  
86 be classified using all available systems. These classification systems also lack  
87 specificity and have low positive predictive value (6). This results in unnecessary  
88 biopsies. Millions of thyroid biopsies are done every year all over the world. It was  
89 estimated that in 2015, more than 600,000 fine needle aspirations (FNA) were done in  
90 the United States alone (7). Evaluation of the increasing number of benign thyroid  
91 incidentalomas is increasing the burden on the healthcare system.

92  
93 Even when FNA of the thyroid is performed it does not always yield a definitive  
94 result. A final diagnosis cannot be made in one out of seven nodules with FNA (8).  
95 Molecular markers were developed to avoid surgery for benign nodules with  
96 indeterminate cytology. The positive predictive value for these molecular tests varies  
97 between 20 to 50% (8, 9). Many times, a repeat biopsy may be required to do molecular  
98 markers. All of this adds to healthcare expense without improving morbidity or mortality.  
99 Another possible outcome of FNA is a non diagnostic cytology report. Management  
100 options for this scenario includes repeat biopsy, surgery or watchful waiting. False  
101 negative results can occur after a biopsy and it is estimated that less than 5% of  
102 cytologically benign nodules are proved to be malignant after surgery. However, only  
103 about 10% of cytologically benign nodules undergo surgery, hence accurate estimates  
104 of false negative results may not be possible (7, 10, 11). Therefore, at present, we do  
105 not have a reliable non subjective method for avoiding invasive procedures for benign  
106 thyroid nodules.

107  
108 Similar problems exist in other medical domains and artificial intelligence (AI)  
109 algorithms have provided solutions. There are FDA cleared AI software to diagnose  
110 diabetic retinopathy, stroke, and breast lesions (12–15). AI algorithms have been used  
111 to classify thyroid nodules objectively (16–19). Given a thyroid ultrasound image, these  
112 algorithms can predict whether a thyroid nodule is benign or malignant. However,  
113 predictions from these algorithms are not generally explainable hence they are called  
114 black box algorithms. In clinical practice, explainable or interpretable deep learning

115 models are needed to gain the trust of physicians (20). Deep learning algorithms may  
116 make predictions based on non-medically relevant information present in the images.  
117 Winkler et al. demonstrated that having gentian violet surgical skin markings in  
118 dermoscopic images increased the chance of melanoma prediction by an algorithm  
119 (21). Images obtained from different imaging machines will have different features. This  
120 can act as a confounding factor while creating deep learning algorithms. For example, if  
121 pneumonia is more common in x-rays obtained in the Emergency Room (ER) when  
122 compared to surgical ward, then there is an increased possibility of false positives on X-  
123 rays performed in ERs. Algorithms created in one institution may underperform when  
124 used in another institution. Algorithms may overlook the actual pathology and instead  
125 may rely on other clinically non relevant features in the image like placement of a  
126 metallic token to mark laterality.(22) Because of these shortcomings we decided to  
127 create an image similarity AI model instead of a classification model. AI image  
128 classification algorithms for thyroid nodules gives a single output, benign or malignant  
129 without any supporting evidence regarding how it reached that conclusion. On the other  
130 hand, an image similarity algorithm will output similar images to the test image with  
131 corresponding diagnosis.

132

133 In this article, we describe the creation of an image similarity deep learning  
134 algorithm for thyroid nodule risk stratification.

## 135 **Materials and methods**

136 The research study was approved by the Mercy Institutional Review Board.

137 Image database

138           Ultrasound images of thyroid nodules from patients who underwent either biopsy  
139 or thyroid surgery from February of 2012 through February of 2017 at Mercy  
140 Endocrinology Clinic or Mercy Hospital in Springfield, Missouri were initially collected for  
141 the study. Cytology and histopathology examinations were done at the Mercy hospital  
142 by one group of pathologists. Cytopathology was reported using the Bethesda system  
143 (23).

