1   **Abstract**

2   Objectives

3   To evaluate the efficacy of combining predictive artificial intelligence and image similarity model to risk

4   stratify thyroid nodules, using a retrospective study.

5   Methods

6   Two datasets were used to determine the efficacy of the algorithm. One was the publicly available

7   Stanford dataset consisting of ultrasound images of 192 nodules between April 2017 to May 2018 and

8   the second one was from a private practice setting consisting of 118 thyroid nodule images from 2018-

9   2023. All the nodules had definitive diagnosis either by biopsy or by surgery.  The software was used to

10   predict the diagnosis and TI-RADS score.

11   Results

12   In the Stanford dataset, the AI algorithm predicted malignancies with a sensitivity of 1.0 and a specificity

13   of 0.55. The PPV was 0.18 and the NPV was 1.0. The AUCROC was 0.78. The AI algorithm did not miss

14   any cases of cancer. TI-RADS based clinical recommendation had a polychoric correlation of 0.67. In the

15   private dataset, the AI algorithm predicted malignancies with a sensitivity of 0.91 and a specificity of

16   0.95. The PPV was 0.8 and NPV was 0.98. AUCROC was 0.93 and accuracy was 0.94. TI-RADS based

17   clinical recommendation had a polychoric correlation of 0.94 for this dataset.

18   Conclusion

19   The AI model demonstrated high negative predictive value with a potential for 60% reduction in the

20   need for biopsy. This could reduce the burden on patients and healthcare costs.

21   Introduction

22

23   Thyroid nodules are commonplace findings in clinical settings, with an estimated prevalence in the

24   general population ranging from 4% to 6.5%.[1] Though the vast majority of these nodules are benign,

25   roughly 10%-15% of them are malignant.[2]  Currently, the best method to evaluate thyroid nodules

26   involve ultrasound-guided fine-needle aspiration biopsy, which is invasive and can be emotionally

27   distressing for patients.[3] Moreover, up to 30% of biopsies lead to indeterminate results, requiring a

28   repeat biopsy or surgery.[4] To effectively distinguish whether a thyroid nodule is benign or malignant is

29   crucial in determining accurate clinical management and reducing the number of unnecessary biopsies.

30

31   AI has been increasingly utilized in various fields of medicine, including radiology and pathology,

32   demonstrating its potential to augment the accuracy of diagnosis.[5] In thyroid nodule evaluation

33   particularly, AI-driven predictive models offer non-invasive strategies to detect malignancies.[6]

34

35   Additionally, image similarity assessment, which involves the comparison of visual characteristics of

36   images, can also be used in medical diagnostics. It offers an efficient analysis of medical images that may

37   exceed the capabilities of the human eye.[7] The potential of combining AI-driven predictive models with

38   image similarity assessment in thyroid nodule evaluation has not been explored for diagnosis and ACR

39   TI-RADS assessment.[8] Therefore, in this paper we are evaluating the benefits of combining these

40   methods. We elucidate the methods of our software, evaluate its performance, and discuss the

41    potential implications of combining image similarity and AI to provide better screening for thyroid

42    nodules.

43    Materials and methods

44    In this study, we used software that integrates AI-driven predictive models with image similarity

45    assessment for thyroid nodule evaluation.  Version 2 of this software also predicts ACR TI-RADS. PEARL

46    IRB determined the study to be exempt. Two diverse datasets were used to evaluate the AI model. The

47    first dataset is an open-source dataset from Stanford University from 2021, which consists of 192 images

48    of thyroid nodules .[9] These images were collected between April 2017 and May 2018. The second data

49    set is from a private practice setting consisting of 118 thyroid nodule images from 2018 to 2023. This

50    data set consists of images from an in-house thyroid ultrasound machine as well as an external radiology

51    ultrasound machine. Both data sets had confirmed cytopathology and a TIRADS score. For the second

52    data set with in-house images, the TIRADS score was assigned by the performing endocrinologist (RV)

53    which was then reviewed and confirmed by a second endocrinologist (RP). Any discrepancies were

54    resolved by a third endocrinologist(JC).

55

56    The inclusion criteria for the study were males and females, aged 18 years with thyroid surgery or biopsy

57    at participating sites with a definitive diagnosis by cytology or pathology. Indeterminate nodules

58    (Bethesda III, IV, and V) upon initial evaluation should have undergone surgery with a definitive

59    diagnosis to be included in the study. Thyroid nodules measuring between 5 mm and 40 mm (4.0cm) in

60    the maximum dimension by ultrasound imaging in transverse dimension. The longest diameter of the

61    thyroid nodule should be less than the length of the ultrasound transducer.

