

1 International evaluation of an AI system 2 for breast cancer screening

3
4 Scott Mayer McKinney^{1†*}, Marcin T. Sieniek^{1†}, Varun Godbole^{1†}, Jonathan Godwin^{2†}, Natasha
5 Antropova², Hutan Ashrafian⁴, Trevor Back², Mary Chesus², Greg C. Corrado¹, Ara Darzi⁵,
6 Mozziyar Etemadi⁶, Florencia Garcia-Vicente⁶, Fiona J Gilbert⁷, Mark Halling-Brown⁸, Demis
7 Hassabis², Sunny Jansen⁹, Alan Karthikesalingam³, Christopher J Kelly³, Dominic King³, Joseph
8 R. Ledsam², David Melnick⁶, Hormuz Mostofi¹, Bernardino Romera-Paredes², Lily Peng¹,
9 Joshua Jay Reicher¹⁰, Richard Sidebottom^{11,12}, Mustafa Suleyman², Daniel Tse¹, Kenneth C.
10 Young⁸, Jeffrey De Fauw^{2‡}, and Shravya Shetty^{1‡*}

11
12 ¹Google Health, Palo Alto, CA, USA

13 ²DeepMind, London, UK

14 ³Google Health, London, UK

15 ⁴Imperial College London, London, UK

16 ⁵Institute of Global Health Innovation, Imperial College London, London, UK

17 ⁶Northwestern Medicine, Chicago, IL, USA

18 ⁷University of Cambridge School of Clinical Medicine, Cambridge, UK

19 ⁸Royal Surrey County Hospital, Guildford, UK

20 ⁹Verily Life Sciences, South San Francisco, CA, USA

21 ¹⁰Stanford Health Care and Palo Alto Veterans Affairs, Palo Alto, CA, USA

22 ¹¹The Royal Marsden Hospital, London, UK

23 ¹²Thirlestaine Breast Centre, Cheltenham, UK

24

25 †These authors contributed equally.

26 ‡These authors supervised jointly.

27 *e-mail: sshetty@google.com; scottmayer@google.com

28

29 **Screening mammography aims to identify breast cancer before symptoms appear,**
30 **enabling earlier therapy for more treatable disease¹. Despite the existence of screening**
31 **programmes worldwide, interpretation of mammograms suffers from suboptimal rates of**
32 **false positives and false negatives². Here we present an AI system capable of surpassing**
33 **expert readers in breast cancer prediction performance. To assess its performance in the**
34 **clinical setting, we curated a large representative dataset from the United Kingdom (UK)**
35 **and a large enriched dataset from the United States (US). We show an absolute reduction**
36 **of 5.7%/1.2% (US/UK) in false positives and 9.4%/2.7% (US/UK) in false negatives. We**
37 **show evidence of the system's ability to generalise from the UK sites to the US site. In an**
38 **independently-conducted reader study, the AI system out-performed all six radiologists**
39 **with an area under the receiver operating characteristic curve greater than the average**
40 **radiologist by an absolute margin of 11.5%. By simulating the AI system's role in the**
41 **double-reading process, we maintain noninferior performance while reducing the second**
42 **reader's workload by 88%. This robust assessment of the AI system paves the way for**
43 **prospective clinical trials to improve the accuracy and efficiency of breast cancer**
44 **screening.**
45

46 Breast cancer is the second leading cause of cancer death in women³, but early detection and
47 treatment can dramatically improve outcomes^{1,4,5}. As a consequence, many developed nations
48 have implemented large-scale mammography screening programmes. Major medical and
49 governmental organisations recommend screening for all women starting between the ages of
50 40 and 50⁶⁻⁸. In the US and UK combined, over 42 million exams are performed each year^{9,10}.

51 Despite mammography's widespread adoption, interpretation of these images remains
52 challenging. There is high variability in experts' cancer detection accuracy, and the performance
53 of even the best clinicians leaves room for improvement^{11,12}. False positives can lead to patient
54 anxiety¹³, unnecessary follow up, and invasive diagnostic procedures. Cancers missed at
55 screening may not be identified until they are more advanced and less amenable to treatment¹⁴.

56 Artificial intelligence (AI) may be uniquely poised to help. Recent studies have demonstrated
57 AI's ability to meet or exceed the performance of human experts on several medical image
58 analysis tasks¹⁵⁻¹⁹. As a shortage of mammography professionals threatens availability and
59 adequacy of breast screening services around the world²⁰⁻²³, the scalability of AI could improve
60 access to high quality care for all.
61

62 Computer-aided detection (CAD) software for mammography was introduced in the 1990s, and
63 multiple assistive tools have been approved for medical use²⁴. Despite early promise^{25,26}, this
64 generation of software failed to improve reader performance in real-world settings^{12,27,28}. More
65 recently, the field has seen a renaissance owing to the success of deep learning. A few studies
66 have shown breast cancer prediction systems with standalone performance approaching that of
67 human experts^{29,30}. Still, existing work has several limitations. Most studies evaluate on small,
68 enriched datasets with limited follow-up, and few have compared performance to readers in
69 actual clinical practice, instead relying on lab-based simulations of the reading environment. To
70 date, there has been little evidence of the ability of AI systems to translate between different

71 screening populations and settings without additional training data³¹. Critically, the pervasive
72 use of follow-up intervals no longer than 12 months^{29,30,32,33} means that more subtle cancers, not
73 identified until the next screen, may be ignored.

74 In this study, we evaluate the performance of a new AI system for breast cancer prediction
75 using two large, clinically-representative datasets from the UK and US. We compare the
76 system's predictions to those made by readers in routine clinical practice and show performance
77 better than individual radiologists. These observations are confirmed with an independently-
78 conducted reader study. We further show how this system might be integrated into screening
79 workflows, and provide evidence that the system can generalise across continents. Figure 1
80 depicts a high-level overview.

81 **Screening programme datasets**

82 A deep learning model for identifying breast cancer in screening mammograms was developed
83 and evaluated using two large datasets from the UK and the US. We report results on test sets
84 withheld from AI development.

85
86 The UK test set consisted of screening mammograms from 25,856 women collected between
87 2012 and 2015 at two screening centers in England, where women are screened every three
88 years. It included 785 women who had a biopsy, and 414 women with cancer diagnosed within
89 39 months (3 years and 3 months) of imaging. This was a random sample of 10% of all women
90 with screening mammograms at these sites during this time period. The UK cohort resembled
91 the broader screening population in age and disease characteristics (Extended Data Table 1a).

92
93 The test set from the US, where women are screened every 1 or 2 years, consisted of screening
94 mammograms from 3,097 women collected between 2001 and 2018 at one academic medical
95 center. We included images from all 1,511 women biopsied during this time period and a
96 random subset of women who never underwent biopsy (Methods). Among the women who
97 received a biopsy, 686 were diagnosed with cancer within 27 months (2 years and 3 months) of
98 imaging.

99
100 Breast cancer outcome was determined on the basis of multiple years of follow up (Figure 1).
101 We chose the follow-up duration based on the screening interval in each dataset's country of
102 origin. Following previous work³⁴, we augment each interval with a 3-month buffer to account for
103 variability in scheduling and latency of follow up. Cases designated as cancer positive were
104 accompanied by a biopsy-confirmed diagnosis within the follow-up period. Cases labeled as
105 cancer negative had at least one follow-up non-cancer screen; cases without this follow up were
106 excluded from the test set.

107 **Retrospective clinical comparison**

108 We used biopsy-confirmed breast cancer to evaluate predictions of the AI system as well as the
109 original decisions made by radiologists in the course of clinical practice. Human performance
110 was computed based on the clinician's decision to recall the patient for further diagnostic
111 investigation. The receiver operating characteristic (ROC) curve of the AI system's cancer
112 prediction is shown in Figure 2.

113
114 In the UK, each mammogram is interpreted by two readers. In cases of disagreement, an
115 arbitration process is used, invoking a third opinion. These interpretations occur serially such
116 that each reader has access to prior readers' opinions. The records of these decisions yield
117 three human performance benchmarks for cancer prediction.

118
119 Compared to the first reader, the AI system demonstrated a statistically significant absolute
120 specificity improvement of 1.2% (95% CI [0.29%, 2.1%]; $p = 0.0096$ for superiority) and an
121 absolute sensitivity improvement of 2.7% (95% CI [-3%, 8.5%]; $p = 0.004$ for noninferiority at a
122 prespecified 5% margin; Extended Data Table 2a).

123
124 Compared to the second reader, the AI system showed non-inferiority (at a 5% margin) for both
125 specificity ($p < 0.001$) and sensitivity ($p = 0.02$). The AI system likewise showed non-inferiority
126 (at a 5% margin) to the consensus judgment for specificity ($p < 0.001$) and sensitivity ($p =$
127 0.0039).

128
129 In the standard US screening protocol, each mammogram is interpreted by a single radiologist.
130 We used the BI-RADS³⁵ score assigned to each case in the original screening context as a
131 proxy for the human cancer prediction (Methods, Interpreting clinical reads). Compared to the
132 typical reader, the AI system demonstrated statistically significant improvements in absolute
133 specificity of 5.7% (95% CI [2.6%, 8.6%]; $p < 0.001$) and sensitivity of 9.4% (95% CI [4.5%,
134 13.9%]; $p < 0.001$; Extended Data Table 2a).

135 **Generalisation across populations**

136 To evaluate the AI system's ability to generalise across populations and screening settings, we
137 trained the same architecture using only the UK dataset and applied it to the US test set (Figure
138 2b). Even without exposure to US training data, the AI system's ROC curve envelops the point
139 indicating the average performance of US radiologists. Once again, the AI system showed
140 superior specificity (+3.5%, $p = 0.0212$) and superior sensitivity (+8.1%, $p = 0.0006$; Extended
141 Data Table 2b).

142 **Reader study comparison**

143 In a reader study conducted by an external clinical research organisation, six US board-certified
144 radiologists compliant with Mammography Quality Standards Act (MQSA) requirements

145 interpreted 500 mammograms randomly sampled from the US test set. Where data were
146 available, readers were equipped with contextual information typically available in the clinical
147 setting, including patient age, breast cancer history, and prior screening mammograms.

148

149 Among the 500 cases selected for this study, 125 had biopsy-proven cancer within 27 months,
150 125 had a negative biopsy within 27 months, and 250 were not biopsied (Extended Data Table
151 3). These proportions were chosen to increase the difficulty of the screening task and increase
152 statistical power; such enrichment is typical in observer studies³⁶.

153

154 Readers rated each case using the forced BI-RADS³⁵ scale. BI-RADS scores were compared to
155 ground truth outcomes to fit an ROC curve for each reader. The scores of the AI system were
156 treated in the same manner (Extended Data Table 2).

157

158 The AI system exceeded average radiologist performance by a significant margin ($\Delta\text{AUC} =$
159 $+0.115$, 95% CI: [0.055, 0.175], $p = 0.0002$). Similar results were observed when 1 year follow-
160 up was used instead of 27 months (Figure 3c, Extended Data Figure 2).

161

162 In addition to producing case-level classification, the AI system was designed to highlight areas
163 of suspicion for malignancy. Likewise, the readers in our study supplied rectangular region-of-
164 interest (ROI) annotations surrounding concerning findings.