144           Nodules were excluded if there was no definitive diagnosis of benignity or  
145 malignancy or if there were no good quality thyroid ultrasound images. 482 nodules met  
146 the inclusion criteria. The area of interest was cropped from the ultrasound images  
147 along with some neighboring tissue. Both sagittal and transverse view images were  
148 used. This image set served as the training database. The testing imaging dataset was  
149 created in a similar fashion from the retrospective collection of ultrasound images of  
150 patients with thyroid nodules who underwent biopsy or surgery between March 2017  
151 through July of 2018. There were 103 thyroid nodules in the testing dataset.

152

153

154 Convolutional Neural Network model

155           A 34 layered Convolutional Neural Network (CNN) - ResNet 34 was trained on  
156 thyroid ultrasound images of 482 thyroid nodules using transfer learning techniques  
157 (24). All images were resized to 224 X 224 pixels before being fed into the CNN. Image  
158 embeddings for these ultrasound images were obtained by taking the output before the  
159 final fully connected layer and stored in a database. Embeddings are N dimensional  
160 vectors representing one unique image. When a query image is received, it is first

161 converted to image embeddings using the CNN. Embeddings from the input image are  
162 used to find embeddings that are similar to it from our training image database using a  
163 nearest neighbor algorithm. Finally, N number of nearest neighbors will be displayed as  
164 the output along with the label of the image. Fig. 1 depicts the schema of image  
165 similarity algorithm.

166

167 Scoring the algorithms.

168 In phase 1 image classification algorithms were used. The ResNet 34 model  
169 trained on 482 nodule images was used to classify test images. Algorithm returned a  
170 prediction for the test image as either benign or malignant. In phase 2, each of the test  
171 images were fed through the image similarity algorithm, AIBx. Image embeddings were  
172 created and the first nearest neighbor / similar image from our training dataset was  
173 identified. If the nearest neighbor for a benign test image is a benign nodule from the  
174 training database, it was considered a true negative. If the algorithm outputs a  
175 malignant nodule as the similar image for a malignant nodule in the test set it was  
176 considered a true positive. Opposite was true for false negative and false positives  
177 respectively.

178

179 Statistical methods

180 A confusion matrix was created from the true positives, true negatives, false  
181 positive and false negatives in both phases. Python programming language was used to  
182 calculate accuracy, sensitivity, specificity, positive predictive value, and negative  
183 predictive value.

184

## 185 **Results**

186           The training dataset consisted of 2025 images from 482 nodules. These included  
187 images with and without square aspect ratio. Testing set had 103 images from 103  
188 nodules. Training and testing set had 66 and 33 malignant nodules respectively. Images  
189 used in the study came from ultrasound machines manufactured by GE, Siemens,  
190 Philips and Sonosite. Of the training set, 6% were subcentimeter nodules. All of the  
191 testing nodules had at least one dimension greater than 1 cm.

192

### 193 Phase 1

194           Image classification using ResNet 34 model resulted in an accuracy of 77.7%.  
195 Sensitivity, specificity, positive predictive value, and negative predictive values were  
196 84.9%, 74.3%, 60.9% and 91.2% respectively. The average time for prediction was 30  
197 milliseconds per image.

198

### 199 Phase 2

200           When the image similarity model was used to classify test Images, accuracy was  
201 81.5%. The sensitivity, specificity, positive predictive value and negative predictive  
202 value of the model was 87.8%, 78.5%, 65.9% and 93.2% respectively. Average time for  
203 prediction was 900 milliseconds.

204

205           In using the image classification algorithm generated in phase 1, 55.3% of the  
206 nodules were determined to be benign. When the image similarity algorithm was used,

207 57.3% of the nodules were determined to be benign. Hence, using image similarity  
208 algorithm will avoid more biopsies. When results of Phase 1 and Phase 2 were  
209 compared (Table 1), image similarity algorithm turned out to be superior.

210

## 211 **Discussion**

212 One of the main challenges in the management of thyroid nodules is risk  
213 stratification. With more than 50% of women over the age of 50 affected by thyroid  
214 nodules and millions of biopsies done every year resulting in less than 5 to 10% cancer  
215 diagnosis, a reliable noninvasive test is needed. Experienced physicians generally  
216 evaluates an ultrasound image and arrives at a decision regarding biopsy based on  
217 their previous experience and heuristics. Most have a mental picture of how a malignant  
218 thyroid nodule should appear. We tried to emulate this by creating an image similarity  
219 CNN model. While the repertoire of representative images stored in a physician's mind  
220 is limited by his experience and memory, AI models can store unlimited images and  
221 query it millions of time.

