

---

**Abstract**

This study evaluated three timing strategies for delivering AI assistance in pathological slide diagnosis—pre-diagnosis (triage), during diagnosis (concurrent), and post-diagnosis (secondary)—and assessed users’ perceptions of AI assistance. All three AI modes improved diagnostic performance and reduced workload versus no AI assistance. Concurrent mode was preferred for its balance between efficiency and reader control; secondary mode was appreciated for minimizing bias and aiding training. Triage mode yielded lower workload and higher performance but raised concerns about trust and transparency. AI was regarded as a valuable tool for initial slide review, but not as a replacement for expert readers. Participants generally trusted the AI for highlighting suspicious areas, not making final decisions. After use, willingness to rely on AI for final diagnosis declined, though trust and usability remained moderate to high. To increase adoption, designers should manage AI information presentation to avoid bias, balance sensitivity and specificity based on user feedback, and improve explainability to enhance reader confidence and trust.

*Keywords:*

Human-AI collaboration, Usability/acceptance measurement and research, AI-assisted diagnosis, Trust in AI

---

## 1. Introduction

Recent advancements in artificial intelligence (AI), particularly in deep learning, have created great opportunities to revolutionize the field of pathological slide diagnosis. AI tools enable readers to achieve greater diagnostic precision, improve the likelihood of early detection of critical conditions, and contribute to predicting disease progression, thereby facilitating the development of more accurate and timely treatment plans (Abraham & Joseph, 2024; Ahmadi, 2024; Peng & Ren, 2024). Additionally, AI assistance enhances efficiency in various ways, such as replacing human proofreaders and automating complex image interpretation tasks, allowing physicians to focus more on higher-level decision-making (Abraham & Joseph, 2024; Peng & Ren, 2024; Avanzo et al., 2024). While the potential applications of AI-assisted pathological slide diagnosis are promising, fostering reader acceptance and effective use of AI tools is essential, which can be facilitated by enhancing the usability of AI systems through understanding readers' trust in AI and the ideal timing of introducing AI assistance during diagnosis.

For AI-assisted slide diagnoses, readers' trust in AI assistance depends on the AI's ability to consistently provide accurate and clinically relevant insights (Hatherley, 2020; Hallowell et al., 2022). However, a challenge associated with AI-assisted diagnosing tools is their "black box" nature, which prevents pathologists from understanding how AI decisions are made. This lack of transparency often leads to skepticism and reluctance to rely on AI tools (Poon & Sung, 2021; Nasarian et al., 2024). Pathologists are more likely to trust AI when clear explanations of its decisions, along with the evidence supporting these decisions are provided (Nasarian et al., 2024; Yang

26 et al., 2023). Trust in AI can be further enhanced when the system as-  
27 sumes the role of offering recommendations rather than acting as the pri-  
28 mary decision-maker. This approach enables readers to confirm or reject the  
29 AI’s suggestions at their discretion, fostering a collaborative relationship be-  
30 tween human expertise and AI assistance (Yang et al., 2023; Sivaraman et al.,  
31 2023; Burgess et al., 2023). Ideally, AI systems must also balance specificity  
32 and sensitivity during classification tasks and demonstrate the capacity to  
33 accurately identify outliers (Kan se et al., 2022; McCague et al., 2023).

34 Although AI assistance has generally proven beneficial, the timing of pre-  
35 senting AI information can further influence system performance and user  
36 experience—an aspect that has yet to be systematically evaluated in most  
37 studies. Depending on when AI input is delivered, there are three distinct  
38 collaboration modes between AI and readers (Chen et al., 2024). In the  
39 **triage assistance mode**, AI acts as a pre-screener by filtering out slides  
40 it identifies as negative, leaving only those diagnosed as positive for human  
41 review (e.g., (Yala et al., 2019; Raya-Povedano et al., 2021; Lauritzen et al.,  
42 2022; Leibig et al., 2022; Dvijotham et al., 2023)). This approach positions  
43 AI at the forefront of the diagnostic process, handling the bulk of the task  
44 while humans focus on reviewing slides with potentially critical issues. Such  
45 collaboration can significantly improve efficiency, allowing users to dedicate  
46 more time and attention to priorities (Farber, 2025; Xu et al., 2023; Baruwal  
47 Chhetri et al., 2024; Lai & Rau, 2024). However, this may also reduce hu-  
48 man control and adaptability in responding to AI errors or oversights, which  
49 can be particularly consequential in high-risk activities like medical diagnosis  
50 (Hauptman et al., 2024; Hassan & El-Ashry, 2024). Additionally, while AI

51 automation provides greater consistency and accuracy (Smith et al., 2025),  
52 especially in repetitive, data-driven tasks such as pathological slide review,  
53 human creativity and intuition remain essential for handling complex or novel  
54 cases (Koivisto & Grassini, 2023; Moura, 2023). AI systems may also face  
55 challenges in adapting to situations that require knowledge beyond its train-  
56 ing data.

57 In the **concurrent assistance mode**, readers review all slides while AI  
58 acts as an accompanying reader that offers real-time diagnostic suggestions  
59 (e.g., (Sung et al., 2021; Duron et al., 2021; Guerhazi et al., 2022; Hendrix  
60 et al., 2023; Ajmera et al., 2023)). This approach grants pathologists greater  
61 control over the diagnostic process. However, as readers are exposed to the  
62 AI’s identified suspicious regions and diagnostic suggestions, their attention  
63 and decisions may become biased toward the AI’s input (Bansal et al., 2021).  
64 Over time, this could lead to reduced critical thinking and over-reliance on  
65 AI, potentially hindering readers’ skill development (Zhai et al., 2024; Mac-  
66 namara et al., 2024), especially if the AI’s performance falls short of that of  
67 a well-trained human expert (He et al., 2023).

68 Finally, in the **secondary assistance mode**, AI acts as a proofreader  
69 by reviewing the readers’ initial diagnosis and providing feedback for recon-  
70 sideration (e.g., (Marinovich et al., 2023; Frazer et al., 2023; McKinney et  
71 al., 2020)). This approach grants readers the most control, as they indepen-  
72 dently review all slides before receiving AI input. However, it may increase  
73 readers’ time and task load compared to working without AI assistance, and  
74 it limits their ability to enhance efficiency with AI features. This mode also  
75 undermines two key advantages of pre-trained AI tools: (1) their ability to be

76 rapidly deployed to address labor shortages and (2) their capacity to support  
77 workers with limited specialized skills (Szajna & Kostrzewski, 2022; Rasheed  
78 et al., 2020).

79 In this paper, we present an experiment in which pathologists reviewed  
80 slides both without AI assistance and with AI performing the three afore-  
81 mentioned roles to address two research questions:

- 82 1. How does the timing of AI information presentation impact perfor-  
83 mance, workload, and user preference in human-AI collaborative patho-  
84 logical slide diagnosis?
- 85 2. What are end-users' general perceptions of AI assistance in terms of  
86 trust and perceived usability?

87 The insights from this study aim to inform the future design of AI-based  
88 pathological slide diagnosis tools and strategies for their effective adoption.

## 89 **2. Methods**

### 90 *2.1. Participants*

91 A total of 108 cytopathologists of varying experience levels (87 females  
92 and 21 males;  $M_{\text{age}} = 37.9$ ,  $SD = 7.5$ ) were recruited from multiple hospitals  
93 in China. Participants had an average of 7.8 years of work experience in  
94 routine cytopathological reading ( $SD = 5.7$  years). Prior to the study, all  
95 participants provided informed consent in accordance with the procedures  
96 approved by the Institutional Review Board of Chinese Academy of Medi-  
97 cal Sciences and Peking Union Medical College (Approval No. CAMS and  
98 PUMC-IEC-2022-022).

99 *2.2. Experimental Design*

100 The experimental procedure is shown in Figure 1. All participants com-  
101 pleted a pre-study survey collecting basic demographic information and years  
102 of experience. The pre-study survey included seven questions assessing trust  
103 and perceived usability of AI tools. These questions were also part of the  
104 post-study survey to see whether there were changes in their opinion after  
105 interacting with AI as part of this study. After completing the survey, each  
106 participant would perform four blocks of cervical cytology slides reading. A  
107 total of 1,620 slides were selected from real patient data based on several  
108 criteria, that included slide quality, patient age, treatment history, and cer-  
109 vical condition, to balance the diagnostic effort required was consistent across  
110 slides. All included slides had definitive diagnostic outcomes and were con-  
111 firmed by two senior cytopathologists. In each block, participants reviewed  
112 15 randomly assigned slides, with no more than three negative cases, to  
113 approximate real-world conditions where negative cases are more prevalent  
114 than positive ones. Each participant was assigned a different set of slides in  
115 each block to avoid repetition. Each block featured a different AI assistance  
116 mode and block order was randomized. The four AI assistance modes were  
117 as follows:

- 118 1. **Unassisted.** Participants would review all slides without any assis-  
119 tance or feedback from AI.
- 120 2. **Triage AI Assistance.** AI would review all slides first and filter out  
121 the negative slides. Only the positive slides identified by AI would be  
122 presented to the participant for independent review.
- 123 3. **Concurrent AI Assistance.** Participants would work side by side

124 with the AI, which means that they could review all slides and assess  
125 AI suggestions for any slide at any time.

