

1 **Title**

2 Deep learning saliency maps do not accurately highlight diagnostically relevant regions
3 for medical image interpretation

4

5 **Authors**

6 Adriel Saporta BA MBA^{1*}, Xiaotong Gui BA^{1*}, Ashwin Agrawal MS^{1*}, Anuj Pareek MD
7 PhD², Steven QH Truong MBA³, Chanh DT Nguyen PhD^{3,4}, Van-Doan Ngo MD⁵, Jayne
8 Seekins DO⁶, Francis G. Blankenberg MD⁶, Andrew Y. Ng PhD¹, Matthew P. Lungren MD
9 MPH^{2†}, Pranav Rajpurkar MS^{1†}

10

11 **Affiliations**

12 ¹Stanford University Department of Computer Science, USA

13 ²Stanford Center for Artificial Intelligence in Medicine and Imaging, USA

14 ³VinBrain, Vietnam

15 ⁴VinUniversity, Vietnam

16 ⁵Vinmec International Hospital, Vietnam

17 ⁶Stanford University School of Medicine, Department of Radiology, USA

18 *These authors contributed equally: Adriel Saporta, Xiaotong Gui, Ashwin Agrawal

19 †These authors contributed equally: Matthew P. Lungren, Pranav Rajpurkar

20 Corresponding author: Pranav Rajpurkar (email: pranavs@cs.stanford.edu)

21

22 Current word count: 3941

23 **Abstract**

24 Deep learning has enabled automated medical image interpretation at a level often
25 surpassing that of practicing medical experts. However, many clinical practices have cited
26 a lack of model interpretability as reason to delay the use of “black-box” deep neural
27 networks in clinical workflows. Saliency maps, which “explain” a model’s decision by
28 producing heat maps that highlight the areas of the medical image that influence model
29 prediction, are often presented to clinicians as an aid in diagnostic decision-making. In
30 this work, we demonstrate that the most commonly used saliency map generating
31 method, Grad-CAM, results in low performance for 10 pathologies on chest X-rays. We
32 examined under what clinical conditions saliency maps might be more dangerous to use
33 compared to human experts, and found that Grad-CAM performs worse for pathologies
34 that had multiple instances, were smaller in size, and had shapes that were more
35 complex. Moreover, we showed that model confidence was positively correlated with
36 Grad-CAM localization performance, suggesting that saliency maps were safer for
37 clinicians to use as a decision aid when the model had made a positive prediction with
38 high confidence. Our work demonstrates that several important limitations of
39 interpretability techniques for medical imaging must be addressed before use in clinical
40 workflows.

41

42 **Introduction**

43 Deep learning has enabled automated medical imaging interpretation at a level shown to
44 surpass that of practicing experts in some settings¹⁻³. While the potential benefits of
45 automated diagnostic models are numerous, lack of model interpretability in the use of

46 “black-box” deep neural networks (DNNs) represents a major barrier to clinical trust and
47 adoption^{4,5,6}. In fact, it has been argued that the European Union’s recently adopted
48 General Data Protection Regulation (GDPR) affirms an individual’s right to an explanation
49 in the context of automated decision-making⁷. Although many DNN interpretability
50 techniques have been proposed, rigorous investigation of the accuracy and reliability of
51 these strategies is lacking and necessary before they are integrated into the clinical
52 setting⁸.

53
54 One type of DNN interpretation strategy widely used in the context of medical imaging is
55 based on saliency (or pixel-attribution) methods^{9–12}. Saliency methods produce heat
56 maps that highlight the areas of the medical image that most influenced the DNN’s
57 prediction. The heat maps help to visualize whether a DNN is concentrating on the same
58 regions of the medical image that a human expert would focus attention on for a given
59 diagnosis, rather than concentrating on a clinically irrelevant part of the medical image or
60 even on confounders in the image^{13–15}. However, recent work has shown that saliency
61 methods used to validate model predictions can be misleading in some cases and may
62 lead to increased bias and loss of user trust with concerning implications for clinical
63 translation efforts¹⁶.