222

223 Multiple artificial intelligence models have been developed for thyroid nodule  
224 classification, but none of them are widely used (17–19). A recent study by Buda et al  
225 suggested that machine learning algorithms could match the performance of  
226 radiologists in classifying thyroid nodules (16). When tested on 99 nodules, their model  
227 achieved a sensitivity of 87% and specificity was 52%. This was comparable to the  
228 performance of three ACR-TIRADS committee members and nine other radiologists  
229 (16). Most image classification algorithms used in the risk stratification of thyroid

230 nodules are black boxes. Therefore, we cannot readily explain why an algorithm yielded  
231 the wrong classification of a thyroid nodule. Heatmaps have been used to explain  
232 outputs of image classification. But this approach does not help with thyroid nodule  
233 classification. Fig. 2A depicts a benign cystic nodule. Image classification algorithm  
234 correctly classified it as a benign nodule. The corresponding heatmap in Fig. 2A shows  
235 that the algorithm is focusing on the cystic area and the posterior enhancement to arrive  
236 at the diagnosis. Fig. 2B depicts a malignant thyroid nodule. This was classified as a  
237 benign nodule by image classification algorithm. But the heatmaps does not help us to  
238 understand the rationale behind this prediction.

239

240 Explainable AI models will increase the trust of physicians and will foster  
241 adoption of these systems into clinical practice (20). AI healthcare team from Google  
242 created an image similarity model (SMILY) to help diagnose histopathology images  
243 (25). Given a histopathology image, SMILY can output images with similar histological  
244 features. They suggested that this approach could also inform us about the outcomes of  
245 patients with similar pathology.

246

247 To our knowledge, there is no published study on the use of image similarity  
248 models for the classification of thyroid. Unlike other AI algorithms, our image similarity  
249 model, AIBx uses physician in the loop (PIL), Fig. 3. During each stage of AIBx,  
250 physicians have an active role. Operating physicians will select the image to be fed into  
251 AIBx. AIBx will output N number of similar images as requested by the user along with a  
252 classification. Physicians can verify the diagnosis by reviewing the similar images and

253 then accept or reject the classification provided by AIBx. This could also be used to  
254 retrieve all available information for a nodule including diagnosis, molecular markers,  
255 treatment received and recurrence status.

256

#### 257 Advantages

258 If an AI classification algorithm misclassified an image, adding this image back to  
259 the database and retraining the AI model may not result in the correct classification.  
260 Hence, retraining the classification model again with incorrectly classified images may  
261 not always result in better prediction of the misclassified images. However, a properly  
262 trained image similarity model has a high chance of reclassifying the image correctly if  
263 embeddings of the misclassified image are added to the database. Therefore, an image  
264 similarity model does not have to be retrained to increase accuracy. Models trained on  
265 images from one ultrasound machine may not generalize well to images classified from  
266 another ultrasound machine. AIBx used images from ultrasound manufacturers  
267 frequently used in clinical settings, including GE, Philips, Siemens and Sonosite. AIBx  
268 can easily be incorporated into current physician workflow. Any ultrasound machine with  
269 an image output port can wirelessly transmit the image of interest to a nearby mobile  
270 computing device where it can be classified using the algorithm. If this method is used,  
271 data never leaves the healthcare facility. AIBx can also be deployed as a local or  
272 remote website. Furthermore, the system could also be used as a teaching tool for  
273 residents and fellows.

274

#### 275 Potential disadvantages

276 Image similarity algorithms can consume more computing resources and time  
277 when compared to classification algorithms. But the difference was under a second for  
278 AIBx when compared to phase 1, image classification algorithm. The test dataset only  
279 had 103 images. It is possible that a larger test set may yield different results. Images  
280 acquired from an ultrasound machine other than the machines used in the study may  
281 not produce the correct response. However, this could be verified by the physician  
282 comparing similar images generated by AIBx to the test image. During testing, AIBx  
283 retrieved images from ultrasound machines other than the one from the test image,  
284 partly alleviating these concerns. Most thyroid nodules evaluated by FNA or surgery  
285 may have worrisome features and nodules that were not biopsied and/or surgically  
286 removed may have a benign appearance. As such, there could be underlying selection  
287 bias in our database.