126 4. **Secondary AI Assistance.** Participants first reviewed and made di-  
127 agnosis on all slides without AI assistance. After, they were presented  
128 with the AI’s diagnoses and at which point they could make adjust-  
129 ments to their diagnoses.

130 There was no preset number of negative or positive slides for each block. After  
131 every block, participants were asked to complete a NASA-TLX questionnaire  
132 to assess their perceived workload.

133 After completing all four blocks of slide review, participants were asked  
134 to fill out a post-study survey consisting of three parts:

- 135 1. A 13-question survey evaluating trust in the AI assistance tool, devel-  
136 oped based on the survey proposed by Hoffman et al. (2018). Responses  
137 were measured on a Likert scale ranging from ”strongly disagree” to  
138 ”strongly agree.”
- 139 2. A ten-question System Usability Scale (SUS) survey (Brooke, 1996),  
140 also using a Likert scale.
- 141 3. Four additional questions involved: ranking the four AI assistance  
142 modes, comparing experiences with and without AI assistance, and  
143 providing open-ended feedback about the AI assistance tool.

### 144 2.3. Task

145 In this study, the main task for the participants was to examine digital  
146 cervical cytology slides using an updated version of the diagnostic system  
147 proposed by Xue et al. (2023). For each trial, the participant needed to

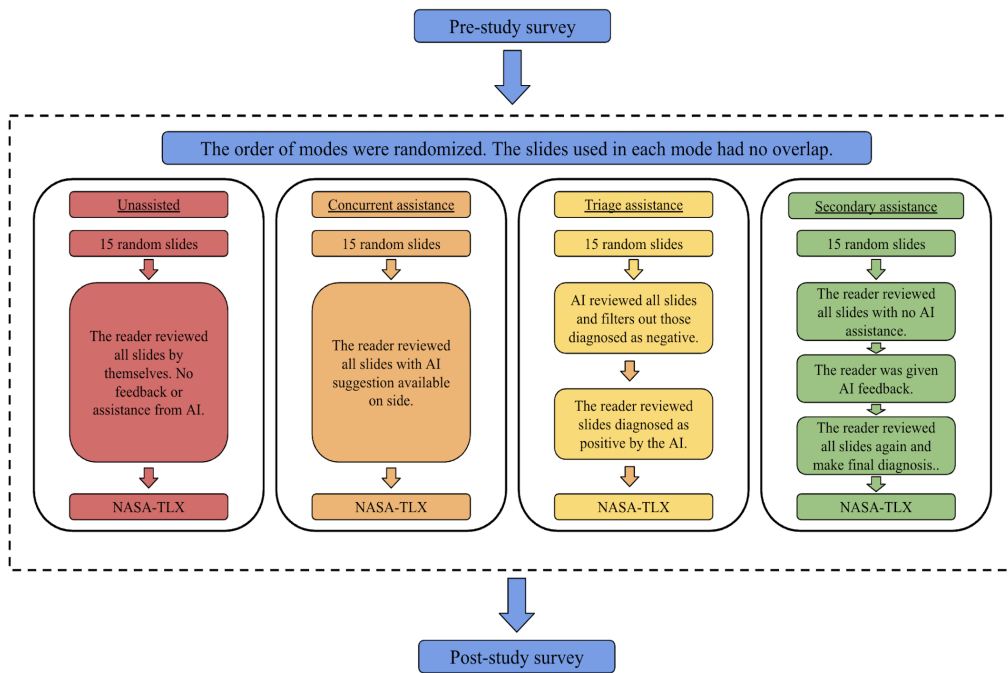


Figure 1: Diagram depicting the full experimental timeline, including the sequence of events across the entire study and a breakdown of the four individual experimental blocks experienced by each participant.

148 review the current slide and diagnose it as negative (NILM) or one of the  
149 seven types of positive cellular lesion (ASC-US, LSIL, ASC-H, HSIL, SCC,  
150 AGC, and AIS) based on the TBS criterion (Nayar & Wilbur, 2015), with or  
151 without AI assistance.

152 When there was no AI assistance, the participant could only see the slide  
153 and the diagnostic options. They could zoom in and out, pan across the slide,  
154 and add annotations using shapes or text in various colors. Additionally,  
155 they could adjust gamma, contrast, brightness, and RGB channel settings as  
156 needed. The interface also allowed them to measure lesion diameters, capture  
157 screenshots, and navigate between slides.

158 When there was AI assistance, participants viewed the AI's general diag-  
159 nosis (positive or negative) with an associated probability (0–100%). The AI  
160 would also cut all suspicious regions from the slide, group them by the types  
161 of abnormality, and sort them from highest to lowest risk. All participants  
162 were introduced to and trained in using the diagnostic system before data col-  
163 lection through an online live lecture and a tutorial video about the use of this  
164 system. About 27% of the participants self-reported that they consistently  
165 used an AI diagnostic assistance tool at least once per month. Regardless of  
166 prior AI experience, to ensure that participants with no prior experience us-  
167 ing AI in cytopathology were adequately prepared, we conducted a practice  
168 session for each participant, during which they reviewed 20 slides in each of  
169 the four AI assistance modes prior to the start of the experimental portion  
170 of the study.

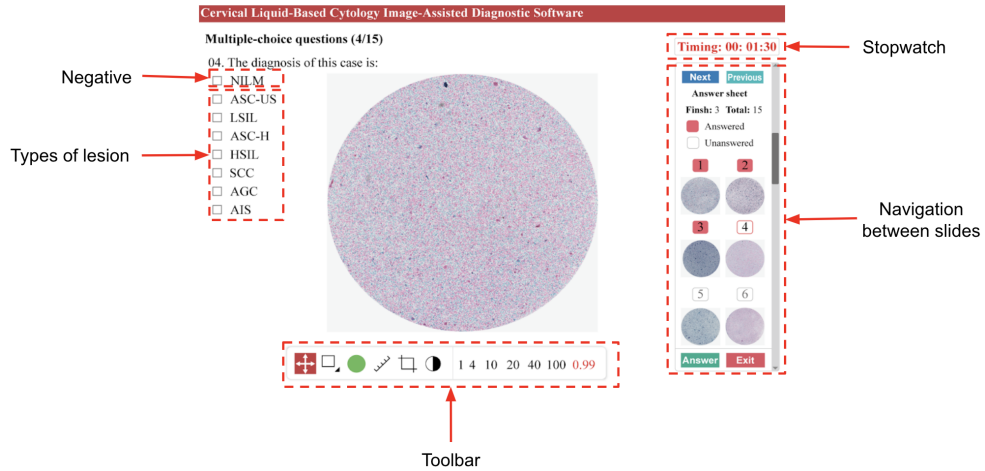


Figure 2: Diagnostic user interface without AI assistance. This interface was used in the unassisted mode, triage assistance mode, and the independent diagnosis stage of the secondary assistance mode.

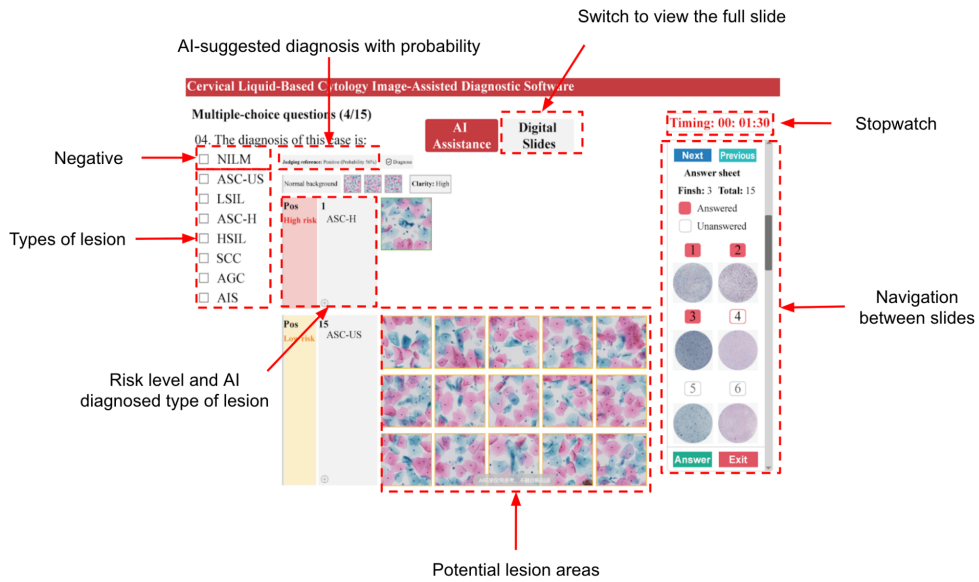


Figure 3: Diagnostic user interface with AI assistance. This interface was used in the concurrent assistance mode and the proof-reading stage of secondary assistance mode.