64
65 The purpose of this work is to perform a systematic evaluation of the most common
66 saliency method, Grad-CAM¹⁷, on multi-label classification models for medical imaging
67 interpretation from chest X-rays. In order to evaluate how well the saliency method
68 identifies critical areas of an image for diagnosis, we compared the saliency method

69 segmentations to human expert benchmark and reference annotations. We evaluated the
70 accuracy of the saliency method segmentations and the human expert benchmark
71 segmentations first by calculating their overlap with the human expert reference
72 segmentations, and then by determining whether the segmentations correctly located the
73 pathology of concern, regardless of the exact bounds of the segmentations. We further
74 conducted statistical analyses to better understand how the localization accuracy of
75 saliency methods is affected both by pathological characteristics and also by model
76 confidence.

77

78 **Results**

79 **Framework for evaluating a saliency method on multi-label classification models**

80 A model can be trained to perform pixel-level localization in one of two ways: either
81 through supervised learning, in which the model is trained directly on pixel-level
82 segmentations, or through weakly supervised learning, in which the model is trained only
83 on image-level class labels. Because ground truth segmentations for medical imaging
84 can be especially time-consuming and expensive to obtain given the domain expertise
85 required to create them¹⁸, one of the most common saliency methods in medical imaging
86 is a weakly supervised localization technique, Grad-CAM, in which the classification
87 model is never exposed to pixel-level segmentations during training. Instead, Grad-CAM
88 generates a heat map corresponding to each image-level task label that highlights the
89 regions of the input image that, in theory, are most indicative of that task label. This
90 saliency method has been widely used for a variety of medical imaging tasks and
91 modalities including but not limited to: visualizing the performance of a convolutional

92 neural network in predicting (1) myocardial infarction¹⁹ and hypoglycemia²⁰ from
93 electrocardiograms, (2) visual impairment²¹, refractive error²², and anaemia²³ from retinal
94 photographs (3) long-term mortality²⁴ and tuberculosis²⁵ from chest x-ray images, and (4)
95 appendicitis²⁶, and pulmonary embolism²⁷ on computed tomography scans.

96

97 We use the following framework, which we call CheXplanation, to evaluate Grad-CAM in
98 a multi-label classification setup. We trained and evaluated an ensemble of 30 CNN
99 models on CheXpert, a large publicly available CXR dataset²⁸. We then passed each of
100 the 668 CXRs in the dataset's holdout test set into the ensemble model to obtain image-
101 level predictions for the following 10 pathologies: Atelectasis, Cardiomegaly,
102 Consolidation, Edema, Enlarged Cardiomeastinum, Lung Lesion, Lung Opacity, Pleural
103 Effusion, Pneumothorax, and Support Devices. For each CXR, we used Grad-CAM to
104 generate 10 heat maps, one for each of the 10 pathologies. We then applied a threshold
105 to the heat maps to produce 10 binary segmentations in order to evaluate their overlap
106 with the human expert reference segmentations (see Fig. 1a). We also extracted the
107 location of the pixel with the largest value from each heat map to determine whether it fell
108 within the bounds of the reference segmentation (see Fig. 1d). In doing so, we could
109 evaluate Grad-CAM's localization performance regardless of the exact bounds of its
110 binary segmentations.