165

166 We used multi-localisation receiver operating characteristic (mLROC) analysis³⁷ to compare the
167 ability of the readers and the AI system to identify malignant lesions within each case (Methods,
168 Localisation analysis).

169

170 We summarised each mLROC plot by computing the partial area under the curve (pAUC) in the
171 false positive fraction interval from 0 to 0.1³⁸ (Extended Data Figure 3). The AI system exceeded
172 human performance by a significant margin ($\Delta\text{pAUC} = +0.0192$, 95% CI: [0.0086, 0.0298], $p =$
173 0.0004).

174 **Potential clinical applications**

175 The AI system's classifications could be used to reduce the workload involved in the UK's
176 double reading process while preserving the standard of care. We explored this scenario
177 through simulation by omitting the second reader and any ensuing arbitration when the AI's
178 decision agreed with the first reader. In these cases, the first reader's opinion was treated as
179 final. In cases of disagreement, the second and consensus opinions were invoked as usual.
180 This combination of human and machine displays performance equivalent to that of the
181 traditional double reading process, while saving 88% of the second reader's effort (Extended
182 Data Table 4a).

183

184 The AI system could also be used to provide automated, immediate feedback in the screening
185 setting.

186
187 In order to identify normal cases with high confidence, we used a very low decision threshold.
188 On the UK data, we achieved a negative predictive value (NPV) of 99.99% while retaining a
189 specificity of 41.15%. Similarly, on the US data, we achieved a NPV of 99.90% while retaining a
190 specificity of 34.79%. These data suggest that it may be feasible to pick out 35–41% of normal
191 cases if we allow for one cancer in every 1,000–10,000 negative predictions (NPV 99.90–
192 99.99% in US–UK). For comparison, consensus double reading in our UK dataset included one
193 cancer in every 182 cases deemed normal.

194
195 To identify cancer cases with high confidence, we used a very high decision threshold. On the
196 UK data, we achieved a positive predictive value (PPV) of 85.6% while retaining a sensitivity of
197 41.2%. Likewise, on the US data, we achieved a PPV of 82.4% while retaining a sensitivity of
198 29.8%. These data suggest that it may be feasible to rapidly prioritise 30–40% of cancer cases
199 with approximately 5 of 6 follow ups leading to cancer diagnosis. By comparison, in our study
200 only 22.8% of UK cases recalled by consensus double reading and 4.9% of US cases recalled
201 by single reading were ultimately diagnosed with cancer.

202 **Performance breakdown**

203 Comparing the errors of the AI system with errors from clinical reads revealed many cases in
204 which the AI system correctly identified cancer while the reader did not and vice versa
205 (Supplementary Table 1). Most of the cases in which only the AI system identified cancer were
206 invasive (Extended Data Table 5). On the other hand, cases in which only the reader identified
207 cancer were split more evenly between in situ and invasive. Further breakdowns by invasive
208 cancer size, grade, and molecular markers show no clear biases (Supplementary Table 2).

209
210 We also considered the disagreement between the AI system and the six radiologists that
211 participated in the US reader study. Figure 4a shows a sample cancer case missed by all six
212 radiologists, but correctly identified by the AI system. Figure 4b shows a sample cancer case
213 caught by all six radiologists but missed by the AI system. While we were unable to determine
214 clear patterns among these instances, the presence of such edge cases suggests potentially
215 complementary roles for the AI system and human readers in reaching accurate conclusions.

216
217 We compared the performance of the 20 individual readers best represented in the UK clinical
218 dataset with that of the AI system (Extended Data Table 7). The results of this analysis suggest
219 that the aggregate comparison presented above is not unduly influenced by any particular
220 readers. Breakdowns by cancer type, grade, and lesion size suggest no apparent difference in
221 the distribution of cancers detected by the AI system and human readers (Extended Data Table
222 6a).

223
224 On the US test set, a breakdown by cancer type (Extended Data Table 6b) shows that the AI
225 system's sensitivity advantage is concentrated on the identification of invasive cancers (e.g.
226 invasive lobular/ductal carcinoma) rather than in situ cancer (e.g. ductal carcinoma in situ). A

227 breakdown by BI-RADS³⁵ breast density category shows that performance gains apply equally
228 across the spectrum of breast tissue types represented in this data set (Extended Data Table
229 6c).

230 **Discussion**

231 In this study we present an AI system that outperforms radiologists on a clinically relevant
232 breast cancer identification task. These results held on two large datasets representative of
233 different country-specific screening populations and practices.

234 In the UK, the AI system showed specificity superior to that of the first reader. Sensitivity at the
235 same operating point was noninferior. Consensus double reading has been shown to improve
236 performance compared to single reading³⁹, and represents the current standard of care in the
237 UK and many European countries⁴⁰. Our system did not outperform this benchmark, but was
238 statistically noninferior to the second reader and consensus opinion.

239 In the US, the AI system displayed specificity and sensitivity superior to that of radiologists
240 practicing in an academic medical center. This trend was confirmed in an externally conducted
241 reader study, which showed that the scores of the AI system stratify cases better than each of
242 the six readers' BI-RADS ratings, the standard scale for mammography assessment in the US.

243 Remarkably, the human readers (both in the clinic and our reader study) had access to patient
244 history and prior mammograms when making screening decisions. The US clinical readers may
245 have also had access to breast tomosynthesis images. In contrast, the AI system only
246 processed the most recent mammogram.

247 These comparisons are not without limitations. While the UK dataset mirrored the nationwide
248 screening population in age and cancer prevalence (Extended Data Table 1a), the same cannot
249 be said of the US data, which was drawn from a single screening centre and was enriched for
250 cancer cases.

251 By chance, the vast majority of images used in this study were acquired on devices made by
252 Hologic, Inc. Future research should assess the AI system's performance across a variety of
253 manufacturers in a more systematic way.

254 In our reader study, all the radiologists were eligible to interpret screening mammograms in the
255 US, but did not uniformly receive fellowship training in breast imaging. It is possible that a higher
256 performance benchmark could have been obtained with more specialised readers⁴¹.

257 To obtain high-quality ground-truth labels, we employed extended follow-up intervals chosen to
258 encompass a subsequent screening round in each country. Although there is some precedent in
259 clinical trials³⁴ and targeted cohort studies⁴², this step is not usually taken when undertaking
260 systematic evaluation of AI systems for breast cancer detection.

261 In retrospective datasets with shorter follow-up intervals, outcome labels tend to be skewed in
262 favour of readers. Since they are gatekeepers for biopsy, asymptomatic cases will only receive
263 a cancer diagnosis if a mammogram raised reader suspicion. A longer follow-up interval
264 decouples the ground truth labels from reader opinions (Extended Data Figure 4) and includes
265 cancers that may have been initially missed by human eyes.

266 The use of an extended interval makes cancer prediction a more challenging task. Cancers
267 diagnosed years later may include new growths for which there could be no mammographic
268 evidence in the original images. Consequently, the sensitivity values presented here are lower
269 than what has been reported for 12 month intervals² (Extended Data Figure 5).

270 We present early evidence of the AI system's ability to generalise across populations and
271 screening protocols. We retrained the system using exclusively UK data, and then measured
272 performance on unseen US data. In this context, the system continued to outperform
273 radiologists, albeit by a smaller margin. This suggests that in future clinical deployments, the
274 system might offer strong baseline performance, but may benefit from fine-tuning with local
275 data.

276 The utility of the AI system within clinical workflows remains to be determined. The specificity
277 advantage exhibited by the AI system suggests it could help reduce recall rates and
278 unnecessary biopsies. The improvement in sensitivity, exhibited in the US data, shows that the
279 AI system may be capable of detecting cancers earlier than the standard of care. An analysis of
280 the AI system's localisation performance suggests the early promise of using it to flag
281 suspicious regions for review by experts. Notably, the additional cancers identified tended to be
282 invasive rather than in situ disease.

283 Beyond augmenting reader performance, the technology described here may have a number of
284 other clinical applications. Through simulation, we suggest how the system could obviate the
285 need for double reading in 88% of UK screening cases, while maintaining similar accuracy to
286 the standard protocol. We also explore how high-confidence operating points can be used to
287 triage high-risk cases and dismiss low-risk cases. These analyses highlight the potential of this
288 technology to deliver screening results in a sustainable manner despite workforce challenges in
289 countries like the UK⁴³. Prospective clinical studies will be required to understand the full extent
290 to which this technology can benefit patient care.

291 **References**

- 292 1. Tabár, L. *et al.* Swedish two-county trial: impact of mammographic screening on breast
293 cancer mortality during 3 decades. *Radiology* **260**, 658–663 (2011).
- 294 2. Lehman, C. D. *et al.* National Performance Benchmarks for Modern Screening Digital
295 Mammography: Update from the Breast Cancer Surveillance Consortium. *Radiology* **283**,
296 49–58 (2017).

- 297 3. Bray, F. *et al.* Global cancer statistics 2018: GLOBOCAN estimates of incidence and
298 mortality worldwide for 36 cancers in 185 countries. *CA Cancer J. Clin.* **68**, 394–424 (2018).
- 299 4. Tonelli, M. *et al.* Recommendations on screening for breast cancer in average-risk women
300 aged 40-74 years. *CMAJ* **183**, 1991–2001 (2011).
- 301 5. Marmot, M. G. *et al.* The benefits and harms of breast cancer screening: an independent
302 review. *Br. J. Cancer* **108**, 2205–2240 (2013).
- 303 6. Lee, C. H. *et al.* Breast cancer screening with imaging: recommendations from the Society
304 of Breast Imaging and the ACR on the use of mammography, breast MRI, breast
305 ultrasound, and other technologies for the detection of clinically occult breast cancer. *J. Am.*
306 *Coll. Radiol.* **7**, 18–27 (2010).
- 307 7. Oeffinger, K. C. *et al.* Breast Cancer Screening for Women at Average Risk: 2015
308 Guideline Update From the American Cancer Society. *JAMA* **314**, 1599–1614 (2015).
- 309 8. Siu, A. L. & U.S. Preventive Services Task Force. Screening for Breast Cancer: U.S.
310 Preventive Services Task Force Recommendation Statement. *Ann. Intern. Med.* **164**, 279–
311 296 (2016).
- 312 9. Center for Devices & Radiological Health. MQSA National Statistics. *U.S. Food and Drug*
313 *Administration* (2019). Available at: [http://www.fda.gov/radiation-emitting-products/mqsa-](http://www.fda.gov/radiation-emitting-products/mqsa-insights/mqsa-national-statistics)
314 [insights/mqsa-national-statistics](http://www.fda.gov/radiation-emitting-products/mqsa-insights/mqsa-national-statistics). (Accessed: 16th July 2019)
- 315 10. Breast screening. *Cancer Research UK* (2017). Available at:
316 <https://www.cancerresearchuk.org/about-cancer/breast-cancer/screening/breast-screening>.
317 (Accessed: 26th July 2019)
- 318 11. Elmore, J. G. *et al.* Variability in interpretive performance at screening mammography and
319 radiologists' characteristics associated with accuracy. *Radiology* **253**, 641–651 (2009).
- 320 12. Lehman, C. D. *et al.* Diagnostic Accuracy of Digital Screening Mammography With and
321 Without Computer-Aided Detection. *JAMA Intern. Med.* **175**, 1828–1837 (2015).