### 171 3. Results

172 F1-score (calculated as the harmonic mean of precision and recall) across  
173 different AI assistance modes were analyzed using a repeated-measures anal-  
174 ysis of variance (ANOVA). Greenhouse-Geisser corrected  $p$ -values were re-  
175 ported to account for violations of sphericity. Post-hoc pairwise comparisons  
176 were conducted using Tukey’s Honestly Significant Difference (HSD) test  
177 with adjusted  $p$ -values. A two-tailed alpha level of 0.05 was used as the  
178 threshold for statistical significance. The same procedure was applied to the  
179 combined score and subscores of the NASA-TLX questionnaire. For subjec-  
180 tive preference rankings, nonparametric statistical methods were employed.  
181 Specifically, a Friedman test was used to assess differences in rankings of the  
182 AI assistance modes. Pairwise comparisons were conducted with Bonferroni  
183 corrections to account for multiple testing. Finally, the Wilcoxon signed-  
184 rank test was conducted to assess the significance of differences in survey  
185 responses before and after the experiment.

#### 186 3.1. Performance

187 Across the four AI assistance modes, a total of 1,620 slides (1414 negatives  
188 and 206 positives) were analyzed across 108 participants. For performance  
189 assessment, we only evaluate whether the final diagnosis was correctly clas-  
190 sified as positive or negative, without verifying the correctness of identifying  
191 the specific type among the seven positive lesion types. The AI-assistance  
192 mode significantly impacted the f1-score ( $F(3,321)=44.17$ ,  $p<.001$ ,  $\eta_g^2=0.20$ ).  
193 The f1-score in the unassisted mode (M=0.41, SD=0.26) was significantly  
194 lower than in concurrent (M=0.59, SD=0.24,  $M_{\text{diff}}=-0.18$ ,  $p<.001$ ), secondary

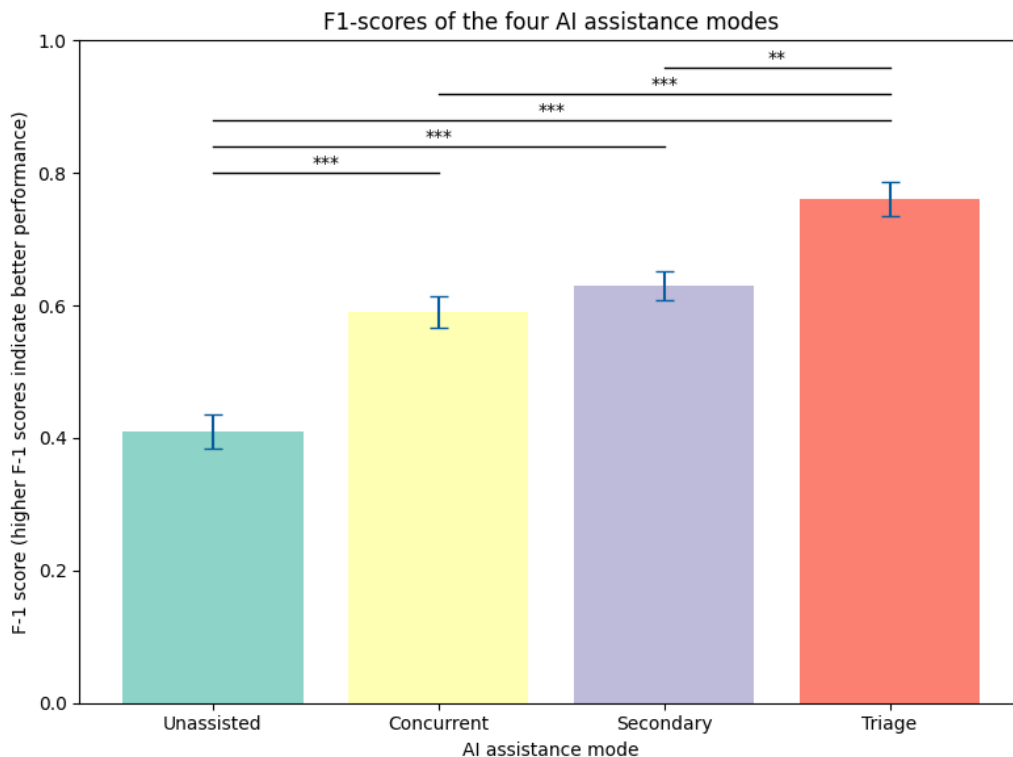


Figure 4: Bar plot depicting the mean values and standard errors of f1-scores as a function of AI assistance modes. \*\* indicates the comparison was significant at the  $p < .01$  level, and \*\*\* indicates the comparison was significant at the  $p < .001$  level.

195 (M=0.63, SD=0.22,  $M_{\text{diff}}=-0.23$ ,  $p < .001$ ), and triage (M=0.76, SD=0.27,  
 196  $M_{\text{diff}}=-0.35$ ,  $p < .001$ ) modes. And the f1-score in the triage mode was also  
 197 significantly higher than in the secondary mode ( $M_{\text{diff}}=0.13$ ,  $p = .002$ ) and the  
 198 concurrent mode ( $M_{\text{diff}}=0.17$ ,  $p < .001$ ).

199 To assess the reliability of AI suggestions from the readers' perspective,  
 200 we evaluated the positive predictive value (PPV, the proportion of true posi-  
 201 tives among all predicted positives) and the negative predictive value (NPV,  
 202 the proportion of true negatives among all predicted negatives) in the three

203 AI-assisted modes. Additionally, we calculated the human-AI diagnostic mis-  
204 match rate, defined as the percentage of cases where the reader's and AI's  
205 diagnoses differed out of the total diagnoses. The repeated-measures ANOVA  
206 showed no significant difference in PPV or NPV, but a significant difference  
207 in human-AI diagnostic mismatch rate ( $F(2,214)=23.87, p<.001, \eta_g^2=0.12$ )  
208 across the three AI-assisted modes. For **concurrent mode**, the AI exhibited  
209 a mean PPV of 51% and a mean NPV of 98%. The mean human-AI diag-  
210 nostic mismatch rate was 22.1%. For **secondary mode**, the AI exhibited a  
211 mean PPV of 52% and a mean NPV of 98%. The mean human-AI diagnostic  
212 mismatch rate after AI assistance was 22%, which was not significantly dif-  
213 ferent than in the concurrent mode. For the **triage mode**, the AI exhibited  
214 a mean PPV of 52% and a mean NPV of 98%. The AI pre-screened 1,216  
215 slides as negative, leaving 404 slides for human review. The mean human-  
216 AI mismatch occurred in 37% of the diagnoses for these slides, which was  
217 significantly higher than in the concurrent mode ( $M_{\text{diff}}=0.15, p<.001$ ) and  
218 secondary mode ( $M_{\text{diff}}=0.15, p<.001$ ).

### 219 3.2. NASA-TLX

220 The changes in AI assistance modes significantly affected overall work-  
221 load ( $F(3,321)=59.71, p<.001, \eta_g^2=0.21$ ). Overall workload in the unas-  
222 sisted mode ( $M=6.72, SD=1.64$ ) was significantly higher than in the concu-  
223 rent ( $M=4.77, SD=1.61, M_{\text{diff}}=1.96, p<.001$ ), secondary ( $M=5.29, SD=1.49,$   
224  $M_{\text{diff}}=1.43, p<.001$ ), and triage ( $M=4.69, SD=1.57, M_{\text{diff}}=2.03, p<.001$ )  
225 modes. And the combined workload in the secondary mode was signifi-  
226 cantly higher than in the triage mode ( $M_{\text{diff}}=0.59, p=0.03$ ). The changes  
227 in AI assistance modes also had significant main effects on all subscores:

228 mental workload ( $F(3, 321)=50.71, p<.001, \eta_g^2=0.18$ ), physical workload  
229 ( $F(3, 321)=35.39, p<.001, \eta_g^2=0.14$ ), temporal workload ( $F(3, 321)=37.76,$   
230  $p<.001, \eta_g^2=0.16$ ), performance ( $F(3, 321)=5.38, p=0.001, \eta_g^2=0.03$ ), ef-  
231 fort ( $F(3, 321)=33.47, p<.001, \eta_g^2=0.12$ ), and frustration ( $F(3, 321)=15.34,$   
232  $p<.001, \eta_g^2=0.06$ ). The effects of AI assistance modes in the subscores were  
233 generally consistent with those observed in the combined scores.