111

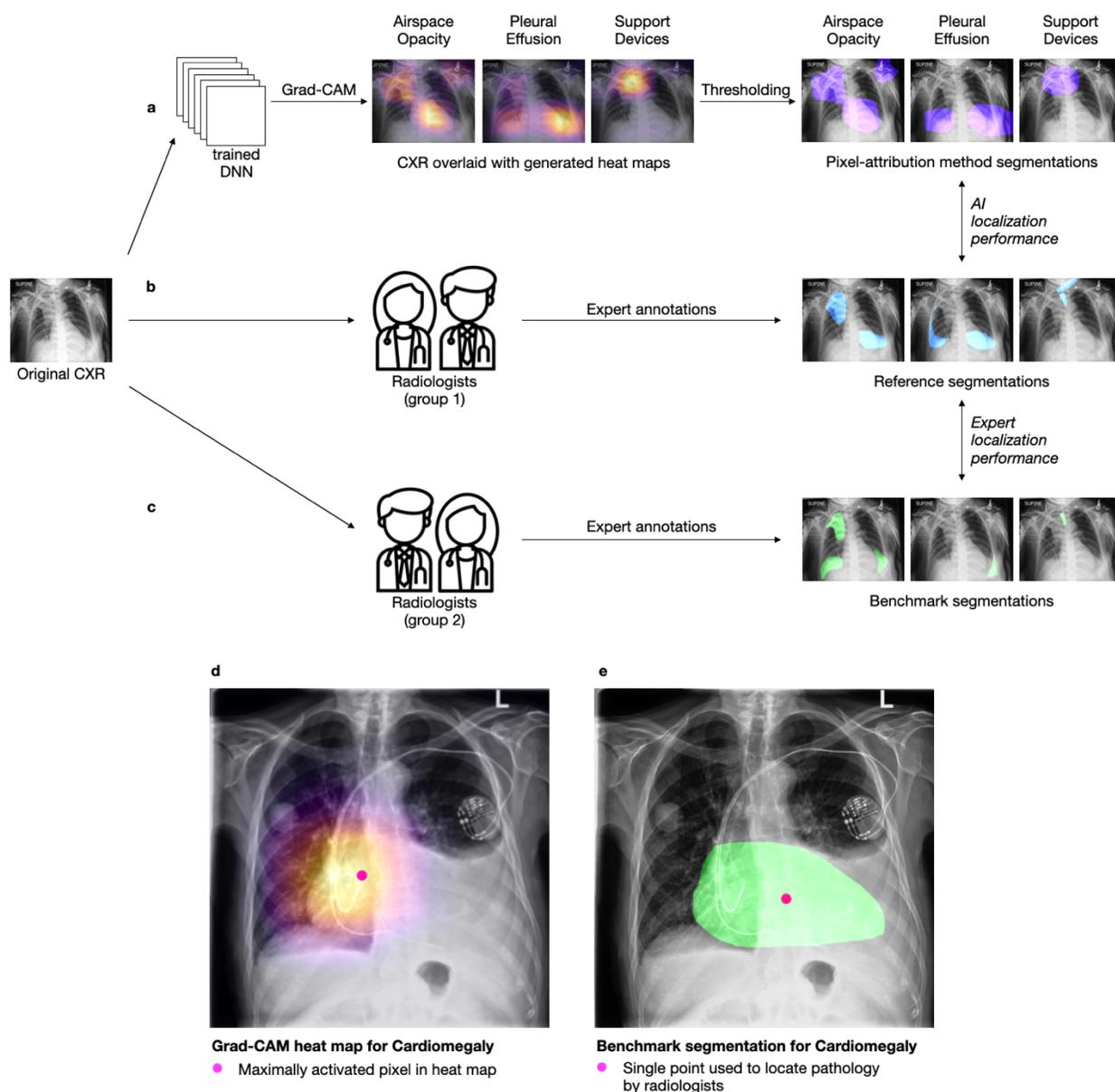
112 To evaluate how well the saliency method segmentations identified clinically relevant
113 pathology regions of input CXRs ("AI localization performance"), we obtained pixel-level
114 reference segmentations on the holdout test set from two board-certified radiologists who

115 were asked to segment any of the 10 pathologies that were present in each CXR as
116 determined by the dataset's ground-truth labels (see Fig. 1b). We also established a
117 human benchmark ("expert localization performance") by collecting segmentations from
118 three additional radiologists who were asked to segment the 10 pathologies of interest
119 present in each CXR as determined by the dataset's ground-truth labels (there was no
120 overlap between these three radiologists and the two who drew the reference
121 segmentations) (see Fig. 1c). This second group of radiologists was also asked to locate
122 each pathology present on each CXR using only a single most representative point for
123 that pathology on the CXR to see whether that point fell within the bounds of the reference
124 segmentation (see Fig. 1e). Which point in the CXR is "most representative" varied by
125 pathology, but the point would always lie inside the pathology's segmentation drawn by
126 the radiologist. For example, the most representative point for Pneumothorax would be
127 wherever the pathology is most evident, whereas the most representative point for
128 Cardiomegaly would be the center of the heart. See Supplementary Figs. 1 through 11
129 for the detailed instructions given to the radiologists.