288  
289 A study by Grani et al. applied classification systems, ACR-TIRADS, ATA, AACE,  
290 EU-TIRADS and K-TIRADS to 502 nodules and reported that 11 malignant nodules  
291 would have been classified as not requiring biopsy by at least one of these systems (4,  
292 6, 26–29). The ATA system could not classify some of the nodules (6). This study  
293 shows the variability in classification systems even with experienced physicians.  
294 According to this study, ACR-TIRADS performed better and recommended the lowest  
295 number of biopsies. The PPV and NPV for ACR-TIRADS was 12.8% and 97.8%  
296 respectively. In another study by Ahmadi et al. NPV for both ATA and ACR-TIRADS  
297 was 90% (30). All of these classification systems are subjective and will yield different  
298 results when applied by different practitioners. Using AIBx will eliminate subjectivity

299 without significant compromise in negative predictive value. Thyroid cancer has a 98.2%  
300 5-year survival rate and low morbidity (31). Combined with the practice of active  
301 surveillance of thyroid cancer by many centers, a system with greater than 90%  
302 negative predictive value will be helpful in avoiding unnecessary biopsies without  
303 increasing morbidity or mortality (32).

304

305

## 306 **Conclusion**

307 Millions of thyroid biopsies are done every year based on very subjective criteria  
308 to find thyroid cancer in a very small percentage of population with an invasive  
309 technique which may not be diagnostic 1 out of 7 times. Here we described an image  
310 similarity algorithm based on deep learning for thyroid nodule risk stratification. When  
311 compared to published results of ACR TIRADS and ATA classification system, AIbX, the  
312 image similarity model had comparable negative predictive value with better sensitivity  
313 specificity and positive predictive value. By using image similarity AI models, we can  
314 eliminate subjectivity and decrease the number of unnecessary biopsies. This algorithm  
315 may also aid in the management of indeterminate and non-diagnostic thyroid nodules.  
316 Using image similarity AI model, we were able to create an explainable AI model which  
317 encourages physician's confidence in model predictions.

318

## 319 **Acknowledgments**

320 We thank Mercy research for facilitating the study.

321

322 **Author Disclosure Statement**

323 JT - No competing financial interests exist.

324 TH – No competing financial interests exist.

325

326

327 **Corresponding author**

328 Johnson Thomas, MD, FACE

329 3231 S National Ave, Suite 440

330 Springfield, Missouri, USA 65807

331 Email: johnson.thomas@mercy.net

332 Phone: 417 888 5660

333

334 **References**

335 1. Haymart MR, Banerjee M, Reyes-Gastelum D, Caoili E, Norton EC 2018 Thyroid

336 Ultrasound and the Increase in Diagnosis of Low-Risk Thyroid Cancer. J Clin

337 Endocrinol Metab.

338 2. Choi SH, Kim E-K, Kwak JY, Kim MJ, Son EJ 2009 Interobserver and

339 Intraobserver Variations in Ultrasound Assessment of Thyroid Nodules. Thyroid.

340 3. Horvath E, Majlis S, Rossi R, Franco C, Niedmann JP, Castro A, Dominguez M

341 2009 An ultrasonogram reporting system for thyroid nodules stratifying cancer risk

342 for clinical management. J Clin Endocrinol Metab **94**:1748–1751.