234 For the **mental dimension**, the unassisted mode ( $M=8.31, SD=1.92$ )  
235 had significantly higher ratings than the concurrent ( $M=5.83, SD=2.42,$   
236  $M_{\text{diff}}=2.47, p<.001$ ), secondary ( $M=6.49, SD=2.38, M_{\text{diff}}=1.81, p<.001$ ),  
237 and triage ( $M=5.42, SD=2.70, M_{\text{diff}}=2.89, p<.001$ ) modes, and the sec-  
238 ondary mode also had significantly higher ratings than the triage ( $M_{\text{diff}}=1.07,$   
239  $p=0.005$ ) mode. For the **physical dimension**, the unassisted mode ( $M=7.13,$   
240  $SD=2.72$ ) had significantly higher ratings than the concurrent ( $M=4.78,$   
241  $SD=2.22, M_{\text{diff}}=2.35, p<.001$ ), secondary ( $M=5.38, SD=2.41, M_{\text{diff}}=1.75,$   
242  $p<.001$ ), and triage ( $M=4.61, SD=2.46, M_{\text{diff}}=2.52, p<.001$ ) modes. For the  
243 **temporal dimension**, the unassisted mode ( $M=6.79, SD=2.67$ ) had sig-  
244 nificantly higher ratings than the concurrent ( $M=4.25, SD=2.35, M_{\text{diff}}=2.54,$   
245  $p<.001$ ), secondary ( $M=5.24, SD=2.59, M_{\text{diff}}=1.55, p<.001$ ), and triage ( $M=4.14,$   
246  $SD=2.30, M_{\text{diff}}=2.65, p<.001$ ) modes, and the secondary mode also had sig-  
247 nificantly higher ratings than the triage ( $M_{\text{diff}}=1.10, p=0.007$ ) and concurrent  
248 ( $M_{\text{diff}}=0.99, p=0.018$ ) modes.

249 For the **performance dimension**, the unassisted mode ( $M=6.22, SD=2.38$ )  
250 had significantly lower ratings than the concurrent ( $M=7.01, SD=1.95, M_{\text{diff}}=-$   
251  $0.79, p=0.026$ ) and secondary ( $M=7.09, SD=1.77, M_{\text{diff}}=-0.87, p=0.010$ )  
252 modes, but not the triage mode ( $M=6.67, SD=2.05$ ). For the **effort dimen-**

253 **sion**, the unassisted mode (M=8.03, SD=2.13) had significantly higher rat-  
254 ings than the concurrent (M=6.09, SD=2.41,  $M_{\text{diff}}=1.94$ ,  $p<.001$ ), and triage  
255 (M=5.81, SD=2.57,  $M_{\text{diff}}=2.22$ ,  $p<.001$ ), secondary (M=6.77, SD=2.31,  $M_{\text{diff}}=1.26$ ,  
256  $p<.001$ ), and the secondary mode also had significantly higher ratings than  
257 the triage mode ( $M_{\text{diff}}=0.96$ ,  $p=0.015$ ). For the **frustration dimension**, the  
258 unassisted mode (M=5.30, SD=2.75) had significantly higher ratings than  
259 the concurrent (M=3.65, SD=2.40,  $M_{\text{diff}}=1.65$ ,  $p<.001$ ), secondary (M=3.94,  
260 SD=2.45,  $M_{\text{diff}}=1.35$ ,  $p<.001$ ), and triage (M=3.86, SD=2.46,  $M_{\text{diff}}=1.44$ ,  
261  $p<.001$ ) modes.

### 262 3.3. Preference on Modes of AI Assistance

263 After aggregating the rankings and calculating the mean for each AI  
264 assistance mode, concurrent was the most preferred (mean ranking=1.95),  
265 followed by triage (2.11), secondary (2.44), and finally, unassisted (3.49) was  
266 the least preferred. The same ranking was obtained using the Kemeny-Young  
267 method, and the consensus ranking aligned with the majority preference in  
268 five out of the six possible pairwise comparisons.

269 Rankings differed significantly with respect to the AI assistance mode,  
270  $\chi^2(3)=92.94$ ,  $p<.001$ . The unassisted mode was ranked significantly lower  
271 than the concurrent mode ( $p<.001$ ), triage mode ( $p<.001$ ), and secondary  
272 ( $p<.001$ ) mode. The concurrent mode was ranked significantly higher than  
273 the secondary mode ( $p=.009$ ), but not significantly higher than the triage  
274 mode. The rank difference between the triage and secondary modes was not  
275 significant.

276 Participants were also asked to compare their preferences between the  
277 unassisted mode and the AI-assisted modes. Ninety percent of the partici-

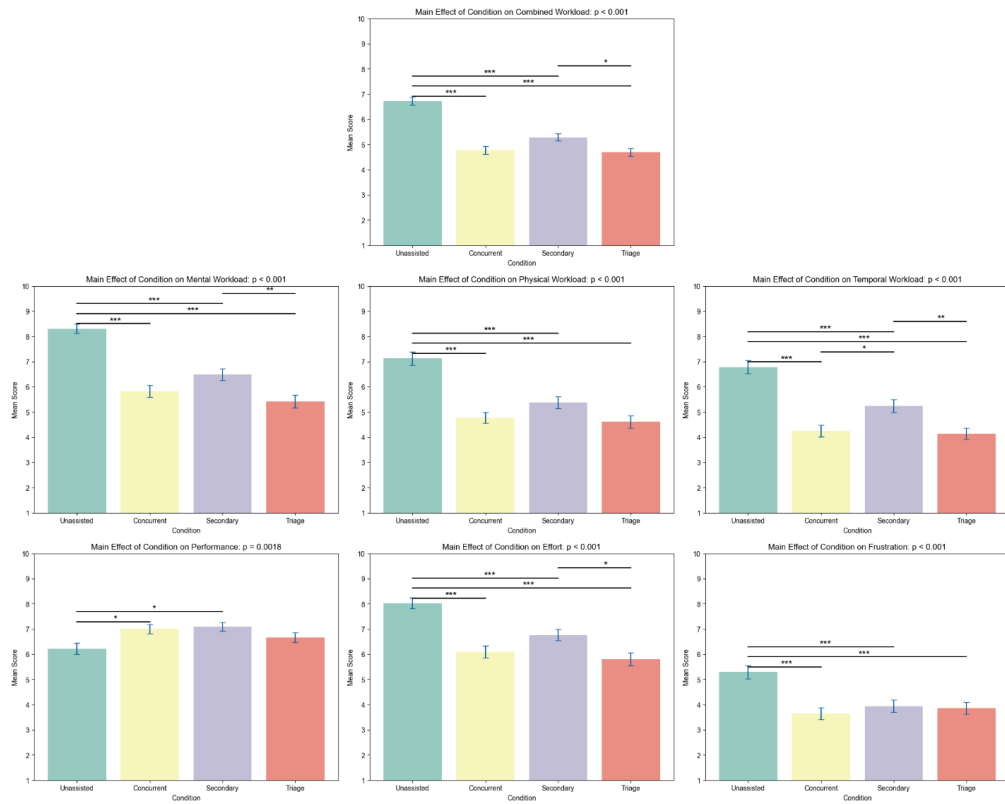


Figure 5: Bar plot depicting the mean values and standard errors of combined score and subscores of NASA-TLX as a function of AI assistance modes. The performance subscore was unreversed, where a higher score indicates better perceived performance. \* indicates the comparison was significant at the  $p < .05$  level, \*\* indicates the comparison was significant at the  $p < .01$  level, and \*\*\* indicates the comparison was significant at the  $p < .001$  level.

278 pants favored the AI-assisted modes for their efficiency, 85% found them to  
279 reduce difficulty of completing their task, and 85% reported greater confi-  
280 dence in their diagnoses when using the AI-assisted modes.

### 281 3.4. Trust on AI Assistance

282 After the experimental portion of the study, participants answered 13  
283 questions to assess their trust on the AI assistance tool. All responses were  
284 measured using the Likert scale, with 1 indicating "strongly disagree" and 5  
285 indicating "strongly agree." After inverting the responses for reverse-scored  
286 questions, participants' mean trust score was 3.27 with a standard deviation  
287 of 0.50, indicating an overall neutral to slightly positive level of trust in the  
288 AI. The score distribution of each question is shown in Figure 6 *Left*.