130

131 Our dataset of expert reference segmentations has been made publicly available to
132 encourage further development and evaluation of CXR interpretation models.

133



134

135 **Fig. 1 | Framework for evaluating a saliency method on multi-label classification**

136 **models.** **a**, Left, a CXR image from the holdout test set is passed into an ensemble DNN

137 trained only on CXR images and their corresponding pathology task labels. Grad-CAM is

138 used to generate 10 heat maps for the example CXR, one for each task. Middle, there

139 are three pathologies present in this CXR (Airspace Opacity, Pleural Effusion, and

140 Support Devices). Right, a threshold is applied to the heat maps to produce binary

141 segmentations for each present pathology. **b**, Two board-certified radiologists were asked

142 to segment the pathologies present in the CXR as determined by the dataset's ground-

143 truth labels. Saliency method segmentations are compared to these reference
144 segmentations to evaluate how well Grad-CAM identifies the clinically relevant areas of
145 the input CXR (“AI localization performance”). **c**, Three radiologists (separate from those
146 in **b**) were asked to segment the pathologies present in the CXR as determined by the
147 dataset’s ground-truth labels. These benchmark segmentations are compared to the
148 reference segmentations to determine a human benchmark (“expert localization
149 performance”). **d**, The location of the pixel with the largest value was extracted from each
150 heat map. **e**, In addition to drawing segmentations, the benchmark radiologists were
151 asked to locate each pathology present on each CXR using only a single point on that
152 CXR.

153

154 **Evaluating the localization performance of the saliency method**

155 We used two evaluation schemes to compare the AI localization performance to the
156 expert localization performance. First, we used mean Intersection over Union (mIoU) to
157 measure how much, on average, either the saliency method segmentations or the human
158 benchmark segmentations overlapped with the reference segmentations. Second, we
159 used the pointing game setup²⁹, in which a “hit” is when the single point used to locate a
160 pathology lies within the reference segmentation and a “miss” is when the single point
161 lies outside the reference segmentation. Localization performance is then calculated as
162 the hit rate across the dataset³⁰. See Fig. 2a for mIoU and hit rate for Grad-CAM
163 segmentations on example CXRs.