343 4. Haugen BR, Alexander EK, Bible KC, Doherty GM, Mandel SJ, Nikiforov YE,

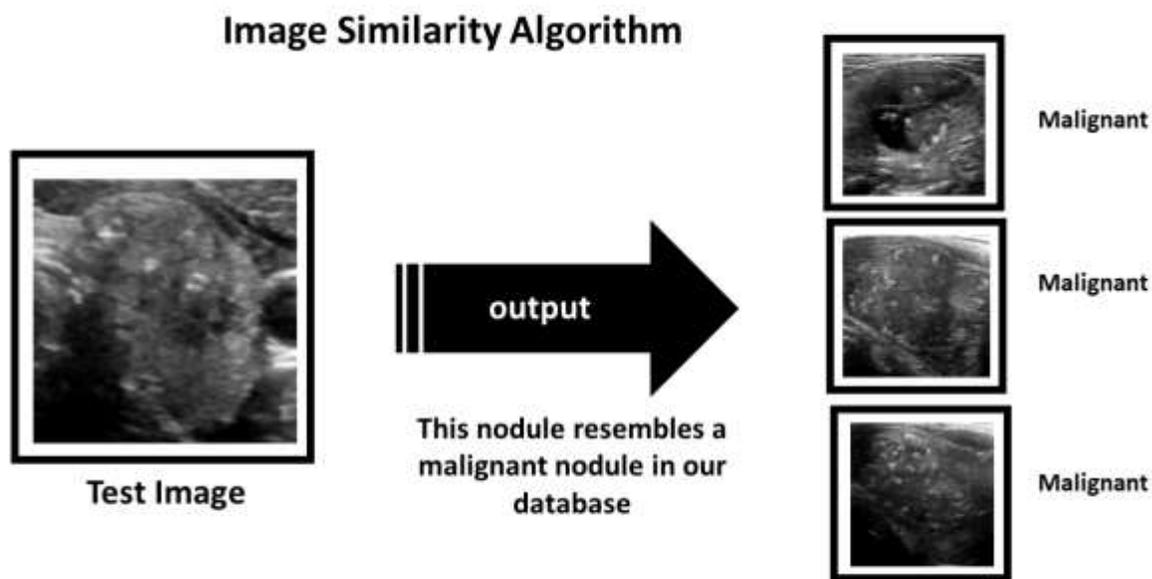
344 Pacini F, Randolph GW, Sawka AM, Schlumberger M, Schuff KG, Sherman SI,

- 345 Sosa JA, Steward DL, Tuttle RM, Wartofsky L 2016 2015 American Thyroid  
346 Association Management Guidelines for Adult Patients with Thyroid Nodules and  
347 Differentiated Thyroid Cancer: The American Thyroid Association Guidelines Task  
348 Force on Thyroid Nodules and Differentiated Thyroid Cancer. *Thyroid* **26**:1–133.
- 349 5. Kim DS, Kim JH, Na DG, Park SH, Kim E, Chang KH, Sohn CH, Choi YH 2009  
350 Sonographic features of follicular variant papillary thyroid carcinomas in  
351 comparison with conventional papillary thyroid carcinomas. *J Ultrasound Med.*
- 352 6. Grani G, Lamartina L, Ascoli V, Bosco D, Biffoni M, Giacomelli L, Maranghi M,  
353 Falcone R, Ramundo V, Cantisani V, Filetti S, Durante C 2019 Reducing the  
354 number of unnecessary thyroid biopsies while improving diagnostic accuracy:  
355 Toward the “Right” TIRADS. *J Clin Endocrinol Metab.*
- 356 7. Dean DS, Gharib H 2015 Fine-Needle Aspiration Biopsy of the Thyroid Gland.
- 357 8. Ali SZ, Siperstein A, Sadow PM, Golding AC, Kennedy GC, Kloos RT, Ladenson  
358 PW 2019 Extending expressed RNA genomics from surgical decision making for  
359 cytologically indeterminate thyroid nodules to targeting therapies for metastatic  
360 thyroid cancer. *Cancer Cytopathol.*
- 361 9. Taye A, Gurciullo D, Miles BA, Gupta A, Owen RP, Inabnet WB, Beyda JN, Marti  
362 JL 2018 Clinical performance of a next-generation sequencing assay (ThyroSeq  
363 v2) in the evaluation of indeterminate thyroid nodules. *Surg (United States).*
- 364 10. Caruso D, Mazzaferri EL 1991 Fine needle aspiration biopsy in the management  
365 of thyroid nodules. *Endocrinologist.*
- 366 11. Gharib H, Goellner JR 1993 Fine-needle aspiration biopsy of the thyroid: An  
367 appraisal. *Ann Intern Med.*