- 322 13. Tosteson, A. N. A. *et al.* Consequences of False-Positive Screening Mammograms. *JAMA*
323 *Internal Medicine* **174**, 954 (2014).
- 324 14. Houssami, N. & Hunter, K. The epidemiology, radiology and biological characteristics of
325 interval breast cancers in population mammography screening. *NPJ Breast Cancer* **3**, 12
326 (2017).
- 327 15. Gulshan, V. *et al.* Development and Validation of a Deep Learning Algorithm for Detection
328 of Diabetic Retinopathy in Retinal Fundus Photographs. *JAMA* **316**, 2402–2410 (2016).
- 329 16. Esteva, A. *et al.* Dermatologist-level classification of skin cancer with deep neural networks.
330 *Nature* **542**, 115–118 (2017).
- 331 17. De Fauw, J. *et al.* Clinically applicable deep learning for diagnosis and referral in retinal
332 disease. *Nat. Med.* **24**, 1342–1350 (2018).
- 333 18. Ardila, D. *et al.* End-to-end lung cancer screening with three-dimensional deep learning on
334 low-dose chest computed tomography. *Nat. Med.* **25**, 954–961 (2019).
- 335 19. Topol, E. J. High-performance medicine: the convergence of human and artificial
336 intelligence. *Nat. Med.* **25**, 44–56 (2019).
- 337 20. Moran, S. & Warren-Forward, H. The Australian BreastScreen workforce: a snapshot.
338 *Radiographer* **59**, 26–30 (2012).
- 339 21. Wing, P. & Langelier, M. H. Workforce shortages in breast imaging: impact on
340 mammography utilization. *AJR Am. J. Roentgenol.* **192**, 370–378 (2009).
- 341 22. Rimmer, A. Radiologist shortage leaves patient care at risk, warns royal college. *BMJ* **359**,
342 j4683 (2017).
- 343 23. Nakajima, Y., Yamada, K., Imamura, K. & Kobayashi, K. Radiologist supply and workload:
344 international comparison. *Radiation Medicine* **26**, 455–465 (2008).
- 345 24. Rao, V. M. *et al.* How widely is computer-aided detection used in screening and diagnostic
346 mammography? *J. Am. Coll. Radiol.* **7**, 802–805 (2010).

- 347 25. Gilbert, F. J. *et al.* Single reading with computer-aided detection for screening
348 mammography. *N. Engl. J. Med.* **359**, 1675–1684 (2008).
- 349 26. Giger, M. L., Chan, H.-P. & Boone, J. Anniversary Paper: History and status of CAD and
350 quantitative image analysis: The role of Medical Physics and AAPM. *Medical Physics* **35**,
351 5799–5820 (2008).
- 352 27. Fenton, J. J. *et al.* Influence of computer-aided detection on performance of screening
353 mammography. *N. Engl. J. Med.* **356**, 1399–1409 (2007).
- 354 28. Kohli, A. & Jha, S. Why CAD Failed in Mammography. *J. Am. Coll. Radiol.* **15**, 535–537
355 (2018).
- 356 29. Rodriguez-Ruiz, A. *et al.* Stand-Alone Artificial Intelligence for Breast Cancer Detection in
357 Mammography: Comparison With 101 Radiologists. *J. Natl. Cancer Inst.* (2019).
358 doi:10.1093/jnci/djy222
- 359 30. Wu, N. *et al.* Deep Neural Networks Improve Radiologists' Performance in Breast Cancer
360 Screening. (2019).
- 361 31. Zech, J. R. *et al.* Variable generalization performance of a deep learning model to detect
362 pneumonia in chest radiographs: A cross-sectional study. *PLoS Med.* **15**, e1002683 (2018).
- 363 32. Becker, A. S. *et al.* Deep Learning in Mammography: Diagnostic Accuracy of a
364 Multipurpose Image Analysis Software in the Detection of Breast Cancer. *Invest. Radiol.*
365 **52**, 434–440 (2017).
- 366 33. Ribli, D., Horváth, A., Unger, Z., Pollner, P. & Csabai, I. Detecting and classifying lesions in
367 mammograms with Deep Learning. *Sci. Rep.* **8**, 4165 (2018).
- 368 34. Pisano, E. D. *et al.* Diagnostic performance of digital versus film mammography for breast-
369 cancer screening. *N. Engl. J. Med.* **353**, 1773–1783 (2005).
- 370 35. American College of Radiology. *ACR BI-RADS Atlas: Breast Imaging Reporting and Data*
371 *System ; Mammography, Ultrasound, Magnetic Resonance Imaging, Follow-up and*

372 *Outcome Monitoring, Data Dictionary*. (2013).

373 36. Gallas, B. D. *et al.* Evaluating imaging and computer-aided detection and diagnosis devices
374 at the FDA. *Acad. Radiol.* **19**, 463–477 (2012).

375 37. Swensson, R. G. Unified measurement of observer performance in detecting and localizing
376 target objects on images. *Med. Phys.* **23**, 1709–1725 (1996).

377 38. Samulski, M. *et al.* Using computer-aided detection in mammography as a decision
378 support. *Eur. Radiol.* **20**, 2323–2330 (2010).

379 39. Brown, J., Bryan, S. & Warren, R. Mammography screening: an incremental cost
380 effectiveness analysis of double versus single reading of mammograms. *BMJ* **312**, 809–
381 812 (1996).

382 40. Giordano, L. *et al.* Mammographic screening programmes in Europe: organization,
383 coverage and participation. *J. Med. Screen.* **19 Suppl 1**, 72–82 (2012).

384 41. Sickles, E. A., Wolverton, D. E. & Dee, K. E. Performance parameters for screening and
385 diagnostic mammography: specialist and general radiologists. *Radiology* **224**, 861–869
386 (2002).

387 42. Ikeda, D. M., Birdwell, R. L., O’Shaughnessy, K. F., Sickles, E. A. & Brenner, R. J.
388 Computer-aided detection output on 172 subtle findings on normal mammograms
389 previously obtained in women with breast cancer detected at follow-up screening
390 mammography. *Radiology* **230**, 811–819 (2004).

391 43. The breast imaging and diagnostic workforce in the United Kingdom. *The Royal College of*
392 *Radiologists* (2016). Available at: [https://www.rcr.ac.uk/publication/breast-imaging-and-](https://www.rcr.ac.uk/publication/breast-imaging-and-diagnostic-workforce-united-kingdom)
393 [diagnostic-workforce-united-kingdom](https://www.rcr.ac.uk/publication/breast-imaging-and-diagnostic-workforce-united-kingdom). (Accessed: 22nd July 2019)

394
395
396

397 **Figures**

398

399 **Figure 1. Development of an AI system to detect cancer in screening mammograms.**

400 Datasets representative of the UK and US breast cancer screening populations were curated
401 from three screening centers in the UK and one center in the US. Outcomes were derived from
402 the biopsy record and longitudinal follow up. An AI system was trained to identify the presence
403 of breast cancer from a screening mammogram; it was evaluated in three primary ways. AI
404 predictions were compared with the historical decisions made in clinical practice. To evaluate
405 the generalisability across populations, a version of the AI system was developed using only the
406 UK data and retested on the US data. Finally, the AI system was compared with six
407 independent radiologists using a subset of the US test set.

408

409 **Figure 2. Breast cancer prediction performance.**

410 **a.** Receiver operating characteristic (ROC) curve of the AI system on the UK screening data.

411 The area under the curve (AUC) is 0.889 (95% CI [0.871, 0.907]; n = 25,856 patients). Also

412 shown are the (sensitivity, specificity) pairs of the human decisions made in clinical practice.

413 Cases were considered positive if they received a biopsy-confirmed cancer diagnosis within 39

414 months (3 years and 3 months) from the time of screening. The consensus decision represents

415 the standard of care in the UK, and will involve input from between 2 and 3 expert readers. The

416 inset shows an enhancement of the gray shaded region. AI system operating points were

417 selected on a separate validation dataset: point (i) was intended to match the sensitivity and

418 exceed the specificity of the first reader; points (ii) and (iii) were selected to attain non-inferiority

419 for both the sensitivity and specificity of the second reader and consensus opinion, respectively.

420 **b.** ROC curve of the AI system on the US screening data. When trained on both datasets, the

421 AUC is 0.8107 (95% CI [0.791, 0.831]; n = 3,097 patients). When trained only on the UK dataset

422 (dotted curve), the AUC is 0.757 (95% CI [0.732, 780]). Also shown are the sensitivity and

423 specificity achieved by radiologists in clinical practice using BI-RADS³⁵. Cases were considered

424 positive if they received a biopsy-confirmed cancer diagnosis within 27 months (2 years and 3

425 months) from the time of screening. AI system operating points were chosen to exceed the

426 average reader's sensitivity and specificity. Negative cases were upweighted to account for the

427 sampling protocol (Methods, Inverse probability weighting). Extended Data Figure 1 shows an

428 unweighted analysis. See Extended Data Table 2a for statistical comparisons of sensitivity and

429 specificity.

430

431 **Figure 3. Breast cancer prediction performance compared to six independent readers.**

432 **a.** Six readers rated each case (n = 465) using the 6-point BI-RADS scale. A fitted ROC curve

433 for each of the readers is compared to the ROC curve of the AI system (Methods, Statistical

434 analysis). For reference, a nonparametric ROC curve is presented in tandem. Cases were

435 considered positive (n = 113) if they received a pathology-confirmed cancer diagnosis within 27

436 months (2 years and 3 months) from the time of screening. Note that this sample of cases was

437 enriched for patients that had received a negative biopsy result (n = 119), making this a more

438 challenging population for screening. The mean reader AUC was 0.625 (s.d. 0.032), while the AI

439 system's AUC was 0.740 (95% CI: [0.696, 0.794]). The AI system exceeded human

440 performance by a significant margin ($\Delta = +0.115$, 95% CI: [0.055, 0.175], $p = 0.0002$, two-sided
441 ORH method). For results using a 12-month interval, see Extended Data Figure 2.

442 **b.** Pooled results from all six readers from (a).

443 **c.** Pooled results ($n = 408$) from all six readers using a 12-month interval for cancer definition.
444 Cases were considered positive ($n = 56$) if they received a pathology-confirmed cancer
445 diagnosis within 1 year (Extended Data Table 3).

446

447 **Figure 4. Discrepancies between the AI system and human readers.**

448 **a.** A sample cancer case missed by all six readers in the US reader study, but correctly
449 identified by the AI system. The images show two views of a small, irregular mass with
450 associated microcalcifications in the lower inner right breast.

451 **b.** A sample cancer case caught by all six readers in the US reader study, but missed by the AI
452 system. The images show two views of a dense mass in the lower inner right breast.
453 (left, mediolateral oblique; right, craniocaudal)

454

455

456 **Methods**

457 **Ethical approval.** Use of the UK dataset for research collaborations by both commercial and
458 non-commercial organisations received ethical approval (REC reference 14/SC/0258). The US
459 data was fully de-identified and released only after an Institutional Review Board approval
460 (STU00206925).

461
462 **UK dataset.** The UK dataset was collected from three breast screening sites in the United
463 Kingdom National Health Service Breast Screening Programme (NHSBSP). The NHSBSP
464 invites women aged between 50 and 70 who are registered with a general practitioner (GP) for
465 mammographic screening every 3 years. Women who are not registered with a GP, or who are
466 older than 70, can self-refer to the screening programme. In the UK, the screening programme
467 uses double reading: each mammogram is read by two radiologists, who are asked to decide
468 whether to recall the woman for additional followup. When there is disagreement, an arbitration
469 process takes place.