164

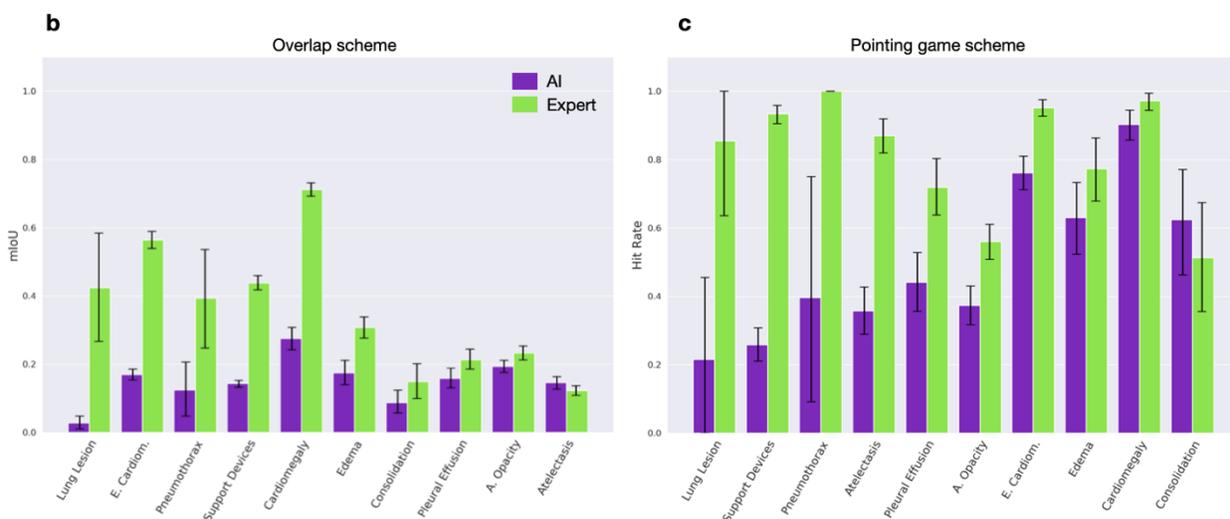
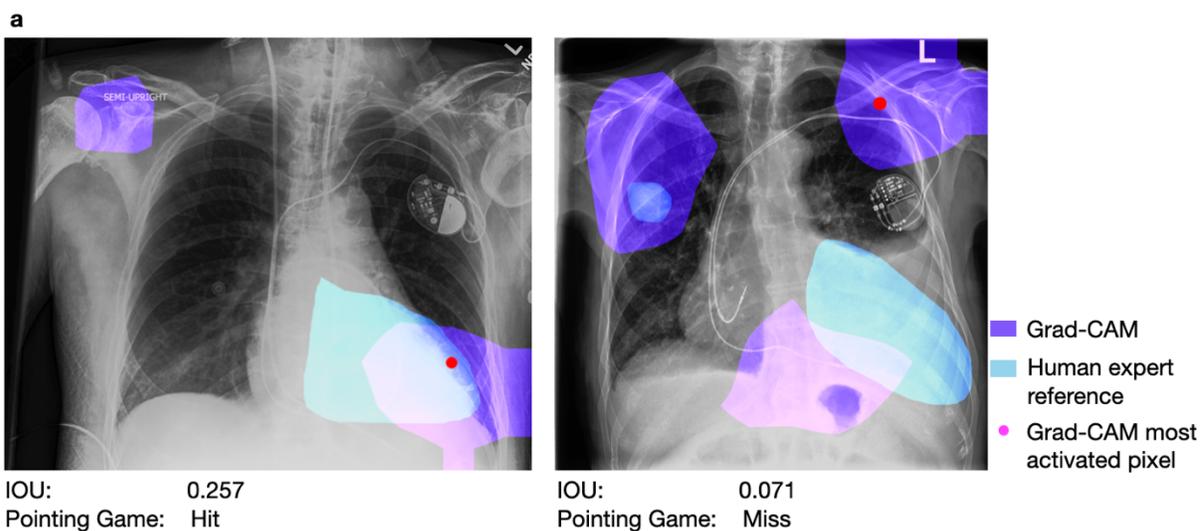
165 For nine of the 10 pathologies, the saliency segmentations had a lower overlap with the
166 reference segmentations than did the human benchmark segmentations (see Fig. 2b).
167 Similarly, for nine of the 10 pathologies, the saliency method had a lower hit rate than the
168 human benchmark (see Fig. 2c). Under both the overlap and the pointing game evaluation

169 schemes, the gap between AI localization performance and expert localization
170 performance was the largest for Lung Lesion: AI mIoU was 93.6% smaller than expert
171 mIoU, and Grad-CAM hit rate was 74.9% smaller than expert hit rate. Lung Lesion is also
172 the pathology for which the saliency method had both the lowest mIoU (0.027) and the
173 lowest hit rate (0.215). Under the hit rate scheme, Support Devices and Pneumothorax
174 displayed the second and third largest gaps, respectively, between AI and expert
175 localization performance (Support Devices: 72.4%; Pneumothorax: 60.4%). Under the
176 overlap scheme, Support Devices and Pneumothorax also displayed the fourth and third
177 largest gaps, respectively, between AI and expert localization performance (Support
178 Devices: 67.4%; Pneumothorax: 68.5%).