- 368 12. Abramoff MD, Lavin PT, Birch M, Shah N, Folk JC 2018 Pivotal trial of an  
369 autonomous AI-based diagnostic system for detection of diabetic retinopathy in  
370 primary care offices. *npj Digit Med*.
- 371 13. Chilamkurthy S, Ghosh R, Tanamala S, Biviji M, Campeau NG, Venugopal VK,  
372 Mahajan V, Rao P, Warier P 2018 Deep learning algorithms for detection of  
373 critical findings in head CT scans: a retrospective study. *Lancet*.
- 374 14. Pan I, Cadrin-Chênevert A, Cheng PM 2019 Tackling the Radiological Society of  
375 North America Pneumonia Detection Challenge. *Am J Roentgenol*.
- 376 15. Kudo S ei, Mori Y, Misawa M, Takeda K, Kudo T, Itoh H, Oda M, Mori K 2019  
377 Artificial intelligence and colonoscopy: Current status and future perspectives. *Dig*  
378 *Endosc*.
- 379 16. Buda M, Wildman-Tobriner B, Hoang JK, Thayer D, Tessler FN, Middleton WD,  
380 Mazurowski MA 2019 Management of Thyroid Nodules Seen on US Images:  
381 Deep Learning May Match Performance of Radiologists. *Radiology* 181343.
- 382 17. Guan Q, Wang Y, Du J, Qin Y, Lu H, Xiang J, Wang F 2019 Deep learning based  
383 classification of ultrasound images for thyroid nodules: a large scale of pilot study.  
384 *Ann Transl Med*.
- 385 18. Yu B, Wang Z, Zhu R, Feng X, Qi M, Li J, Zhao R, Huang L, Xin R, Li F, Zhou F  
386 2019 The Transverse Ultrasonogram of Thyroid Papillary Carcinoma Has a Better  
387 Prediction Accuracy Than the Longitudinal One. *IEEE Access*.
- 388 19. Verburg F, Reiners C 2019 Sonographic diagnosis of thyroid cancer with  
389 support of AI. *Nat Rev Endocrinol*.
- 390 20. E. Khan C 2019 *Radiology: Artificial Intelligence. Explain AI*.

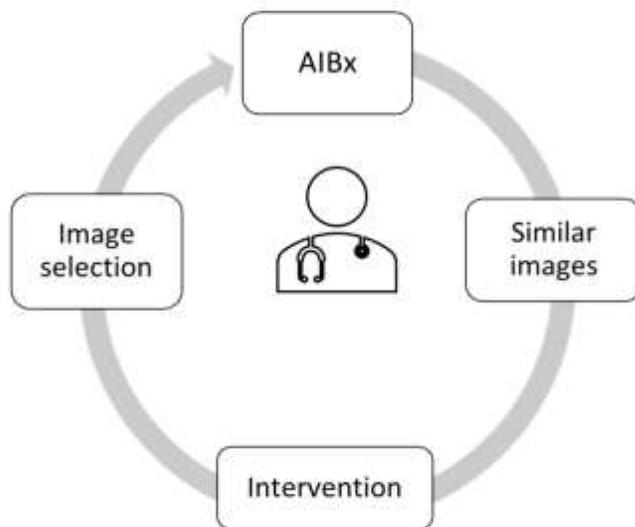
- 391 21. Winkler JK, Fink C, Toberer F, Enk A, Deinlein T, Hofmann-Wellenhof R, Thomas  
392 L, Lallas A, Blum A, Stolz W, Haenssle HA 2019 Association Between Surgical  
393 Skin Markings in Dermoscopic Images and Diagnostic Performance of a Deep  
394 Learning Convolutional Neural Network for Melanoma Recognition. *JAMA*  
395 *Dermatology*.
- 396 22. Zech JR, Badgeley MA, Liu M, Costa AB, Titano JJ, Oermann EK 2018 Variable  
397 generalization performance of a deep learning model to detect pneumonia in  
398 chest radiographs: A cross-sectional study. *PLoS Med*.
- 399 23. Cibas ES, Ali SZ, of the Science Conference NCITFNAS 2009 The Bethesda  
400 system for reporting thyroid cytopathology. *Am J Clin Pathol* **132**:658–665.
- 401 24. He K, Zhang X, Ren S, Sun J 2016 Deep residual learning for image  
402 recognition Proceedings of the IEEE Computer Society Conference on Computer  
403 Vision and Pattern Recognition.
- 404 25. Hegde N, Hipp JD, Liu Y, Emmert-Buck M, Reif E, Smilkov D, Terry M, Cai CJ,  
405 Amin MB, Mermel CH, Nelson PQ, Peng LH, Corrado GS, Stumpe MC 2019  
406 Similar image search for histopathology: SMILY. *npj Digit Med* **2**:56.
- 407 26. Tessler FN, Middleton WD, Grant EG, Hoang JK, Berland LL, Teefey SA, Cronan  
408 JJ, Beland MD, Desser TS, Frates MC, Hammers LW, Hamper UM, Langer JE,  
409 Reading CC, Scoutt LM, Stavros AT 2017 ACR Thyroid Imaging, Reporting and  
410 Data System (TI-RADS): White Paper of the ACR TI-RADS Committee. *J Am Coll*  
411 *Radiol*.
- 412 27. Gharib H, Papini E, Garber JR, Duick DS, Harrell RM, Hegedüs L, Paschke R,  
413 Valcavi R, Vitti P, Balafouta ST, Baloch Z, Crescenzi A, Dralle H, Frasoldati A,