470
471 The data was initially compiled by OPTIMAM, a Cancer Research UK effort, between the years
472 of 2010 and 2018 from St. George's Hospital (London, UK), Jarvis Breast Centre (Guildford,
473 UK) and Addenbrooke's Hospital (Cambridge, UK). The collected data included screening and
474 follow-up mammograms (comprising mediolateral oblique "MLO" and craniocaudal "CC" views
475 of the left and right breast), all radiologist opinions (including the arbitration result, if applicable)
476 and metadata associated with follow-up treatment.

477
478 The mammograms and associated metadata of 137,291 women were considered for inclusion
479 in the study. Of these, 123,964 had both screening images and uncorrupted metadata. Exams
480 that were recalled for reasons other than radiographic evidence of malignancy, or episodes that
481 were not part of routine screening were excluded. In total, 121,850 women had at least one
482 eligible exam. Women who were aged below 47 at the time of the screen were excluded from
483 validation and test sets, leaving 121,455 women. Finally, women for whom there was insufficient
484 follow up for any scan were excluded from validation and test. This last step resulted in the
485 exclusion of 5,990 of 31,766 test set cases (19%). See Supplementary Figure 1.

486
487 The test set is a random sample of 10% of all women screened at two sites, St. George's and
488 Jarvis, between the years 2012 and 2015. Insufficient data was provided to apply the sampling
489 procedure to the third site. In assembling the test set, we randomly selected a single eligible
490 screening mammogram from each woman's record. For women with a positive biopsy, eligible
491 mammograms were those conducted in the 39 months (3 years and 3 months) prior to the
492 biopsy date. For women that never had a positive biopsy, eligible mammograms were those
493 with a non-suspicious mammogram at least 21 months later.

494
495 The final test set consisted of 25,856 women (see Supplementary Figure 1). When compared to
496 the UK national breast cancer screening service we see a very similar cancer prevalence, age
497 and cancer type distribution (see Extended Data Table 1a). Digital mammograms were acquired

498 predominantly on devices manufactured by Hologic, Inc. (95%), followed by General Electric
499 (4%) and Siemens (1%).

500

501 **US dataset.** The US dataset was collected from Northwestern Memorial Hospital (Chicago, IL)
502 between the years of 2001 and 2018. In the US, each screening mammogram is typically read
503 by a single radiologist, and screens are conducted annually or biannually. The breast
504 radiologists at this hospital are fellowship-trained and only interpret breast imaging studies.
505 Their experience levels ranged from 1-30 years. The American College of Radiology (ACR)
506 recommends that women start routine screening at the age of 40, while other organizations
507 including the US Preventive Services Task Force (USPSTF) recommend initiation at 50 for
508 women with average breast cancer risk⁶⁻⁸.

509

510 The US dataset included records from all women that underwent a breast biopsy between 2001
511 and 2018. It also included a random sample of approximately 5% of all women who participated
512 in screening, but were never biopsied. This heuristic was employed in order to capture all
513 cancer cases (to enhance statistical power) and to curate a rich set of benign findings on which
514 to train and test the AI system.

515

516 Supplementary Figure 2 distills the data processing steps involved in constructing the dataset.
517 Among women with a completed mammogram order, we collected the records from all women
518 with a pathology report containing the term “breast”. Among those that lacked such a pathology
519 report, women whose records bore an International Classification of Diseases (ICD) code
520 indicative of breast cancer were excluded. Approximately 5% of this unbiopsied negative
521 population was sampled. After deidentification and transfer, women were excluded if their
522 metadata was either unavailable or corrupted. The women in the dataset were split randomly
523 among train (55%), validation (15%) and test (30%). For testing, a single case was chosen for
524 each woman following a similar procedure as in the UK dataset. In women who underwent
525 biopsy, we randomly chose a case from the 27 months preceding the date of biopsy. For
526 women who did not undergo biopsy, one screening mammogram was randomly chosen from
527 among those with a follow up event at least 21 months later.

528

529 Cases were considered complete if they possessed the four standard screening views
530 (mediolateral oblique “MLO” and craniocaudal “CC” views of the left and right breast) acquired
531 for screening intent. Here too, the vast majority of the studies were acquired using Hologic
532 (including Lorad-branded) devices (99%) while manufacturers Siemens and General Electric
533 together constituted less than 1% of studies.

534

535 The radiology reports associated with cases in the test set were used to flag and exclude cases
536 in the test set which depicted breast implants or were recalled for technical reasons. To
537 compare the AI system against the clinical reads performed at this site, we employed clinicians
538 to manually extract BI-RADS scores from the original radiology reports. There were some cases
539 for which the original radiology report could not be located, even if a subsequent cancer
540 diagnosis was biopsy-confirmed. This might have happened, for example, if the screening case

541 was imported from an outside institution. Such cases were excluded from the clinical reader
542 comparison.

543

544 **Inverse probability weighting.** The US test set includes images from all biopsied women, but
545 only a random subset of women who never underwent biopsy. This enrichment allowed us to
546 accrue more positives in light of the low baseline prevalence of breast cancer, but led to
547 underrepresentation of normal cases. We accounted for this sampling process by using inverse
548 probability weighting to obtain unbiased estimates of human and AI system performance in the
549 natural screening population^{44,45}.

550

551 We acquired images from 7,522 of the 143,238 women who underwent mammography
552 screening but had no cancer diagnosis or biopsy record. Accordingly, we upweighted cases
553 from women who never underwent biopsy by a factor of 19.04. Further sampling occurred when
554 selecting one case per patient: to enrich for difficult cases, we preferentially chose cases from
555 the timeframe preceding a biopsy, if one occurred. Although this sampling increases the
556 diversity of benign findings, it again shifts the distribution from what would be observed in a
557 typical screening interval. To better reflect the prevalence resulting when negative cases are
558 randomly selected, we estimated additional factors by Monte Carlo simulation. When choosing
559 one case per patient with our preferential sampling mechanism, we got 872 cases that were
560 biopsied within 27 months, and 1,662 cases that were not (Supplementary Figure 2). However,
561 100 trials of pure random sampling yielded on average 557.54 and 2,056.46 cases,
562 respectively. Accordingly, cases associated with negative biopsies were down-weighted by
563 $557.54 / 872 = 0.64$. Cases that were not biopsied were up-weighted by another $2,056.46 /$
564 $1,662 = 1.24$, leading to a final weight of $19.04 \times 1.24 = 23.61$. Cancer positive cases carried a
565 weight of 1.0. The final sample weights were used in sensitivity, specificity and ROC
566 calculations.

567

568 **Histopathological outcomes.** In the UK dataset, benign and malignant classifications, given
569 directly in the metadata, followed NHSBSP definitions⁴⁶. To derive the outcomes labels for the
570 US dataset, pathology reports were reviewed by US board-certified pathologists and
571 categorized according to the findings they contained. An effort was made to make this
572 categorization consistent with UK definitions. Malignant pathologies included ductal carcinoma
573 in situ, microinvasive carcinoma, invasive ductal carcinoma, invasive lobular carcinoma, special
574 type invasive carcinoma (including tubular, mucinous, and cribriform carcinomas), intraductal
575 papillary carcinoma, non-primary breast cancers (including lymphoma and phyllodes), and
576 inflammatory carcinoma. Any woman who received a biopsy resulting in any of these malignant
577 pathologies was considered to have a diagnosis of cancer.

578

579 Benign pathologies included lobular carcinoma in situ, radial scar, columnar cell changes,
580 atypical lobular hyperplasia, atypical ductal hyperplasia, cyst, sclerosing adenosis,
581 fibroadenoma, papilloma, periductal mastitis, and usual ductal hyperplasia. None of these
582 findings qualified a woman for a cancer diagnosis.

583

584 **Interpreting clinical reads.** In the UK screening setting, readers categorise mammograms from
585 asymptomatic women as normal or abnormal, with a third option for technical recall due to
586 inadequate image quality. An abnormal result at the conclusion of the double reading process
587 results in further diagnostic workup. We treat mammograms deemed abnormal as a prediction
588 of malignancy. Cases in which the consensus judgment recalled the patient for technical
589 reasons were excluded from analysis, as the images were presumed incomplete or unreliable.
590 Cases in which any single reader recommended technical recall were excluded from the
591 corresponding reader comparison.

592
593 In the US screening setting, radiologists attach a BI-RADS³⁵ score to each mammogram. A
594 score of 0 is deemed "incomplete", and will be later refined based on follow up imaging or
595 repeat mammography to address technical issues. For computation of sensitivity and specificity,
596 we dichotomized the BI-RADS assessments in line with previous work³⁴. Scores of 0, 4 and 5
597 were treated as positive predictions if recall was not based on technical grounds and the
598 recommendation was based on mammographic findings, not solely patient symptoms. Cases of
599 technical recall were excluded from analysis, as the images were presumed incomplete or
600 unreliable. BI-RADS scores were manually extracted from the free-text radiology reports. Cases
601 for which the BI-RADS score was unavailable were excluded from the reader comparison.

602
603 In both datasets, the original readers had access to contextual information normally available in
604 clinical practice. This includes the patient's family history of cancer, prior screening and
605 diagnostic imaging, and radiology or pathology notes from past examinations. In contrast, only
606 the patient's age was made available to the AI system.

607
608 **Overview of the AI system.** The AI system consisted of an ensemble of three deep learning
609 models, each operating on a different level of analysis (individual lesions, individual breasts, and
610 the full case). Each model produces a cancer risk score between 0 and 1 for the entire
611 mammography case. The final prediction of the system was the mean of the predictions from
612 the three independent models. A detailed description of the AI system is available in
613 Supplementary Methods and Supplementary Figure 3.

614
615 **Operating point selection.** The AI system natively produces a continuous score representing
616 the likelihood that cancer is present. To support comparisons with the predictions of human
617 readers, we thresholded this score to produce analogous binary screening decisions. For each
618 clinical benchmark, we used the validation set to choose a distinct operating point; this amounts
619 to a score threshold separating positive and negative decisions. To better simulate prospective
620 deployment, the test sets were never used in selecting operating points.

621
622 The UK dataset contains three clinical benchmarks--the first reader, second reader, and
623 consensus. This last decision is the outcome of the double reading process and represents the
624 standard of care in the UK. For the first reader, we chose an operating point aimed at
625 demonstrating statistical superiority in specificity and non-inferiority for sensitivity. For the
626 second reader and consensus reader, we chose an operating point aimed at demonstrating
627 statistical non-inferiority to the human reader for both sensitivity and specificity.

628
629 The US dataset contains a single operating point for comparison, corresponding to the single
630 radiologist using the BI-RADS rubric for evaluation. In this case, we used the validation set to
631 choose an operating point aimed at achieving superiority on both sensitivity and specificity.
632

633 **Reader study.** For the reader study, 6 US board-certified radiologists interpreted a sample of
634 500 cases from 500 women in the test set. All radiologists were compliant with MQSA
635 requirements for interpreting mammography and had an average of 10 years of clinical
636 experience (Extended Data Table 7). Two of them were fellowship-trained in breast imaging.
637 The sample of cases was stratified to contain 50% normal cases, 25% biopsy negative cases
638 and 25% of biopsy positive cases. A detailed description of the reader study case composition
639 can be found in Extended Data Table 3. Readers were not informed of the enrichment levels in
640 the dataset.