289 Four of the questions were asked both before and after the experiments  
290 to understand how the use of AI affected participants' trust in it. For the  
291 question "*I am confident in the performance of AI,*" the post-experiment  
292 score (M=4.07, SD=0.65) was similar to the pre-experiment score (M=4.01,  
293 SD=0.72,  $M_{\text{diff}}=0.06$ ). For the question "*I feel safe to let the AI help me*  
294 *skip the preliminary reading with a microscope,*" the post-experiment score  
295 (M=3.52, SD=0.91) was significantly lower than the pre-experiment score  
296 (M=3.83, SD=0.79,  $M_{\text{diff}}=-0.31$ ,  $p=0.002$ ). For the question "*I feel safe to let*  
297 *the AI make the diagnosis on its own,*" the post-experiment score (M=2.81,  
298 SD=1.08) was significantly lower than the pre-experiment score (M=3.33,  
299 SD=1.03,  $M_{\text{diff}}=-0.52$ ,  $p<.001$ ). Finally, for the question "*I am wary of the*  
300 *AI tool,*" the post-experiment score (M=3.64, SD=0.89) was similar to the  
301 pre-experiment score (M=3.66, SD=0.87,  $M_{\text{diff}}=-0.02$ ).

302 *3.5. Perceived usability of AI Assistance*

303 After the experimental portion of the study, the SUS questionnaire with  
304 ten questions were answered by each participant to assess their perceived us-  
305 ability of the AI assistance tool. All responses were measured using the Likert  
306 scale, with 1 indicating "strongly disagree" and 5 indicating "strongly agree."  
307 After inverting the responses for reverse-scored questions, participants' mean  
308 usability score was 3.65 with a standard deviation of 0.46, indicating an over-  
309 all neutral to slightly positive level of perceived usability of the AI. The score  
310 distribution of each question is shown in Figure 6 *Right*.

311 Three of the questions were asked both before and after the experiments  
312 to understand how the use of AI affected the perceived usability. For the  
313 question "I think that I would like to use this system frequently," the post-  
314 experiment score (M=4.06, SD=0.73) was similar to the pre-experiment score  
315 (M=4.1, SD=0.61,  $M_{\text{diff}}=-0.04$ ). For the question "I think that I would  
316 need the support of a technical person to be able to use this system," the  
317 post-experiment score (M=3.03, SD=1.22) was significantly lower than the  
318 pre-experiment score (M=3.62, SD=1.16,  $M_{\text{diff}}=-0.59$ ,  $p<.001$ ). And for the  
319 question "I would imagine that most people would learn to use this system  
320 very quickly," the post-experiment score (M=4.17, SD=0.59) was similar to  
321 the pre-experiment score (M=4.21, SD=0.61,  $M_{\text{diff}}=-0.04$ ).

322 *3.6. Open-ended Feedback*

323 After the experimental portion of the study, participants were asked to  
324 describe how their reading and diagnostic strategies differed when using AI  
325 assistance versus not using it, as well as to provide general feedback on the  
326 study and the AI tool.

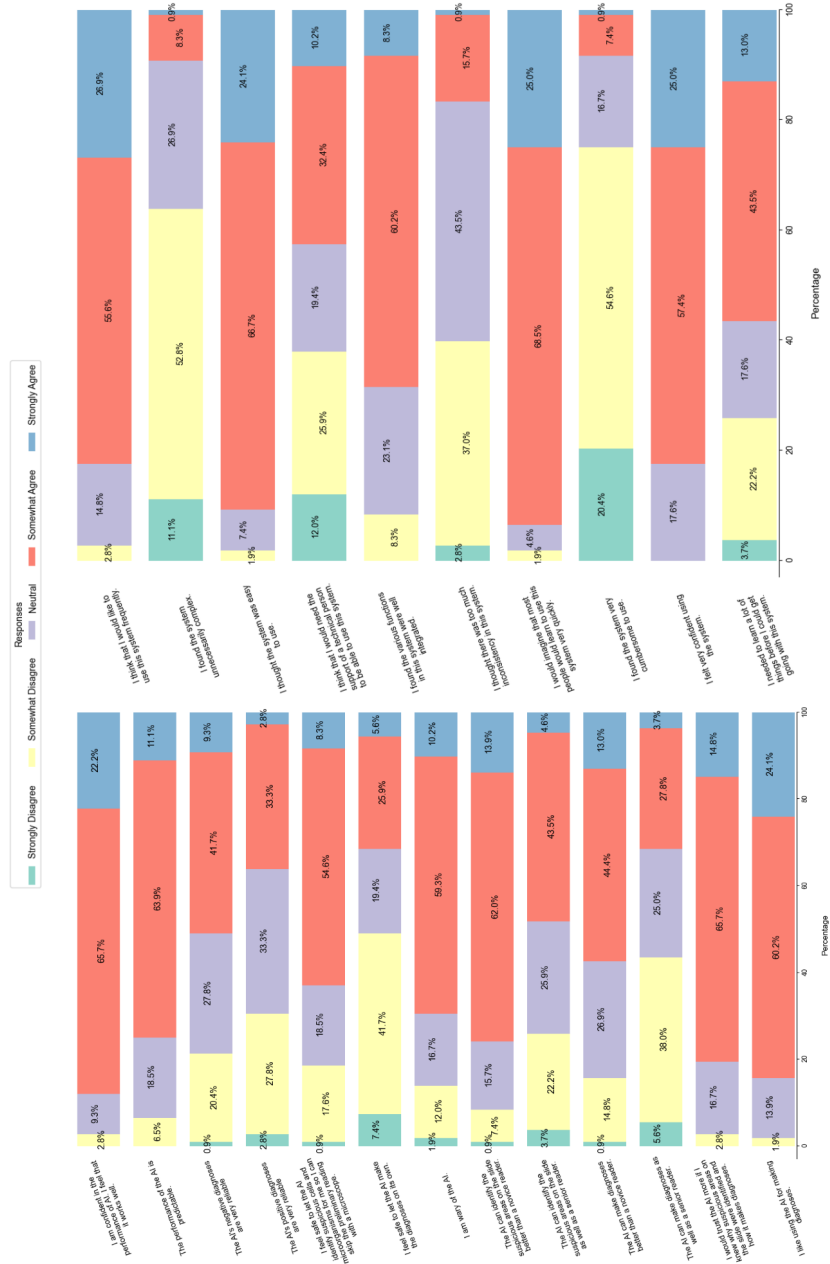


Figure 6: Post-experiment survey responses. *Left*: Responses to questions assessing trust in the AI system. *Right*: Responses to questions evaluating system usability.

327 Compared to self-diagnosis, the three most frequently mentioned advan-  
328 tages of AI-assisted slide reading were higher efficiency, shorter diagnostic  
329 time, and greater confidence in the diagnosis. With AI pre-filtering nega-  
330 tive slides or providing suggested diagnoses, participants were able to review  
331 fewer slides or approach them more strategically, guided by AI recommenda-  
332 tions. Participants mentioned that the use of AI-enhanced their confidence  
333 in their diagnoses, especially when AI confirmed their findings. Some partic-  
334 ipants expressed comfort in relying on AI when available, highlighting that  
335 the time saved by AI-assisted diagnostic suggestions allowed them to exam-  
336 ine slides in greater detail and reduced their stress levels. Additionally, some  
337 participants noted that AI support could help prevent missed positive cases  
338 and believed that its integration improved the overall accuracy in identifying  
339 positive slides.

340 However, concerns regarding AI usage were also raised. Some participants  
341 observed that receiving AI suggestions prior to diagnosis could occasionally  
342 be misleading. When their diagnoses conflicted with the AI's, they reported  
343 feelings of anxiety, reduced confidence, and uncertainty about the appro-  
344 priate course of action. A few participants criticized the AI tool for being  
345 overly conservative, resulting in more errors in predicted positives than antic-  
346 ipated. And they perceived the AI as more reliable when suggesting negative  
347 diagnoses. Others mentioned that relying on AI made them less cautious,  
348 thorough, or alert when reviewing slides. One participant noted that their  
349 worry about AI's reliability increased their anxiety compared to diagnosing  
350 independently. Finally, 26 participants either did not report any differences  
351 or explicitly stated that there was no discernible difference between diagnos-

352 ing with and without AI assistance.

353 Participants provided many valuable suggestions and feedback regarding  
354 the AI assistance tool. A common theme among participants was the desire  
355 for greater transparency in AI's diagnostic logic and interpretability of the  
356 probabilities shown, which they believed would enable them to better eval-  
357 uate its suggestions and calibrate their trust accordingly. Participants also  
358 proposed several usability enhancements, such as incorporating features to  
359 highlight areas with pathological changes, displaying additional patient infor-  
360 mation (e.g., age and medical history), and developing a mobile application.  
361 Moreover, some participants emphasized the importance of training the AI  
362 on a broader dataset covering diverse pathological variations to enhance its  
363 robustness. One participant suggested that the AI should focus on highlight-  
364 ing areas of pathological changes rather than delivering definitive diagnoses.  
365 Several participants raised questions regarding how readers could compare  
366 the performance of different AI tools, identify the weaknesses of the AI, and  
367 calibrate its false-positive rate. Lastly, many participants expressed support  
368 for future promotion of this AI tool and noted its potential to be particularly  
369 beneficial in contexts where large volumes of slides must be analyzed with  
370 limited manpower.