- 414 Gärtner R, Guglielmi R, Mechanick JI, Reiners C, Szabolcs I, Zeiger MA, Zini M  
415 2016 American association of Clinical Endocrinologists, American college of  
416 endocrinology, and Associazione Medici Endocrinologi medical guidelines for  
417 clinical practice for the diagnosis and management of thyroid nodules - 2016  
418 update. *Endocr Pract*.
- 419 28. Russ G, Bonnema SJ, Erdogan MF, Durante C, Ngu R, Leenhardt L 2017  
420 European Thyroid Association Guidelines for Ultrasound Malignancy Risk  
421 Stratification of Thyroid Nodules in Adults: The EU-TIRADS. *Eur Thyroid J*.
- 422 29. Shin JH, Baek JH, Chung J, Ha EJ, Kim JH, Lee YH, Lim HK, Moon WJ, Na DG,  
423 Park JS, Choi YJ, Hahn SY, Jeon SJ, Jung SL, Kim DW, Kim EK, Kwak JY, Lee  
424 CY, Lee HJ, Lee JH, Lee JH, Lee KH, Park SW, Sung JY 2016 Ultrasonography  
425 diagnosis and imaging-based management of thyroid nodules: Revised Korean  
426 society of thyroid radiology consensus statement and recommendations. *Korean J*  
427 *Radiol*.
- 428 30. Ahmadi S, Oyekunle T, Sara X, Scheri R, Perkins J, Stang M, Roman S, Sosa JA  
429 2019 A direct comparison of the ATA and TI-RADS ultrasound scoring systems.  
430 *Endocr Pract*.
- 431 31. 2019 Cancer of the Thyroid - Cancer Stat Facts. SEER.
- 432 32. Tuttle RM, Fagin JA, Minkowitz G, Wong RJ, Roman B, Patel S, Untch B, Ganly I,  
433 Shaha AR, Shah JP, Pace M, Li D, Bach A, Lin O, Whiting A, Ghossein R, Landa  
434 I, Sabra M, Boucai L, Fish S, Morris LGT 2017 Natural history and tumor volume  
435 kinetics of papillary thyroid cancers during active surveillance. *JAMA Otolaryngol -*  
436 *Head Neck Surg*.