641
642 Readers recorded their assessments on a 21CFR11-compliant electronic case report form
643 within the Ambra Health (New York, NY) viewer v3.18.7.0R. They interpreted the images using
644 5MP MSQA-compliant displays. Each reader interpreted the cases in a unique randomized
645 order.

646
647 For each study, readers were asked to first report a BI-RADS³⁵ 5th edition score among 0, 1,
648 and 2, as if they were interpreting the screening mammogram in routine practice. They were
649 then asked to render a forced diagnostic BI-RADS score using one of the following values: 1, 2,
650 3, 4A, 4B, 4C or 5. Readers also gave a finer-grained score between 0 and 100 indicating their
651 suspicion that the case contains a malignancy.

652
653 In addition to the 4 standard mammographic screening images, clinical context was provided to
654 better simulate the screening setting. Readers were presented with the preamble of the
655 deidentified radiology report produced by the radiologist originally interpreting the study. This
656 contained information such as the patient's age and family history of cancer. The information
657 was manually reviewed to ensure that no impression or findings were included.

658
659 Where possible (in 43% of cases), prior imaging was made available to the readers. Readers
660 could review up to four sets of prior screening exams, acquired between 1 and 4 years earlier,
661 accompanied by deidentified radiologist reports. If prior imaging was available, the study was
662 read twice by each reader--first without the prior information and immediately after, with prior
663 information present. The system ensured that readers could not update their initial assessment
664 after the prior information was presented. For cases where prior exams were available, reader
665 assessment after having reviewed priors was used for the analysis.

666
667 Cases for which at least half of the readers indicated image quality concerns were excluded
668 from analysis. Cases in which breast implants were noted were excluded as well. The final
669 analysis was performed on the remaining 465 cases.
670

671 **Localisation analysis.** For this purpose, we considered all screening exams from the reader
672 study for which cancer developed within 12 months. See Extended Data Table 3 for a detailed
673 description of how the dataset was constructed. To collect ground truth localisations, two board-
674 certified radiologists inspected each case, using follow-up data to identify the location of
675 malignant lesions. Instances of disagreement were resolved by one radiologist with fellowship
676 training in breast imaging. To identify the precise location of the cancerous tissue, radiologists
677 consulted subsequent diagnostic mammograms, radiology reports, biopsy notes, pathology
678 reports, and post-biopsy mammograms. Rectangular bounding boxes were drawn around the
679 locations of subsequent positive biopsies in all views in which the finding was visible. In cases
680 where no mammographic finding was visible, the location where the lesion later appeared was
681 highlighted. Of the 56 cancers considered for analysis, location information could be obtained
682 with confidence in 53 cases. Three cases were excluded due to ambiguity in the index
683 examination and the absence of follow-up images. On average, there were 2.018 ground truth
684 regions per cancer-positive case.

685
686 In the reader study, readers supplied rectangular region-of-interest (ROI) annotations
687 surrounding suspicious findings in all cases they rated BI-RADS 3 or higher. A limit of 6 ROIs
688 per case was enforced. On average, the readers supplied 2.04 annotations per suspicious case.
689 In addition to an overall cancer likelihood score, the AI system emits a ranked list of rectangular
690 bounding boxes for each case. To conduct a fair comparison, we allowed the AI system only its
691 top two bounding boxes to match the number of ROIs produced by the readers.

692
693 To compare the localisation performance of the AI system with that of the readers, we used a
694 method inspired by location receiver operating characteristic (LROC) analysis³⁷. LROC analysis
695 differs from traditional ROC analysis in that the ordinate is a sensitivity measure that factors in
696 localisation accuracy. Although LROC analysis traditionally involves a single finding per
697 case^{37,47}, we permitted multiple unranked findings to match the format of our data. We use the
698 term multi-localization ROC analysis (mLROC) to describe our approach. For each threshold, a
699 cancer case was considered a true positive if its casewise score exceeded this threshold and at
700 least one culprit area was correctly localised in any of the four mammogram views. Correct
701 localisation required an intersection-over-union (IoU) of 0.1 with the ground truth ROI. False
702 positives were defined as usual.

703
704 CAD systems are often evaluated on the basis of whether the center of their marking falls within
705 the boundary of a ground truth annotation⁴⁸. This is potentially problematic since it doesn't
706 properly penalize predicted bounding boxes that are so large as to be nonspecific, but whose
707 center nevertheless happens to fall within the target region. Similarly, large ground truth
708 annotations associated with diffuse findings might be overly generous to the CAD system. We
709 prefer the IoU metric because it balances these considerations. We chose a threshold of 0.1 to
710 account for the fact that indistinct margins on mammography findings lead to region-of-interest
711 annotations of vastly different sizes depending on subjective factors of the annotator. See
712 Supplementary Figure 4. Similar work in 3D chest computed tomography¹⁸ used *any* pixel
713 overlap to qualify for correct localisation. Likewise, an FDA-approved software device for wrist

714 fracture detection reports statistics in which true positives require at least one pixel of overlap⁴⁹.
715 An IoU value of 0.1 is strict by these standards.

716

717 **Statistical analysis.** To evaluate standalone AI system performance, the area under the ROC
718 curve was estimated using the normalized Wilcoxon (Mann-Whitney) U statistic⁵⁰. This is the
719 standard nonparametric method employed by most modern software libraries. For the UK
720 dataset, nonparametric confidence intervals on the AUC were computed with DeLong's method
721^{51,52}. For the US dataset, in which each sample carried a scalar weight, the bootstrap was used
722 with 1000 replications.

723

724 On both datasets, we compared the sensitivity and specificity of the readers with that of a
725 thresholded score from the AI system. For the UK dataset, we knew the identity of each reader,
726 so statistics were adjusted for the clustered nature of the data using Obuchowski's method for
727 paired binomial proportions^{53,54}. Confidence intervals on the difference are Wald intervals⁵⁵ and
728 a Wald test was used for noninferiority⁵⁶. Both used the Obuchowski variance estimate.

729

730 For the US dataset, in which each sample carried a scalar inverse probability weight⁴⁵, we used
731 resampling methods⁵⁷ to compare the AI system's sensitivity and specificity with that of the pool
732 of radiologists. Confidence intervals on the difference were generated with the bootstrap method
733 with 1000 replications. A p -value on the difference was generated through the use of a
734 permutation test⁵⁸. In each of 10000 trials, the reader and AI system scores were randomly
735 interchanged for each case, yielding a reader-AI system difference sampled from the null
736 distribution. A two-sided p -value was computed by comparing the observed statistic to the
737 empirical quantiles of the randomization distribution.

738

739 In the reader study, each reader graded each case using a forced BI-RADS protocol (a score of
740 0 was not permitted), and the resulting values were treated as a 6-point index of suspicion for
741 malignancy. Scores of 1 and 2 were collapsed into the lowest category of suspicion; scores 3,
742 4a, 4b, 4c, and 5 were treated independently as increasing levels of suspicion. Because none of
743 the BI-RADS operating points reached the high sensitivity regime (see Figure 3), to avoid bias
744 from nonparametric analysis⁵⁹ we fit parametric ROC curves to the data using the proper
745 binormal model⁶⁰. This issue was not alleviated by using the readers' malignancy suspicion
746 ratings, which showed very strong correspondence with the BI-RADS scores (Supplementary
747 Figure 5). Since BI-RADS is used in actual screening practice, we elected to focus on these
748 scores for their superior clinical relevance. In a similar fashion, we fit a parametric ROC curve to
749 discretized AI system scores on the same data.

750

751 The performance of the AI system was compared to that of the panel of radiologists using
752 methods for the analysis of multi-reader multi-case (MRMC) studies standard in the radiology
753 community⁶¹. More specifically, we compared the AUC-ROC and pAUC-mLROC for the AI
754 system to that of the average radiologist using the ORH procedure, which was proposed in⁶²
755 and updated in⁶³. Originally formulated for the comparison of multiple imaging modalities, this
756 analysis has been adapted to the setting in which the population of radiologists operate on a
757 single modality and interest lies in comparing their performance to that of a standalone

758 algorithm⁶¹. The jackknife method was used to estimate the covariance terms in the model. The
759 p -value and confidence interval computation was conducted in Python using the numpy and
760 scipy packages and benchmarked against a reference implementation in the RJafroc library for
761 the R computing language⁶⁴.

762

763 Our primary comparisons numbered seven in total: sensitivity and specificity for the UK first
764 reader; sensitivity and specificity for the US clinical radiologist; sensitivity and specificity for the
765 US clinical radiologist using a model trained using only UK data; and the AUC-ROC in the
766 reader study. For comparisons with the clinical reads, the choice of superiority or non-inferiority
767 was based on what seemed attainable from simulations conducted on the validation set. For
768 non-inferiority comparisons, a 5% absolute margin was prespecified before inspecting the test
769 set. We employed a statistical significance threshold of 0.05. All seven p -values survived
770 correction for multiple comparisons using the Holm-Bonferroni method⁶⁵.

771

772 **Code availability.** The code used for training the models has a large number of dependencies
773 on internal tooling, infrastructure and hardware, and its release is therefore not feasible.

774 However, all experiments and implementation details are described in sufficient detail in the
775 Supplementary Methods section to allow independent replication with non-proprietary libraries.

776 Several major components of our work are available in open source repositories: Tensorflow:

777 <https://www.tensorflow.org>; Tensorflow Object Detection API:

778 https://github.com/tensorflow/models/tree/master/research/object_detection

779

780 **Data availability.** The dataset from Northwestern Medicine was used under license for the
781 current study, and is not publicly available. Applications for access to the OPTIMAM database
782 can be made at <https://medphys.royalsurrey.nhs.uk/omidb/getting-access/>.

783 **Methods References**

- 784 44. Pinsky, P. F. & Gallas, B. Enriched designs for assessing discriminatory performance--
785 analysis of bias and variance. *Stat. Med.* **31**, 501–515 (2012).
- 786 45. Mansournia, M. A. & Altman, D. G. Inverse probability weighting. *BMJ* **352**, i189 (2016).
- 787 46. Pathology reporting of breast disease in surgical excision specimens incorporating the
788 dataset for histological reporting of breast cancer. *Royal College of Pathologists* (2016).
789 Available at: <https://www.evidence.nhs.uk/document?id=1777849>. (Accessed: 22nd July
790 2019)
- 791 47. Chakraborty, D. P. & Yoon, H.-J. Operating characteristics predicted by models for
792 diagnostic tasks involving lesion localization. *Medical physics* **35**, 435–445 (2008).
- 793 48. Ellis, R. L., Meade, A. A., Mathiason, M. A., Willison, K. M. & Logan-Young, W. Evaluation
794 of computer-aided detection systems in the detection of small invasive breast carcinoma.
795 *Radiology* **245**, 88–94 (2007).
- 796 49. U.S. Food & Drug Administration. Evaluation of automatic class III designation for
797 OsteoDetect. (2018). Available at:
798 https://www.accessdata.fda.gov/cdrh_docs/reviews/DEN180005.pdf. (Accessed: 2nd
799 October 2019)
- 800 50. Hanley, J. A. & McNeil, B. J. The meaning and use of the area under a receiver operating
801 characteristic (ROC) curve. *Radiology* **143**, 29–36 (1982).
- 802 51. DeLong, E. R., DeLong, D. M. & Clarke-Pearson, D. L. Comparing the areas under two or
803 more correlated receiver operating characteristic curves: a nonparametric approach.
804 *Biometrics* **44**, 837–845 (1988).
- 805 52. Gengsheng Qin & Hotilovac, L. Comparison of non-parametric confidence intervals for the
806 area under the ROC curve of a continuous-scale diagnostic test. *Stat. Methods Med. Res.*
807 **17**, 207–221 (2008).