#### 371 **4. Discussion**

372 Compared to the unassisted mode, the three AI-assisted modes signifi-  
373 cantly improved performance and reduced workload. The AI-assisted modes  
374 were significantly more preferred than the unassisted mode, and most par-  
375 ticipants believed that AI assistance improved diagnostic efficiency and con-

376 fidence while decreasing diagnostic difficulty. Furthermore, this study ex-  
377 amined the timing of the AI suggestions, and our findings suggested that  
378 the choice of timing should take into account the specific use case and needs  
379 of the implementation. The concurrent and secondary modes had similar  
380 performance and workload scores. There was no significant difference in  
381 f1-score or human-AI diagnostic mismatch rate between providing AI infor-  
382 mation side by side and using it only for a second round of review. The  
383 combined score and subscores of the NASA-TLX also showed no significant  
384 differences between the two modes, except for temporal workload, where  
385 readers reported feeling more rushed using the secondary mode. This may  
386 be attributed to the additional time required for cross-checking AI recom-  
387 mendations after the task was completed. The concurrent mode was ranked  
388 significantly higher in preference than the secondary mode. However, adding  
389 an independent review round prior to presenting the AI recommendation  
390 may still be meaningful in certain situations. Some participants noted that  
391 viewing AI suggestions before starting the diagnosis could unintentionally  
392 bias their diagnostic approach. This anchoring effect of the initial informa-  
393 tion presented has been documented in psychological research (Englich et  
394 al., 2006) and healthcare settings (Branch et al., 2022). Additionally, partic-  
395 ipants mentioned that prolonged AI use may lead to negative consequences  
396 associated with overreliance, such as reduced caution and thoroughness (Li &  
397 Little, 2023; Hatherley, 2020). Alternatively, as some participants suggested,  
398 AI could ideally highlight potential lesion areas without providing diagnostic  
399 information, or its outputs could be hidden for users to reveal manually.

400 In the triage mode, AI pre-filtered 75% of the samples it diagnosed as

401 negative, substantially reducing readers' task load. This likely explains why  
402 the mental workload, temporal workload, and effort scores in this mode  
403 were significantly lower than in the secondary mode. However, none of  
404 the NASA-TLX subscores differed significantly between the concurrent and  
405 triage modes. This may be because the number of samples each reader re-  
406 viewed was relatively small. Given the AI's negative diagnoses were highly  
407 reliable (with a negative predictive value exceeding 98%), the pre-filtering  
408 led to significantly better human-AI collaboration performance, as reflected  
409 by the higher f1-score compared to the other two AI-assisted modes. These  
410 results suggest that AI can substantially reduce readers' task load while en-  
411 hancing performance.

412 However, several issues remain with the triage mode. The positive slides  
413 identified by the AI for readers to review were only around 50% accurate,  
414 resulting in significantly higher human-AI diagnostic mismatch rate than in  
415 the other two AI-assisted modes. This could undermine readers' confidence  
416 in their performance (Kempton et al., 2023), potentially explaining why the  
417 triage mode was the only AI-assisted mode where the NASA-TLX perfor-  
418 mance subscore was not significantly higher than in the unassisted mode.  
419 The findings also highlighted concerns regarding the AI's reliability. Readers  
420 reported increased anxiety when their diagnoses conflicted with the AI's, par-  
421 ticularly when the AI's diagnosis was positive, given the severe consequences  
422 of false negatives in healthcare settings. Additionally, in the triage mode,  
423 readers could not see the suspicious areas marked by the AI. While this was  
424 intended to reduce bias that could result from the AI recommendation, it  
425 also diminished the interpretability and utility of the AI's diagnoses. More

426 importantly, readers had no control over the slides filtered by the AI, even  
427 though they could still be held accountable for false negatives missed by the  
428 AI (Wang et al., 2024). Finally, while the triage mode was ranked between  
429 the concurrent and secondary modes, its ranking did not significantly differ  
430 from either, making it unclear whether readers preferred the triage mode  
431 over the other AI-assisted modes.

432 The post-study survey results revealed that the participants had a neutral  
433 to slightly positive level of trust in and perceived usability of the AI assis-  
434 tance tool after the experiment, and they expressed support for its future  
435 promotion based on the responses to the open-ended questions. The survey  
436 results also indicated that the participants felt more comfortable using the  
437 AI to identify suspicious regions on slides than allowing it to make final di-  
438 agnostic decisions. Users being wary of granting full control to the AI has  
439 also been reported in previous studies (Yang et al., 2023; Sivaraman et al.,  
440 2023; Burgess et al., 2023). Responses to the four sub-questions indicated  
441 that participants agreed more strongly that the AI could outperform novice  
442 readers than that it could perform as well as senior readers, both in initial  
443 slide review and final diagnosis. Average ratings for the four sub-questions  
444 were all above neutral, except for the one concerning the AI performing as  
445 well as senior readers in final decision-making. Participants also agreed more  
446 that the AI’s suggestions were more reliable for negative samples than for  
447 positive ones, which aligned with the fact that NPV was consistently nearly  
448 twice as high as PPV across all three AI-assisted modes. Additionally, most  
449 readers agreed that they would trust the AI more if it provided more expla-  
450 nations of its working mechanism and decision process, which aligned with

451 findings in previous work regarding AI explainability (Nasarian et al., 2024;  
452 Yang et al., 2023). Regarding usability, the majority of the readers found  
453 the system easy to use, with well-integrated functions and minimal unneces-  
454 sary complexity or inconsistency. They also believed that while substantial  
455 training would be required initially to learn the nuances associated with the  
456 system, users could quickly become proficient with the system. By com-  
457 paring the results from the pre- and post-study surveys, we observed that  
458 readers recalibrated their expectation of reliance on the AI tool after the ex-  
459 perimental portion of the study. The calibrated level of trust could depend  
460 on the AI’s performance, task difficulty, and whether there were human-AI  
461 diagnostic mismatches identified (Cao & Huang, 2022; Lee, 2024). Partici-  
462 pants’ agreement with relying on the AI for both initial diagnosis and final  
463 decision-making significantly decreased after the study.. However, the study  
464 did not significantly alter their general attitude toward the AI system after  
465 using it for this study. This was reflected by the two questions assessing con-  
466 fidence and wariness toward the AI before and after the experimental portion  
467 of the study that remained statistically the same. For usability, the belief  
468 that technicians’ support would be needed was significantly lower after the  
469 experiment. Additionally, the belief that most people could quickly learn to  
470 use the system, as well as the desire to use the system frequently, remained  
471 unchanged.

## 472 **5. Conclusion**

473 In this study, we assessed three timings for providing AI assistance to  
474 readers during pathological slide diagnosis: pre-diagnosis for pre-filtering

475 negative slides (triage mode), concurrent assistance during diagnosis (concur-  
476 rent mode), and post-diagnosis for proofreading (secondary mode). All three  
477 modes demonstrated significant improvements in diagnostic performance and  
478 workload reduction compared to unassisted diagnosis; however, the ideal  
479 mode used should take into account the clinical needs. For scenarios priori-  
480 tizing substantial workload reduction and high performance, such as during  
481 person-power shortages, we would recommend the triage mode. Conversely,  
482 when prioritizing maintaining the readers' control over decision-making, the  
483 concurrent or secondary modes are more suitable. Between these two, the  
484 secondary mode may better fit in contexts requiring thorough reader deliber-  
485 ation or minimizing AI bias, in which may be paramount for training. How-  
486 ever, the concurrent mode is generally more efficient, as it provides real-time  
487 assistance without requiring additional review rounds. And the concurrent  
488 mode was preferred to the secondary mode. Notably, participants did not  
489 significantly favor the triage mode over the other two AI modes, likely be-  
490 cause they saw the AI as better than a novice reader, but not a substitute  
491 for a senior reader. Additionally, they viewed AI as more appropriate for  
492 preliminary slide review rather than final diagnosis. These results speak di-  
493 rectly to broader debates in clinical decision support, where efficiency gains  
494 must be balanced against bias, explainability, and clinician accountability.  
495 Our findings on timing modes suggest that AI integration should be tailored  
496 to diagnostic contexts—for example, triage for high-volume screening set-  
497 tings versus concurrent or secondary support for contexts that emphasize  
498 individualized decision-making and thorough review. This will hinge on ap-  
499 propriately including the AI decision support tool into the workflow so that

500 users know when, where, and how to use the tool.