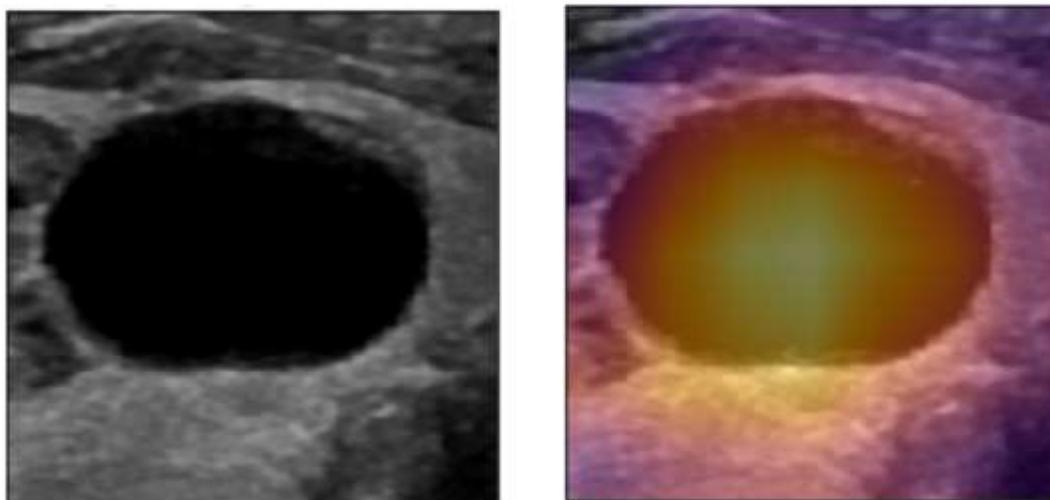


437

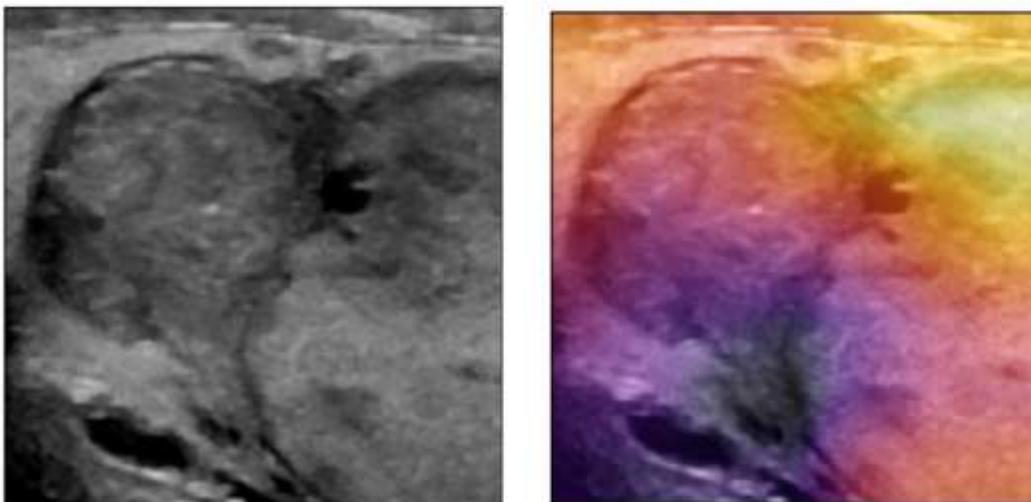
Preprint V2



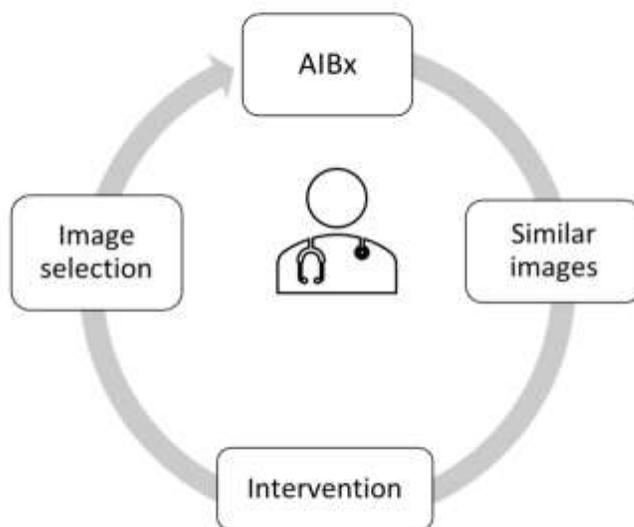
438



439



440



441