62

63    Exclusion criteria for the study were patients below the age of 18 years; indeterminate thyroid nodules

64    without a definitive diagnosis; ultrasound images of thyroid nodules containing annotations, markings,

65    writings, or crosshair within the nodule and whole thyroid nodule not visible in the ultrasound section;

66    metastasis to the thyroid from other malignancies as well as lymphoma of the thyroid were also

67    excluded; multinodular goiters without a clearly separable nodule on ultrasound images and nodules

68    that underwent radioactive iodine treatment, ethanol ablation, radiofrequency ablation or laser

69    ablation.

70

71    The software uses static images in the AP dimension. It automatically identifies regions of interest. By

72    comparing these regions to images in the training dataset, the software predicts whether the nodule is

73    benign or malignant and also provides an ACR Thyroid Imaging Reporting and Data System (TI-RADS)

74    score (Figure 1).

75    We used Python language with Sckit-learn library to do the statistical analysis.[10]

76

77    Results

78

79    In the Stanford public dataset, there were 17 malignant nodules and 175 benign nodules. The

80    prevalence of malignancy in this dataset was 8 percent. Compared to ground truth, the AI algorithm

81    predicted malignancies with a sensitivity of 1.0 and a specificity of 0.55. The positive predictive value

82    (PPV) was 0.18 and the negative predictive value (NPV) was 1.0. The AUCROC was 0.78. The AI algorithm

83    did not miss any cases of cancer. ACR TI-RADS based clinical recommendation had a polychoric

84    correlation of 0.67.

85   In the private dataset there were 96 benign nodules and 22 malignant nodules. The prevalence of

86   malignancy was 18.85%. In this dataset, the AI algorithm predicted malignancies with a sensitivity of

87   0.91 and a specificity of 0.95. The PPV was 0.8 and NPV was 0.98. AUCROC was 0.93 and accuracy was

88   0.94. In this dataset, 5 out of the 99 benign nodules were read as malignant by the AIBx algorithm.

89   These nodules had high risk features with TI-RADS scores 4-5 for 4 out of the 5 nodules and TI-RADS

90   score of 3 for 1 out of the 5 nodules. The AIBx algorithm also predicted 2 out of the 22 malignant

91   nodules as benign. These nodules had a TI-RADS score of 3. TI-RADS-based clinical recommendation had

92   a polychoric correlation of 0.94 for this dataset. Table 1, shows comparison of AI predictions on both

93   datasets.

94   The Pearson correlation coefficient between ground truth cytopathology diagnosis and AI diagnosis was

95   0.824 with a p-value of $2.29 \times 10^{-31}$, indicating a strong positive correlation that is statistically

96   significant. The AI program and ground truth diagnoses exhibit high agreeability with a concordance rate

97   of 94.26 percent and an F1 score of 85.21 percent.

98   Regarding the TI-RADS score by a physician vs that was predicted by AI algorithm, the Pearson

99   Correlation Coefficient was 0.877 with p< 0.001 indicating a strong linear relationship between the two

100  readings. Cohen's Kappa for physician readings vs AI reading was 0.753. This indicates substantial

101  agreement between the physician and the AI system.

102  Discussion

103  In recent years, artificial intelligence tools have become increasingly prevalent across multiple

104  disciplines.

105  AI can be particularly useful in evaluating thyroid nodules, typically for risk stratification.[6,11]

106    Recent studies suggested that the performance of artificial intelligence models was better or at par

107    with radiologists.[12,13] These studies postulated that artificial intelligence software can be especially

108    beneficial for physicians with less experience. Currently, the United States Food and Drug Administration

109    has approved four AI platforms for thyroid disease. Despite the reported efficacy of artificial intelligence,

110    common concerns exist with its usability, such as the proper integration of AI and radiologist

111    interpretations and assessment of productivity. Furthermore, the authors concluded that the successful

112    adoption of AI platforms requires that the software be incorporated into the physician's workflow

113    seamlessly and should have external validation studies.[6]  Our software addresses some of these

114    concerns. By generating human-understandable descriptors and explanations for its decisions, our

115    software's interpretations can be verified by physicians. Having a high negative predictive value and

116    decreasing biopsy need by 60%, this software demonstrated its ability to reduce healthcare spending.

117    Coupled with its easy-to-use nature, this software ensures practicality, workflow efficiency, and

118    demonstrable performance, all of which are critical for acceptance in clinical settings.