- 808 53. Obuchowski, N. A. On the comparison of correlated proportions for clustered data. *Stat.*
809 *Med.* **17**, 1495–1507 (1998).
- 810 54. Yang, Z., Sun, X. & Hardin, J. W. A note on the tests for clustered matched-pair binary
811 data. *Biom. J.* **52**, 638–652 (2010).
- 812 55. Fagerland, M. W., Lydersen, S. & Laake, P. Recommended tests and confidence intervals
813 for paired binomial proportions. *Stat. Med.* **33**, 2850–2875 (2014).
- 814 56. Liu, J.-P., Hsueh, H.-M., Hsieh, E. & Chen, J. J. Tests for equivalence or non-inferiority for
815 paired binary data. *Stat. Med.* **21**, 231–245 (2002).
- 816 57. Efron, B. & Tibshirani, R. J. *An Introduction to the Bootstrap*. (Springer US, 1993).
- 817 58. Chihara, L. M., Hesterberg, T. C. & Dobrow, R. P. *Mathematical Statistics with Resampling*
818 *and R & Probability with Applications and R Set*. (Wiley, 2014).
- 819 59. Gur, D., Bandos, A. I. & Rockette, H. E. Comparing areas under receiver operating
820 characteristic curves: potential impact of the ‘Last’ experimentally measured operating
821 point. *Radiology* **247**, 12–15 (2008).
- 822 60. Metz, C. E. & Pan, X. ‘Proper’ Binormal ROC Curves: Theory and Maximum-Likelihood
823 Estimation. *Journal of Mathematical Psychology* **43**, 1–33 (1999).
- 824 61. Chakraborty, D. P. *Observer Performance Methods for Diagnostic Imaging: Foundations,*
825 *Modeling, and Applications with R-Based Examples*. (CRC Press, 2017).
- 826 62. Obuchowski, N. A. & Rockette, H. E. Hypothesis testing of diagnostic accuracy for multiple
827 readers and multiple tests an anova approach with dependent observations.
828 *Communications in Statistics - Simulation and Computation* **24**, 285–308 (1995).
- 829 63. Hillis, S. L. A comparison of denominator degrees of freedom methods for multiple observer
830 ROC analysis. *Stat. Med.* **26**, 596–619 (2007).
- 831 64. CRAN - Package RJafroc. Available at: [https://cran.r-](https://cran.r-project.org/web/packages/RJafroc/index.html)
832 [project.org/web/packages/RJafroc/index.html](https://cran.r-project.org/web/packages/RJafroc/index.html). (Accessed: 29th January 2019)

833 65. Aikin, M. & Gensler, H. Adjusting for multiple testing when reporting research results: the
834 Bonferroni vs Holm methods. *Am. J. Public Health* **86**, 726–728 (1996).

835 66. Breast Screening Programme - NHS Digital. *NHS Digital* Available at:
836 [https://digital.nhs.uk/data-and-information/publications/statistical/breast-screening-](https://digital.nhs.uk/data-and-information/publications/statistical/breast-screening-programme)
837 [programme.](https://digital.nhs.uk/data-and-information/publications/statistical/breast-screening-programme) (Accessed: 17th July 2019)

838

839

840

841 **Acknowledgments**

842 We would like to acknowledge multiple contributors to this international project: Cancer
843 Research UK, the OPTIMAM project team, and staff at the Royal Surrey County Hospital who
844 developed the UK mammography imaging database; Sandra Tymms and Suzanne Steer for
845 providing patient perspectives and invaluable advice and input throughout the project; Robin
846 Wilson for providing screening programme and expert clinical perspective; all members of the
847 Etemadi Research Group for their efforts in data aggregation and deidentification; members of
848 Northwestern Medicine leadership without whom this couldn't be possible: Mark Schumacher,
849 Carl Christensen, Doug King, and Charles Hogue. We also thank the tireless efforts of everyone
850 at NMIT, including Mark Lombardi, Dalila Fridi, Paul Lendman, Bob Slavicek, Stavroula Xinos,
851 Bob Milfajt and others; Victoria Cornelius who provided expert advice on statistical planning;
852 Ross West & T Saensuksopa for assistance with data visualisation; Ali Eslami & Olaf
853 Ronneberger for machine learning expertise; Hannah Forbes and Claire Zaleski for assistance
854 with project management; Jonny Wong and Fraser Tan for coordinating labeling resources;
855 Razia Ahmed, Rory Pilgrim, Anthony Phalen and Michelle Bawn for meticulous work on
856 partnership formation; Roger Eng, Vashita Dhir, Rajshri Shah for data annotation and
857 interpretation; Cameron Chen for critical review of the manuscript; Diego Ardila for
858 infrastructure development; Cían Hughes and Diogo Moitinho de Almeida for early engineering
859 work; Jessica Yoshimi, Xiang Ji, Will Chen, Teagan Daly, Huy Doan, Eric Lindley, Quang Duong
860 for development of the labeling infrastructure. Professor Fiona Gilbert receives funding from the
861 National Institute for Health Research (Senior Investigator award). The views expressed are
862 those of the authors and not necessarily those of the NIHR or the Department of Health and
863 Social Care.

864 **Author contributions**

865 A.K., A.D., D.H., D.K., H.M, G.C.C., J.D.F., J.R.L., K.C.Y., L.P., M.D.H.B., M.T.S., M.S., R.S.,
866 S.M.M., S.S, and T.B. contributed to conception; A.K., B.R.P., C.J.K., D.H., D.T., F.J.G., J.D.F.,
867 J.R.L., K.C.Y., L.P., M.D.H.B., M.C., M.E., M.T.S., M.S., N.A., R.S., S.J., S.M.M., S.S., T.B. and
868 V.G. contributed to design; D.M. D.T., F.G.V., G.C.C., H.M., J.D.F., J.G., K.C.Y., L.P., M.D.H.B.,
869 M.C., M.E., M.T.S., S.M.M., S.S., and V.G. contributed to acquisition; A.K., A.D, B.R.P., C.J.K.,
870 F.J.G., H.A., J.D.F., J.G., J.J.R., M.S., N.A., R.S., S.J., S.M.M., S.S. and V.G. contributed to
871 analysis and interpretation; A.K., C.J.K., D.T., F.J.G., J.D.F., J.G., J.J.R, M.T.S., N.A., R.S, S.J.,
872 S.M.M., S.S., and V.G. contributed to drafting and revising the manuscript.

873

874 **Competing interests**

875 This study was funded by Google LLC and/or a subsidiary thereof (“Google”). S.M.M., M.T.S.,
876 V.G., J.G., N.A., T.B., M.C., G.C.C., D.H., S.J., A.K., C.J.K, D.K., J.R.L., H.M., B.R.-P., L.P.,
877 M.Su., D.T., J.D.F., and S.S. are employees of Google and own stock as part of the standard
878 compensation package. J.J.R., R.S., F.J.G., and A.D. are paid consultants of Google. M.E.,
879 F.G.-V., D.M., K.C.Y., and M.H.-B received funding from Google to support the research
880 collaboration. The authors have no other competing interests to disclose.

881

882

883 **Extended Data Tables**

884 **Extended Data Table 1. Characteristics of the UK and UK test sets.** For each feature, we
885 constructed a joint 95% confidence interval on the proportions in each category. **a**, The UK test
886 set was drawn from two sites in the UK over a four-year period. For reference, we present the
887 corresponding statistics from the broader UK Breast Screening Programme (BSP)⁶⁶. For
888 comparison with national numbers, only screen-detected cancers are reported here. **b**, The US
889 test was drawn from one academic medical center over an eighteen-year period. For reference,
890 we present the corresponding statistics from the broader US screening population, as reported
891 by the Breast Cancer Surveillance Consortium (BCSC)². Cancers reported here occurred within
892 12 months of screening.

893

894 **Extended Data Table 2. Detailed comparison between human clinical decisions and AI**
895 **predictions.**

896 **a.** Comparison of sensitivity and specificity between human benchmarks, derived retrospectively
897 from the clinical record, and the predictions of the AI system. Score thresholds were chosen,
898 based on separate validation data, to match or exceed the performance of each human
899 benchmark (Methods, Operating point selection). These points are depicted graphically in
900 Figure 2a. Bolded quantities represent estimated differences which are statistically significant
901 for superiority; all others are statistically noninferior at a prespecified 5% margin. Note that the
902 number of cases (N) differs from Figure 2a because a radiologist opinion was not available for
903 all images. We also note that sensitivity and specificity metrics are not easily comparable to
904 most prior publications in breast imaging (eg. the DMIST Trial³⁴) given differences in follow up
905 interval. Negative cases in the US dataset were upweighted to account for the sampling protocol
906 (Methods, Inverse probability weighting).

907 **b.** Same columns as A, but using a version of the AI system trained exclusively on the UK
908 dataset. It was tested on the US dataset to show generalisability of the AI across different
909 populations and healthcare systems. Superiority comparisons on the UK data were conducted
910 using Obuchowski's extension of the two-sided McNemar test for clustered data. Noninferiority
911 comparisons were Wald tests using the Obuchowski correction. Comparisons on the US data
912 were performed with a two-sided permutation test. All p -values survived correction for multiple
913 comparisons (Methods, Statistical analysis).