501 The AI assistance system was generally well received by users, with mod-  
502 erate positive feedback on both usability and trust following the experiment.  
503 Designers must carefully calibrate the amount of AI-generated information  
504 presented to avoid undue bias and balance the system’s conservatism to mini-  
505 mize false negatives without excessive false positives, which could undermine  
506 reader confidence. Additionally, enhancing the explainability of the AI’s  
507 decision-making processes and recommendations could further foster trust in  
508 the system. Our results also showed that clinicians’ trust and reliance on  
509 AI were dynamic and shaped by experience. This meant that, in addition  
510 to maintaining and improving performance, trust in AI-enabled clinical de-  
511 cision support systems should be continuously monitored to guard against  
512 both over-trust, which risks complacency, and under-trust, which reduces  
513 potential benefits.

514 Future studies may explore the factors contributing to individual differ-  
515 ences in subjective measures, such as preferences on AI assistance modes and  
516 trust in AI, observed in this study. One potential consideration is the user’s  
517 AI literacy. Previous work has found that having a better understanding of  
518 AI may enhance confidence in these systems, particularly with targeted train-  
519 ing (Hua et al., 2024). Additionally, because slide assignments were random-  
520 ized, readers’ experiences collaborating with AI varied, potentially affecting  
521 their perceptions and feedback provided after the study. Another avenue for  
522 future research is the role of cultural factors in shaping readers’ perceptions of  
523 AI. As this study was conducted in China, cross-cultural comparisons could  
524 provide valuable insights into how cultural norms and attitudes toward AI

525 influence performance, assistance mode preference, and overall trust in AI  
526 assistance. Finally, our findings suggest that readers calibrated their trust  
527 and expectations of AI after the experiment. Future work could investigate  
528 how readers' attitudes toward AI fluctuate throughout the study, particu-  
529 larly in response to specific types of human-AI diagnostic mismatches that  
530 may significantly weaken or reinforce trust in AI. Understanding these dy-  
531 namics is essential for designing AI tools that promote appropriate trust and  
532 effective human-AI collaboration.

533 **References**

- 534 Abraham, S., & Joseph, S., *Medical Imaging and Artificial Intelligence:*  
535 *Transforming the Nature of Diagnostics and Treatment*, In H. Joshi et  
536 al. (Eds.), *Advances in Medical Technologies and Clinical Practice* (pp.  
537 127–158), IGI Global, 2024.
- 538 Ajmera, P., Onkar, P., Desai, S., Pant, R., Seth, J., Gupte, T., et al., *Val-*  
539 *idation of a Deep Learning Model for Detecting Chest Pathologies from*  
540 *Digital Chest Radiographs*, *Diagnostics*, 13(3), Article 3, 2023.
- 541 Ahmadi, Dr. R., *Emerging Innovations in AI-Driven Medical Imaging: Ele-*  
542 *vating Diagnostic Precision and Therapeutic Decision-Making*, *Journal of*  
543 *AI-Powered Medical Innovations*, 1(1), 1–28, 2024.
- 544 Avanzo, M., Stancanello, J., Pirrone, G., Drigo, A., & Retico, A., *The Evolu-*  
545 *tion of Artificial Intelligence in Medical Imaging: From Computer Science*  
546 *to Machine and Deep Learning*, *Cancers*, 16(21), 3702, 2024.
- 547 Baruwal Chhetri, M., Tariq, S., Singh, R., Jalalvand, F., Paris, C., & Nepal,  
548 S., *Towards human-AI teaming to mitigate alert fatigue in security op-*  
549 *erations centres*, *ACM Transactions on Internet Technology*, 24(3), 1–22,  
550 2024.
- 551 Bansal, G., Wu, T., Zhou, J., Fok, R., Nushi, B., Kamar, E., et al., *Does the*  
552 *whole exceed its parts? The effect of AI explanations on complementary*  
553 *team performance*, *Proc. CHI Conf. Human Factors in Computing Systems*,  
554 1–16, 2021.

- 555 Branch, F., Santana, I., & Hegdé, J., *Biasing influence of ‘Mental Shortcuts’*  
556 *on diagnostic decision-making: radiologists can overlook breast cancer in*  
557 *mammograms when prior diagnostic information is available*, *Diagnostics*,  
558 12(1), 105, 2022.
- 559 Brooke, J., *SUS – A quick and dirty usability scale*, *Usability Evaluation in*  
560 *Industry*, 189(194), 4–7, 1996.
- 561 Burgess, E. R., Jankovic, I., Austin, M., Cai, N., Kapuścińska, A., Currie,  
562 S., et al., *Healthcare AI Treatment Decision Support: Design Principles to*  
563 *Enhance Clinician Adoption and Trust*, *Proc. CHI Conf. Human Factors*  
564 *in Computing Systems*, 1–19, 2023.
- 565 Cao, S., & Huang, C. M., *Understanding user reliance on AI in assisted*  
566 *decision-making*, *Proceedings of the ACM on Human-Computer Interac-*  
567 *tion*, 6(CSCW2), 1–23, 2022.
- 568 Chen, M., Wang, Y., Wang, Q., Shi, J., Wang, H., Ye, Z., et al., *Impact*  
569 *of human and artificial intelligence collaboration on workload reduction in*  
570 *medical image interpretation*, *NPJ Digital Medicine*, 7(1), 349, 2024.
- 571 Duron, L., Ducarouge, A., Gillibert, A., Lainé, J., Allouche, C., Cherel, N.,  
572 et al., *Assessment of an AI Aid in Detection of Adult Appendicular Skeletal*  
573 *Fractures by Emergency Physicians and Radiologists: A Multicenter Cross-*  
574 *sectional Diagnostic Study*, *Radiology*, 300(1), 120–129, 2021.
- 575 Dvijotham, K. (Dj), Winkens, J., Barsbey, M., Ghaisas, S., Stanforth,  
576 R., Pawlowski, N., et al., *Enhancing the reliability and accuracy of AI-*

577 *enabled diagnosis via complementarity-driven deferral to clinicians*, Nature  
578 Medicine, 29(7), 1814–1820, 2023.

579 English, B., Mussweiler, T., & Strack, F., *Playing dice with criminal sen-*  
580 *tences: The influence of irrelevant anchors on experts' judicial decision*  
581 *making*, Personality and Social Psychology Bulletin, 32(2), 188–200, 2006.

582 Farber, S., *Comparing human and AI expertise in the academic peer review*  
583 *process: towards a hybrid approach*, Higher Education Research & Devel-  
584 opment, 1–15, 2025.

585 Frazer, H. M., Peña-Solorzano, C. A., Kwok, C. F., Elliott, M. S., Chen,  
586 Y., Wang, C., et al., *Integrated AI reader development and evaluation pro-*  
587 *vides clinically-relevant guidance for human-AI collaboration in population*  
588 *mammographic screening*, medRxiv, 2023.

589 Guermazi, A., Tannoury, C., Koppel, A. J., Murakami, A. M., Ducarouge,  
590 A., Gillibert, A., et al., *Improving Radiographic Fracture Recognition Per-*  
591 *formance and Efficiency Using Artificial Intelligence*, Radiology, 302(3),  
592 627–636, 2022.

593 Hallowell, N., Badger, S., Sauerbrei, A., Nellåker, C., & Kerasidou, A., “*I*  
594 *don't think people are ready to trust these algorithms at face value*”: *Trust*  
595 *and the use of machine learning algorithms in the diagnosis of rare disease*,  
596 BMC Medical Ethics, 23(1), 112, 2022.

597 Hassan, E. A., & El-Ashry, A. M., *Leading with AI in critical care nursing:*  
598 *challenges, opportunities, and the human factor*, BMC Nursing, 23(1), 752,  
599 2024.

- 600 Hatherley, J. J., *Limits of trust in medical AI*, Journal of Medical Ethics,  
601 46(7), 478–481, 2020.
- 602 Hauptman, A. I., Mallick, R., Flathmann, C., & McNeese, N. J., *Human*  
603 *factors considerations for the context-aware design of adaptive autonomous*  
604 *teammates*, Ergonomics, 1–17, 2024.
- 605 He, G., Kuiper, L., & Gadiraju, U., *Knowing about knowing: An illusion of*  
606 *human competence can hinder appropriate reliance on AI systems*, Proc.  
607 CHI Conf. Human Factors in Computing Systems, 1–18, 2023.
- 608 Hendrix, N., Hendrix, W., van Dijke, K., Maresch, B., Maas, M., Bollen, S.,  
609 et al., *Musculoskeletal radiologist-level performance by using deep learning*  
610 *for detection of scaphoid fractures on conventional multi-view radiographs*  
611 *of hand and wrist*, European Radiology, 33(3), 1575–1588, 2023.
- 612 Hoffman, R. R., Mueller, S. T., Klein, G., & Litman, J., *Measuring trust*  
613 *in the XAI context (Technical Report)*, DARPA Explainable AI Program,  
614 2018.
- 615 Hua, D., Petrina, N., Young, N., Cho, J.-G., & Poon, S. K., *Understand-*  
616 *ing the factors influencing acceptability of AI in medical imaging domains*  
617 *among healthcare professionals: A scoping review*, Artificial Intelligence in  
618 Medicine, 147, 102698, 2024.
- 619 Kanse, A. S., Kurian, N. C., Aswani, H. P., Khan, Z., Gann, P. H., Rane, S.,  
620 & Sethi, A., *Cautious Artificial Intelligence Improves Outcomes and Trust*  
621 *by Flagging Outlier Cases*, JCO Clinical Cancer Informatics, 6, e2200067,  
622 2022.