119

120    Need for explainability in medical AI models.

121    Explainable AI or interpretable AI, is a set of tools and methods that help people understand and

122    interpret predictions made by their machine learning algorithms.[14] This consists of an explainable model

123    and an explanation interface so human users can understand what caused the model to make a certain

124    conclusion or prediction, which helps characterize model accuracy, fairness, transparency, and

125    outcomes in decision-making powered by AI.[15] However, there is a reluctance to use medical AI due to a

126    combination  of lack of focus on the end-user by developers of the AI leading to a subjective difficulty of

127    understanding the algorithm and more comfort with human decision making. [16,17] Therefore, focusing on

128    the end user by developers of medical AI as well as interventions to increase the understanding of a

129     medical algorithmic decision process would be important to increase utilization. This is especially crucial

130     in medicine because medical professionals need to understand the basis for an algorithm's diagnosis. A

131     false negative could mean that a patient doesn't receive life-saving treatment, and a false positive could

132     result in a patient receiving expensive and invasive treatments when it isn't necessary to do so.[18] A level

133     of explainability is essential for medical professionals to have comfort in integrating medical AI into

134     practice. Our AI algorithm took these factors into consideration with its easy to use interface and

135     transparency in decision making that makes it very user-friendly and easy to integrate into daily practice

136     with confidence.

137

138     Due to a lack of validation, many AI technologies are not applied in clinical decision-making 15. External

139     validation is used to evaluate predictive capabilities for target clinical implementations in different

140     populations and settings.[19] Predictive models often perform well under training datasets. However,

141     there is a discrepancy between training and validation performance. This discrepancy even appears

142     when training and validation datasets are from the same populations and settings. Poorly developed

143     models lead to exacerbated disparities in healthcare provisions and outcomes. Thus, external validation

144     is necessary to avoid the consequences of a model with low adaptability. External validation is critical to

145     understanding the clinical utility of prediction models.[20]   Hence we undertook external validation on

146     two widely different datasets and demonstrated good performance.

147

148     One of the unique aspects of our research is its integration of image similarity assessment and TI-RADS

149     scoring to produce diagnostic outcomes, a combination that has not been explored before in thyroid

150     nodules. Image similarity assessment uses visual pattern recognition to compare and contrast features

151    of a nodule against a repository of images already classified as malignant or benign. This results in a

152    more accurate evaluation while simultaneously allowing medical professionals to verify the algorithm's

153    conclusions. A TI-RADS score aids in this endeavor by providing human-understandable descriptors to fill

154    the gap between the novelty of AI algorithms and the traditional use of clinical assessment. Our

155    software identifies similar images from its database when compared to the test image. The diagnosis of

156    the most similar image is displayed as the output of the AIBx algorithm.  A TI-RADS score description and

157    recommendation is then produced by the model to enable verification by medical professionals.

158

159    Some limitations of our study were the small sample size, use of static images, and the low number of

160    malignant cases. These could have contributed to the low positive predictive value. In the future, we

161    could test it on databases with a higher prevalence of malignancy. But the average prevalence of

162    malignancy in the combined dataset was very similar to the general population. Furthermore, this

163    software was not prospectively evaluated in a clinical setting.

164

165    The results from the study showed a high negative predictive value, meaning if our algorithm predicted

166    that a nodule is benign, it had a very low probability of being malignant. This would translate into

167    observation as opposed to undergoing a biopsy. The AI algorithm missed only 2 malignant nodules. Both

168    of these nodules were follicular carcinomas of the thyroid and had benign characteristics isoechoic, clear

169    borders, and small central cystic spaces. However, feedback to the AI with these types of nodules as

170    malignant could lead to better predictions in the future. Our AI model performed well with ultrasound

171    images across multiple institutions using different ultrasound machines and showed no bias across

172    nodules of various types and sizes and age groups.

173 Conclusion

174 The combined image similarity and AI model demonstrated high negative predictive value with a

175 potential for a 60% reduction in the need for biopsy. This holds significant clinical implications, as the

176 integration of image similarity and AI-driven predictive models could revolutionize the approach to

177 thyroid nodule evaluation. Not only does this pave the way for non-invasive screening, but it also has

178 the potential to greatly reduce the burden on patients and healthcare costs alike.

179

180 References

181 1.  Popoveniuc G, Jonklaas J. Thyroid Nodules. *Med Clin North Am*. 2012;96(2):329-349.
182     doi:10.1016/j.mcna.2012.02.002