914

915

916 **Extended Data Table 3. Detailed description of reader study case composition.**
917 **Row 1.** 500 cases were selected for the reader study. The case mixture was enriched for
918 positives as well as challenging negatives.
919 **Row 2.** Cases containing breast implants and those for which at least half of the readers
920 indicated image quality concerns were excluded from analysis. The remaining 465 cases are
921 represented in Figure 3a and b.
922 **Row 3.** We also restricted the cancers to those that developed within 12 months. Those that
923 developed cancer later (but within 27 months) were excluded because they did not meet the
924 follow-up criteria to be considered negative. The remaining 408 cases are represented in
925 Extended Data Table 2c and Extended Data Figure 2.
926 **Row 4.** To perform localisation analysis, the areas of malignancy were determined using follow-
927 up biopsy data. In three instances, ground truth could not reliably be determined. The remaining
928 405 cases are represented in Extended Data Figure 3.
929
930 **Extended Data Table 4. Potential utility of the AI system in two clinical applications.**
931 **a.** Simulation using the UK test set in which the AI system is used in place of the second reader
932 when it concurs with the first reader. In cases of disagreement (12.02%) the consensus opinion
933 was invoked. The high performance of this combination of human and machine suggests that
934 approximately 88% of the second reader's effort can be eliminated while maintaining the
935 standard of care produced by double reading. The AI system's decision was generated using
936 operating point (i) in Figure 2a. Confidence intervals are Wald intervals computed with the
937 Obuchowski correction for clustered data.
938 **b.** Evaluation of the AI system for low-latency triage. Operating points were set to perform with
939 high NPV and PPV for detecting cancer in 12 months.
940
941 **Extended Data Table 5. Discrepancies between the AI system and human readers.**
942 For the UK comparison, we used the first reader operating point (i) shown in Figure 2a. For the
943 US comparison, we used the operating point shown in Figure 2b.
944
945 **Extended Data Table 6. Performance breakdowns.** Analysis excludes technical recalls and
946 US cases for which BI-RADS scores were unavailable.
947 **a.** Sensitivity across cancer subtypes in the UK data. We used the AI system operating point (i)
948 in Figure 2a. Also shown is the first reader performance on the same subset.
949 **b.** Sensitivity across cancer subtypes in the US data. We used the AI system operating point (i)
950 in Figure 2a. Reader performance was derived from the clinical BI-RADS scores on the same
951 subset. ILC = Invasive lobular carcinoma, IDC = invasive ductal carcinoma, DCIS = ductal
952 carcinoma in situ.
953 **c.** Performance across breast density categories. BI-RADS³⁵ breast density was extracted from
954 the radiology report rendered at the time of screening, only available in the US dataset. We
955 used the AI system operating point shown in Figure 2b. Adjusted specificities were computed
956 using inverse probability weighting (Methods).
957
958

959 **Extended Data Table 7.** Reader experience from the UK clinical dataset (a) and the
960 independent reader study (b).
961

962 **Extended Data Figures**

963 **Extended Data Figure 1. Unweighted evaluation of breast cancer prediction on the US**
964 **test set.** Unlike in Figure 2b, the sensitivity and specificity were computed without the use of
965 inverse probability weights to account for the spectrum-enrichment of the study population.
966 Since hard negatives are overrepresented, the specificity of both the AI system and the human
967 readers is reduced. The unweighted human sensitivity and specificity are 48.10% (n = 553) and
968 69.65% (n = 2,185), respectively.
969

970 **Extended Data Figure 2. Breast cancer prediction performance compared to six**
971 **independent readers with a 12-month follow up for cancer status.** While the mean reader
972 AUC was 0.750 (s.d. 0.049), the AI system achieved an AUC of 0.871 (95% CI: [0.785, 0.919]).
973 The AI system exceeded human performance by a significant margin ($\Delta = +0.121$, 95% CI:
974 [0.070, 0.173], $p = 0.0018$, two-sided ORH method). In this analysis, there were 56 positives of
975 408 total cases; see Extended Data Table 3. Note that this sample of cases was enriched for
976 patients that had received a negative biopsy result (n=119), making this a more challenging
977 population for screening. Since these external readers were not gatekeepers for follow up and
978 eventual cancer diagnosis, there was no bias in favour of reader performance at this shorter
979 time horizon. See Figure 3A for a comparison with a time interval chosen to encompass a
980 subsequent screening exam.
981

982 **Extended Data Figure 3. Multi-location receiver operating characteristic (mLROC)**
983 **analysis.**

984 Similar to Extended Data Figure 2, but true positives require localisation of a malignancy in any
985 of the four mammogram views (Methods, Localisation analysis). Here, the cancer interval was
986 12 months (n = 53 positives of 405 cases; see Extended Data Table 3). The dotted line
987 indicates a false positive rate of 10%, which was used as the right-hand boundary for the partial
988 AUC (pAUC) calculation. The mean reader pAUC was 0.029 (s.d. 0.005), while the AI system's
989 pAUC was 0.048 (95% CI: [0.035, 0.061]). The AI system exceeded human performance by a
990 significant margin ($\Delta = +0.0192$, 95% CI: [0.0086, 0.0298], $p = 0.0004$, two-sided ORH method).
991

992 **Extended Data Figure 4. Evidence for the gatekeeper effect in retrospective datasets.**

993 These figures show the change in observed reader sensitivity in the UK (a) and the US (b) as
994 the cancer follow-up interval is extended. At short intervals, measured reader sensitivity is
995 extremely high, owing to the fact that biopsies are only triggered based on radiological
996 suspicion. As the time interval is extended, the task becomes more difficult and measured
997 sensitivity declines. Part of this decline stems from the development of new cancers that were
998 impossible to detect at initial screening. However, more precipitous drops occur when the
999 follow-up window encompasses the screening interval (36 months in the UK, 12 and 24 months

1000 in the US). This is suggestive of what happens to reader metrics when gatekeeper bias is
1001 mitigated by another screening examination.

1002

1003 **Extended Data Figure 5. Quantitative evaluation of reader and AI system performance**
1004 **with a 12-month follow-up interval for ground truth cancer positive status.**

1005 Because a 12-month follow-up interval is unlikely to encompass a subsequent screening exam
1006 in either country, reader-model comparisons on retrospective clinical data may be contaminated
1007 by the gatekeeper effect (Extended Data Figure 4). See Figure 2 for comparison with longer
1008 time intervals.

1009 **a.** AI system performance on UK data. This plot was derived from a total of 25,717 eligible
1010 examples including 274 positives. The AI system achieved an AUC of 0.966 ([0.954, 0.977],
1011 95% CI).

1012 **b.** AI system performance on US data. This plot was derived from a total of 2,770 eligible
1013 examples including 359 positives. The AI system achieved an AUC of 0.883 ([0.859, 0.903],
1014 95% CI).

1015 **c.** Reader performance. In computing reader metrics on the UK data, we excluded cases for
1016 which the reader recommended repeat mammography to address technical issues. In the US
1017 data, radiologist performance could only be assessed on the subset of cases for which a BI-
1018 RADS grade was available.

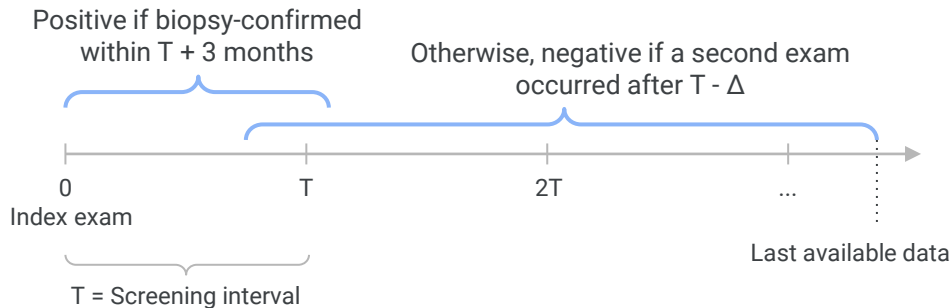
1019

Evaluation data sets



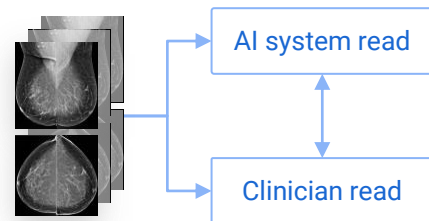
Number of women	25,856	3,097
Clinical evaluation	Double reader	Single reader
Screening interval	3 years	1 or 2 years
Cancer follow-up	39 months	27 months
Number of cancers	414 (1.6%)	686 (22.2%)

Ground truth determination

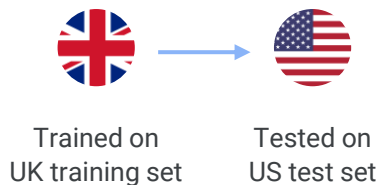


Evaluation

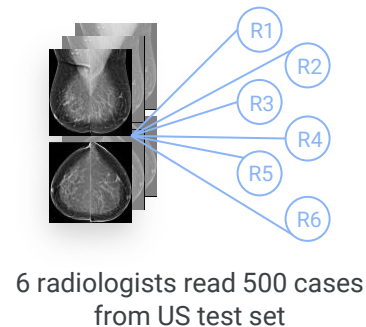
Comparison with retrospective clinical performance