- 623 Kempt, H., Heilinger, J. C., & Nagel, S. K., *“I’m afraid I can’t let you do*  
624 *that, Doctor”*: meaningful disagreements with AI in medical contexts, *AI*  
625 *& Society*, 38(4), 1407–1414, 2023.
- 626 Koivisto, M., & Grassini, S., *Best humans still outperform artificial intelli-*  
627 *gence in a creative divergent thinking task*, *Scientific Reports*, 13(1), 13601,  
628 2023.
- 629 Lai, X., & Rau, P. L. P., *AI as Co-Leader: Effects of Human Considera-*  
630 *tion and AI Structure Behaviors on Leadership Effectiveness and Neural*  
631 *Activation*, *International Journal of Human–Computer Interaction*, 1–19,  
632 2024.
- 633 Lauritzen, A. D., Rodríguez-Ruiz, A., von Euler-Chelpin, M. C., Lyngø,  
634 E., Vejborg, I., Nielsen, M., et al., *An Artificial Intelligence–based Mam-*  
635 *mography Screening Protocol for Breast Cancer: Outcome and Radiologist*  
636 *Workload*, *Radiology*, 304(1), 41–49, 2022.
- 637 Lee, E., *The Power of Perception in Human-AI Interaction: Investigating*  
638 *Psychological Factors and Cognitive Biases that Shape User Belief and*  
639 *Behavior*, arXiv preprint arXiv:2409.15328, 2024.
- 640 Leibig, C., Brehmer, M., Bunk, S., Byng, D., Pinker, K., & Umutlu, L.,  
641 *Combining the strengths of radiologists and AI for breast cancer screening:*  
642 *A retrospective analysis*, *The Lancet Digital Health*, 4(7), e507–e519, 2022.
- 643 Li, M. D., & Little, B. P., *Appropriate reliance on artificial intelligence in*  
644 *radiology education*, *Journal of the American College of Radiology*, 20(11),  
645 1126–1130, 2023.

- 646 Macnamara, B. N., Berber, I., Çavuşoğlu, M. C., Krupinski, E. A., Nalla-  
647 pareddy, N., Nelson, N. E., et al., *Does using artificial intelligence assis-*  
648 *tance accelerate skill decay and hinder skill development without perform-*  
649 *ers' awareness?*, *Cognitive Research: Principles and Implications*, 9(1), 46,  
650 2024.
- 651 Marinovich, M. L., Wylie, E., Lotter, W., Lund, H., Waddell, A., Made-  
652 ley, C., et al., *Artificial intelligence (AI) for breast cancer screen-*  
653 *ing: BreastScreen population-based cohort study of cancer detection*,  
654 *eBioMedicine*, 90, 2023.
- 655 McCague, C., MacKay, K., Welsh, C., Constantinou, A., Jena, R., Crispin-  
656 Ortuzar, M., et al., *Position statement on clinical evaluation of imaging*  
657 *AI*, *The Lancet Digital Health*, 5(7), e400–e402, 2023.
- 658 McKinney, S. M., Sieniek, M., Godbole, V., Godwin, J., Antropova, N.,  
659 Ashrafian, H., et al., *International evaluation of an AI system for breast*  
660 *cancer screening*, *Nature*, 577(7788), 89–94, 2020.
- 661 Tigre Moura, F., *Artificial intelligence, creativity, and intentionality: The*  
662 *need for a paradigm shift*, *The Journal of Creative Behavior*, 57(3),  
663 336–338, 2023.
- 664 Nasarian, E., Alizadehsani, R., Acharya, U. R., & Tsui, K.-L., *Designing in-*  
665 *terpretable ML system to enhance trust in healthcare: A systematic review*  
666 *to proposed responsible clinician-AI-collaboration framework*, *Information*  
667 *Fusion*, 108, 102412, 2024.

- 668 Nayar, R., & Wilbur, D. C., *The Pap test and Bethesda 2014*, *Cancer Cy-*  
669 *topathology*, 123(5), 271–281, 2015.
- 670 Peng, Z., & Ren, X., *Application and Development of Artificial Intelligence-*  
671 *based Medical Imaging Diagnostic Assistance System*, *International Jour-*  
672 *nal of Biology and Life Sciences*, 6(1), 39–43, 2024.
- 673 Poon, A. I. F., & Sung, J. J. Y., *Opening the black box of AI-Medicine*,  
674 *Journal of Gastroenterology and Hepatology*, 36(3), 581–584, 2021.
- 675 Rasheed, J., Jamil, A., Hameed, A. A., Aftab, U., Aftab, J., Shah, S. A., &  
676 Draheim, D., *A survey on artificial intelligence approaches in supporting*  
677 *frontline workers and decision makers for the COVID-19 pandemic*, *Chaos,*  
678 *Solitons & Fractals*, 141, 110337, 2020.
- 679 Raya-Povedano, J. L., Romero-Martín, S., Elías-Cabot, E., Gubern-Mérida,  
680 A., Rodríguez-Ruiz, A., & Álvarez-Benito, M., *AI-based Strategies to Re-*  
681 *duce Workload in Breast Cancer Screening with Mammography and To-*  
682 *mosynthesis: A Retrospective Evaluation*, *Radiology*, 300(1), 57–65, 2021.
- 683 Sivaraman, V., Bukowski, L. A., Levin, J., Kahn, J. M., & Perer, A., *Ign-*  
684 *ore, Trust, or Negotiate: Understanding Clinician Acceptance of AI-*  
685 *Based Treatment Recommendations in Health Care*, *Proc. CHI Conf. Hu-*  
686 *man Factors in Computing Systems*, 1–18, 2023.
- 687 Smith, A., van Wagoner, H. P., Keplinger, K., & Celebi, C., *Navigating AI*  
688 *Convergence in Human–Artificial Intelligence Teams: A Signaling Theory*  
689 *Approach*, *Journal of Organizational Behavior*, 2025.

- 690 Szajna, A., & Kostrzewski, M., *AR-AI tools as a response to high employee*  
691 *turnover and shortages in manufacturing during regular, pandemic, and*  
692 *war times*, Sustainability, 14(11), 6729, 2022.
- 693 Sung, J., Park, S., Lee, S. M., Bae, W., Park, B., Jung, E., et al., *Added Value*  
694 *of Deep Learning-based Detection System for Multiple Major Findings on*  
695 *Chest Radiographs: A Randomized Crossover Study*, Radiology, 299(2),  
696 450–459, 2021.
- 697 Wang, H., Ye, Z., Zhang, P., Cui, X., Chen, M., Wu, A., et al., *Chinese*  
698 *colposcopists' attitudes toward the colposcopic artificial intelligence aux-*  
699 *iliary diagnostic system (CAIADS): A nation-wide, multi-center survey*,  
700 DIGITAL HEALTH, 10, 20552076241279952, 2024.
- 701 Xu, C., Lien, K. C., & Höllerer, T., *Comparing zealous and restrained AI*  
702 *recommendations in a real-world human-AI collaboration task*, Proc. CHI  
703 Conf. Human Factors in Computing Systems, 1–15, 2023.
- 704 Xue, P., Xu, H. M., Tang, H. P., Weng, H. Y., Wei, H. M., Wang, Z., et al.,  
705 *Improving the accuracy and efficiency of abnormal cervical squamous cell*  
706 *detection with cytologist-in-the-loop artificial intelligence*, Modern Pathol-  
707 ogy, 36(8), 100186, 2023.
- 708 Yala, A., Schuster, T., Miles, R., Barzilay, R., & Lehman, C., *A Deep Learn-*  
709 *ing Model to Triage Screening Mammograms: A Simulation Study*, Radi-  
710 ology, 293(1), 38–46, 2019.
- 711 Yang, Q., Hao, Y., Quan, K., Yang, S., Zhao, Y., Kuleshov, V., & Wang,  
712 F., *Harnessing Biomedical Literature to Calibrate Clinicians' Trust in AI*

713 *Decision Support Systems*, Proc. CHI Conf. Human Factors in Computing  
714 Systems, 1–14, 2023.

715 Zhai, C., Wibowo, S., & Li, L. D., *The effects of over-reliance on AI dia-*  
716 *logue systems on students' cognitive abilities: a systematic review*, Smart  
717 Learning Environments, 11(1), 28, 2024.