183 2.  Kamran SC, Marqusee E, Kim MI, et al. Thyroid nodule size and prediction of cancer. *J Clin Endocrinol*
184     *Metab*. 2013;98(2):564-570. doi:10.1210/jc.2012-2968
185 3.  Jasim S, Dean DS, Gharib H. Fine-Needle Aspiration of the Thyroid Gland. In: Feingold KR, Anawalt B,
186     Blackman MR, et al., eds. *Endotext*. MDText.com, Inc.; 2000. Accessed February 19, 2024.
187     http://www.ncbi.nlm.nih.gov/books/NBK285544/
188 4.  Yip L, Farris C, Kabaker AS, et al. Cost Impact of Molecular Testing for Indeterminate Thyroid Nodule
189     Fine-Needle Aspiration Biopsies. *J Clin Endocrinol Metab*. 2012;97(6):1905-1912.
190     doi:10.1210/jc.2011-3048
191 5.  Topol EJ. High-performance medicine: the convergence of human and artificial intelligence. *Nat*
192     *Med*. 2019;25(1):44-56. doi:10.1038/s41591-018-0300-7
193 6.  Tessler FN, Thomas J. Artificial Intelligence for Evaluation of Thyroid Nodules: A Primer. *Thyroid Off J*
194     *Am Thyroid Assoc*. 2023;33(2):150-158. doi:10.1089/thy.2022.0560
195 7.  Krupinski EA. Current perspectives in medical image perception. *Atten Percept Psychophys*.
196     2010;72(5):10.3758/APP.72.5.1205. doi:10.3758/APP.72.5.1205
197 8.  Tessler FN, Middleton WD, Grant EG, et al. ACR Thyroid Imaging, Reporting and Data System (TI-
198     RADS): White Paper of the ACR TI-RADS Committee. *J Am Coll Radiol JACR*. 2017;14(5):587-595.
199     doi:10.1016/j.jacr.2017.01.046
200 9.  Yamashita R, Kapoor T, Alam MN, et al. Toward Reduction in False-Positive Thyroid Nodule Biopsies
201     with a Deep Learning-based Risk Stratification System Using US Cine-Clip Images. *Radiol Artif Intell*.
202     2022;4(3):e210174. doi:10.1148/ryai.210174
203 10. Pedregosa F, Varoquaux G, Gramfort A, et al. Scikit-learn: Machine Learning in Python. *J Mach Learn*
204     *Res*. 2011;12(85):2825-2830.
205 11. Wildman-Tobriner B, Taghi-Zadeh E, Mazurowski MA. Artificial Intelligence (AI) Tools for Thyroid
206     Nodules on Ultrasound, From the AJR Special Series on AI Applications. *AJR Am J Roentgenol*.
207     2022;219(4):1-8. doi:10.2214/AJR.22.27430

208    12. Park VY, Han K, Seong YK, et al. Diagnosis of Thyroid Nodules: Performance of a Deep Learning
209          Convolutional Neural Network Model vs. Radiologists. *Sci Rep*. 2019;9(1):17843.
210          doi:10.1038/s41598-019-54434-1

211    13. He LT, Chen FJ, Zhou DZ, et al. A Comparison of the Performances of Artificial Intelligence System
212          and Radiologists in the Ultrasound Diagnosis of Thyroid Nodules. *Curr Med Imaging*.
213          2022;18(13):1369-1377. doi:10.2174/1573405618666220422132251

214    14. Barredo Arrieta A, Díaz-Rodríguez N, Del Ser J, et al. Explainable Artificial Intelligence (XAI):
215          Concepts, taxonomies, opportunities and challenges toward responsible AI. *Inf Fusion*. 2020;58:82-
216          115. doi:10.1016/j.inffus.2019.12.012

217    15. What is Explainable AI? - Unite.AI. Accessed February 19, 2024. https://www.unite.ai/what-is-
218          explainable-ai/

219    16. Cadario R, Longoni C, Morewedge CK. Understanding, explaining, and utilizing medical artificial
220          intelligence. *Nat Hum Behav*. 2021;5(12):1636-1642. doi:10.1038/s41562-021-01146-0

221    17. Chen H, Gomez C, Huang CM, Unberath M. Explainable medical imaging AI needs human-centered
222          design: guidelines and evidence from a systematic review. *Npj Digit Med*. 2022;5(1):1-15.
223          doi:10.1038/s41746-022-00699-2

224    18. McNamara M. Explainable AI: What is it? How does it work? And what role does data play?
225          Published February 22, 2022. Accessed February 19, 2024.
226          https://www.netapp.com/blog/explainable-ai/

227    19. Tsopra R, Fernandez X, Luchinat C, et al. A framework for validating AI in precision medicine:
228          considerations from the European ITFoC consortium. *BMC Med Inform Decis Mak*. 2021;21(1):274.
229          doi:10.1186/s12911-021-01634-3

230    20. Riley RD, Archer L, Snell KIE, et al. Evaluation of clinical prediction models (part 2): how to undertake
231          an external validation study. *BMJ*. 2024;384:e074820. doi:10.1136/bmj-2023-074820

232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252

253    *Table 1: Comparison of AI predictions on both datasets*

|  | Stanford Data | Private Data |
|---|---|---|
| Sensitivity | 1 | 0.91 |
| Specificity | 0.55 | 0.95 |
| PPV | 0.18 | 0.8 |
| NPV | 1 | 0.98 |
| AUC ROC | 0.78 | 0.93 |

254

255

256

257

258

259

260

261

262

263

264

265

266    Figure 1: AI software result interface.



267

268