Generalization across data sets

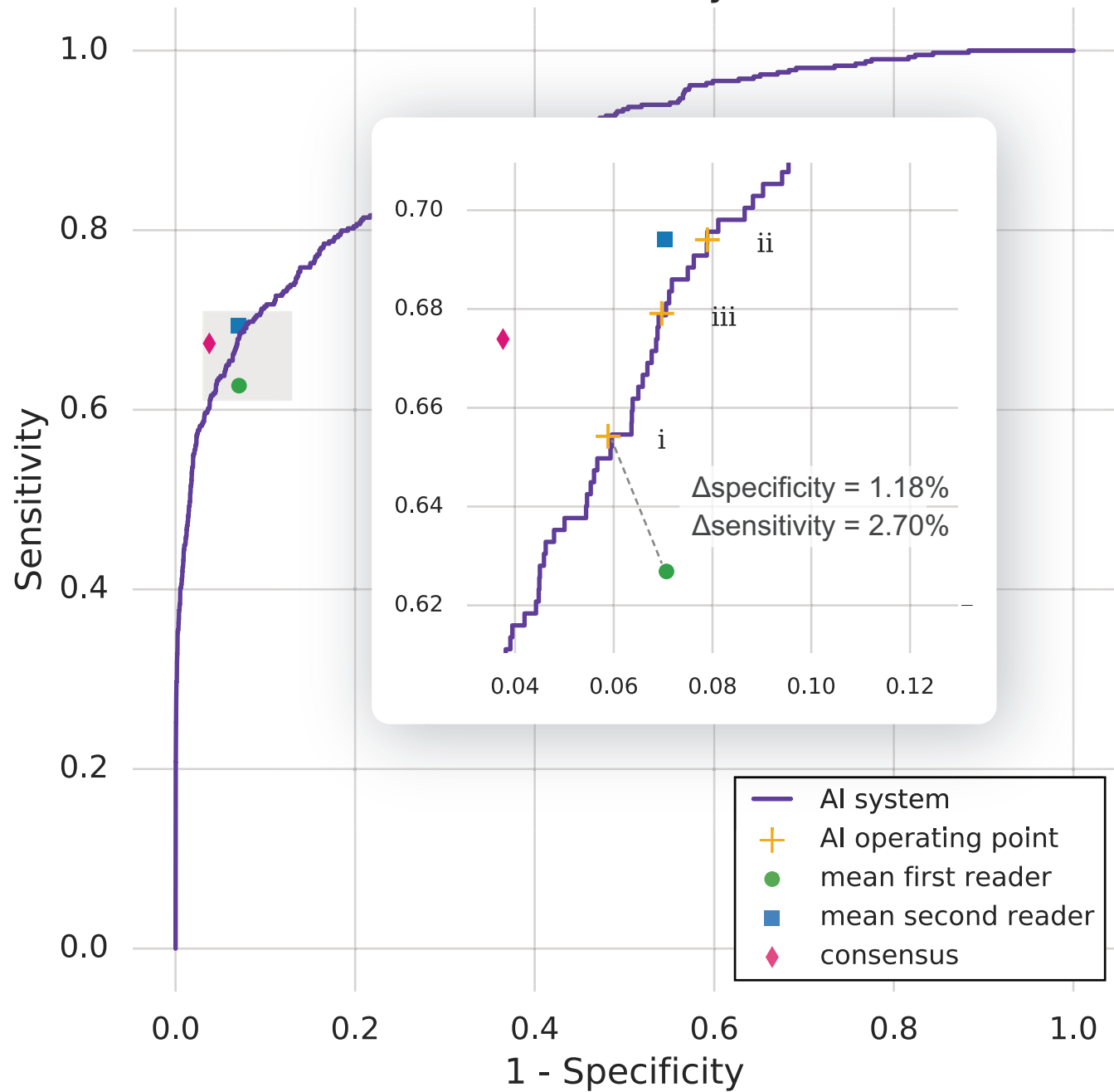


Independently-conducted reader study

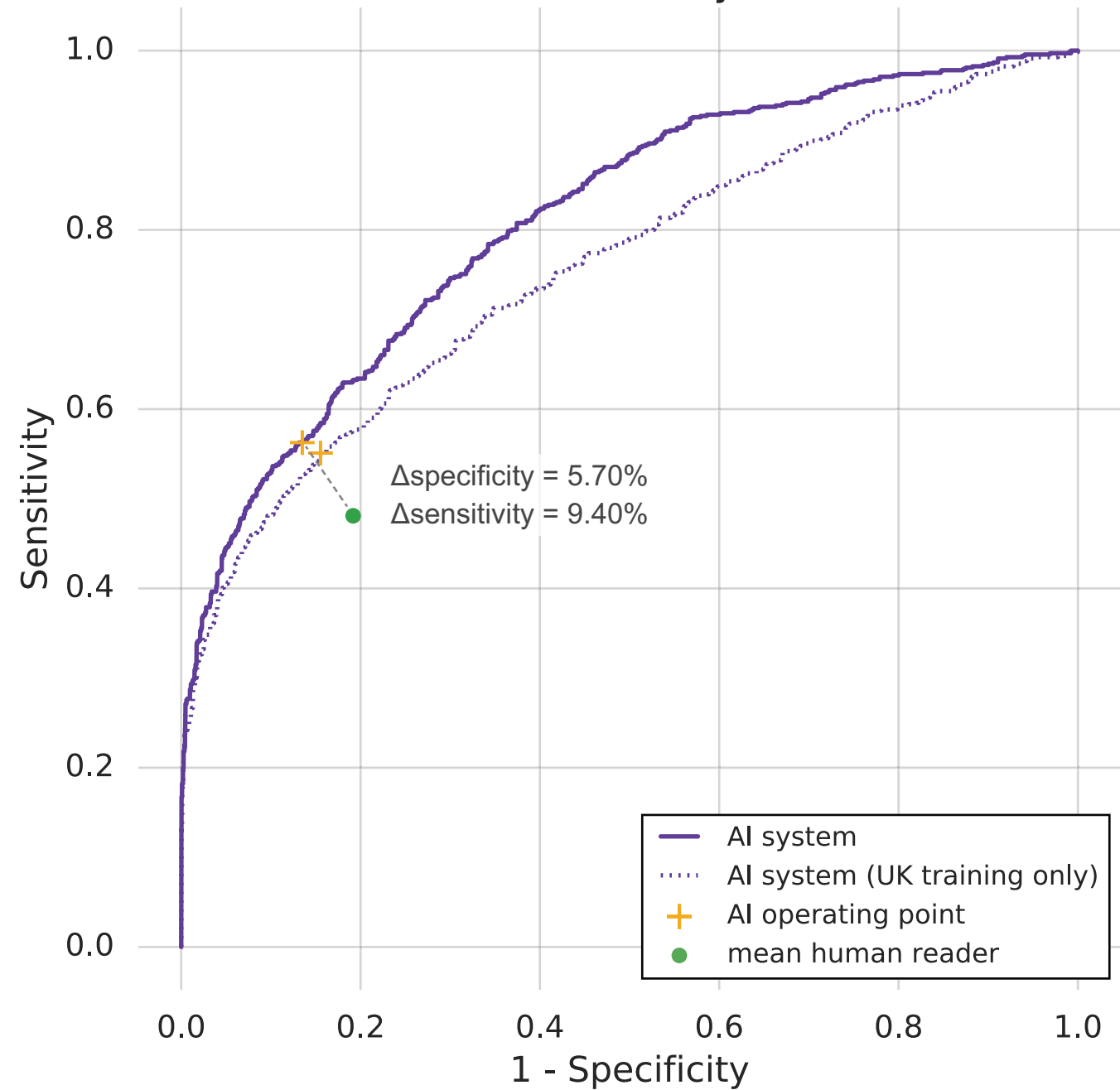


a

Breast cancer in 3 years (UK)

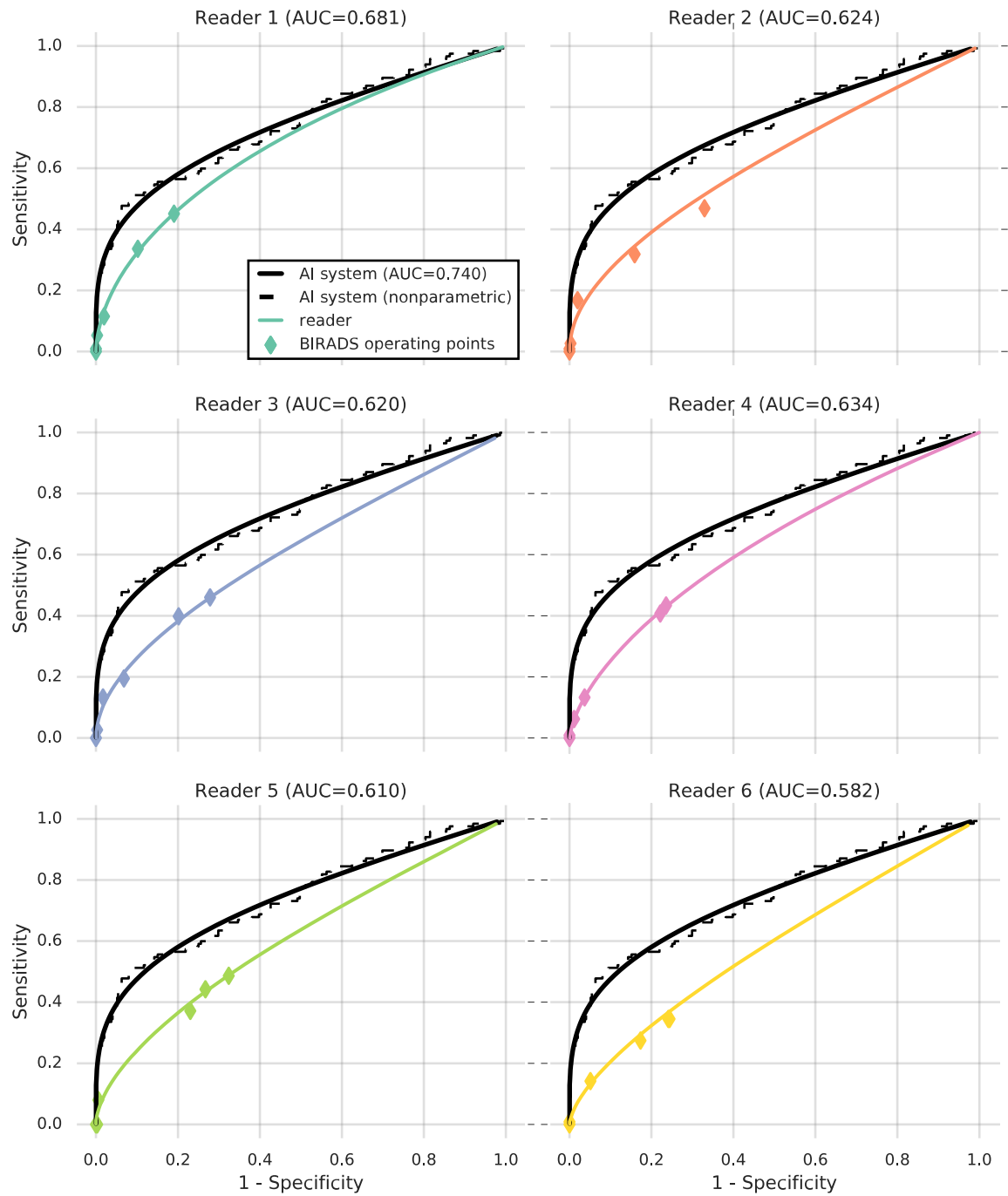
**b**

Breast cancer in 2 years (US)

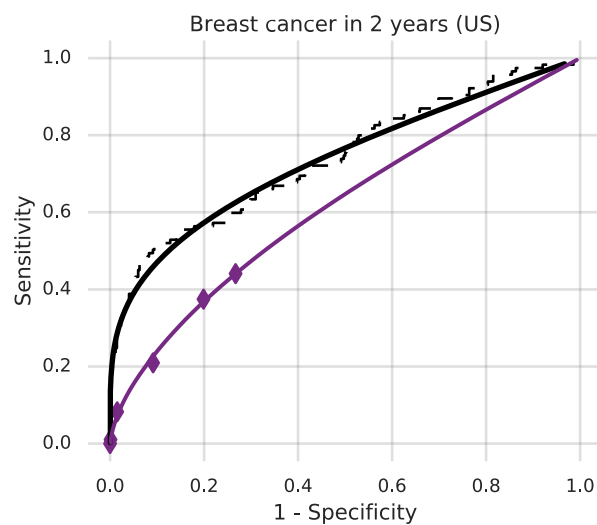


Breast cancer in 2 years (US)

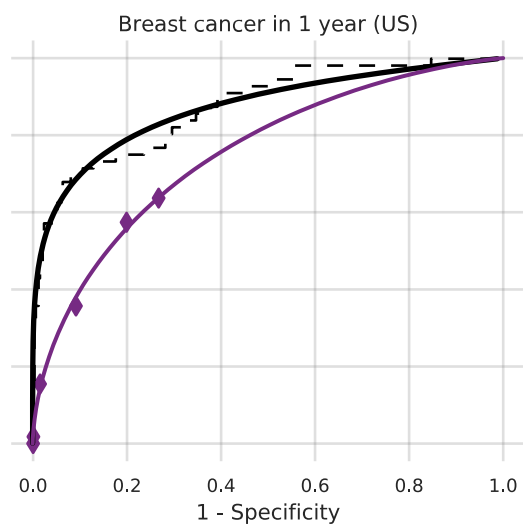
a



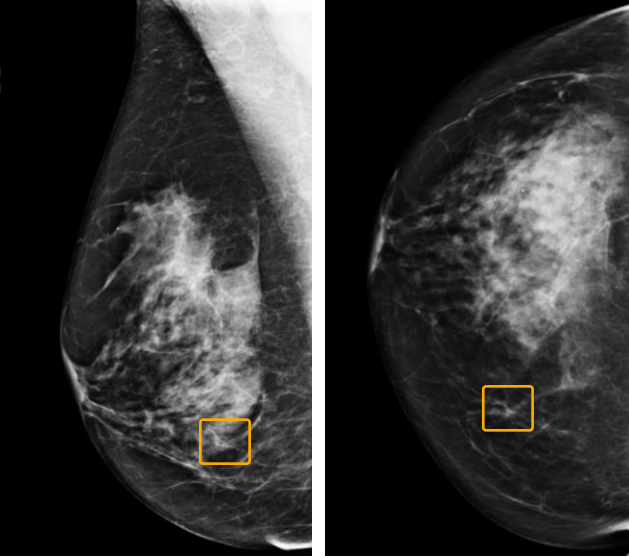
b



c



a



b