179
180 Under the overlap evaluation scheme, the only pathology for which the saliency method
181 (mIoU 0.145 95% CI [0.126, 0.163]) outperformed the human benchmark (mIoU 0.122,
182 95% CI [0.108, 0.136]) was Atelectasis. Under the hit rate evaluation scheme, the only
183 pathology for which the saliency method (hit rate 0.624 95% CI [0.462, 0.771])
184 outperformed the human benchmark (hit rate 0.513 95% CI [0.355, 0.674]) was for
185 Consolidation. However, the differences were statistically insignificant. Both the saliency
186 method and the human benchmark achieved their highest mIoUs for Cardiomegaly (AI:
187 0.275 95% CI [0.242, 0.30]); expert: 0.712 95% CI [0.692, 0.731]). AI hit rate was largest
188 for Cardiomegaly (0.903 95% CI [0.857, 0.945]) and Enlarged Cardiomedastinum (0.761
189 95% CI [0.712, 0.81]). Expert hit rate was above 0.95 for Pneumothorax (1.0 95% CI [1.0,
190 1.0]), Cardiomegaly (0.971 95% CI [0.944, 0.994]), and Enlarged Cardiomedastinum
191 (0.953 95% CI [0.927, 0.975]).

192

193 We also evaluated whether and how the gap between AI and expert localization
 194 performance changed from when overlap was used as an evaluation metric to when hit
 195 rate was used as an evaluation metric. Percentage decrease from expert to AI localization
 196 performance fell the most from mIoU to hit rate for Consolidation (41.2% - (-21.6%) =
 197 62.8%), Cardiomegaly (61.4% - 7.1% = 54.3%), and Enlarged Cardiomedastinum (70.0%
 198 - 20.1% = 49.9%). Percentage decrease from expert to AI localization performance
 199 increased by far the most from mIoU to hit rate for Atelectasis (-18.9% - 59.0% = -77.9%).



200

201 **Fig. 2 | Evaluating the localization performance of the attribution method. a**, Grad-
 202 CAM and human expert reference segmentations for two CXRs with Airspace Opacity.
 203 Left, IoU score is 0.257, and pointing game is a “hit” since Grad-CAM’s most activated
 204 pixel is inside of the reference segmentation. Right, IoU score is 0.071, and pointing game
 205 is a “miss” since Grad-CAM’s most activated pixel is outside of the reference
 206 segmentation. **b**, Comparing AI and expert localization performances under the overlap
 207 evaluation scheme. Pathologies are sorted on the x-axis in descending order of
 208 percentage decrease from expert mIoU to AI mIoU. **c**, Comparing AI and expert
 209 localization performances under the hit rate evaluation scheme. Pathologies are sorted
 210 on the x-axis in descending order of percentage decrease from expert hit rate to AI hit
 211 rate for each pathology. The black error bars indicate 95% bootstrap confidence interval.

Table 1 Percentage decrease from expert localization to AI localization for each pathology						
Pathology	Expert mIoU	AI mIoU	% decrease mIoU	Expert hit rate	AI hit rate	% decrease hit rate
Lung Lesion	0.424	0.027	93.6	0.855	0.215	74.9
E. Cardiom.	0.564	0.169	70.0	0.953	0.761	20.1
Pneumothorax	0.393	0.124	68.4	1.000	0.396	60.4
Support Devices	0.438	0.143	67.4	0.934	0.258	72.4
Cardiomegaly	0.712	0.275	61.4	0.971	0.903	7.1
Edema	0.307	0.174	43.3	0.773	0.630	18.5
Consolidation	0.148	0.087	41.2	0.513	0.624	-21.6
Pleural Effusion	0.213	0.158	25.8	0.719	0.441	38.7
Airspace Opacity	0.232	0.193	16.8	0.560	0.317	33.4
Atelectasis	0.122	0.145	-18.9	0.869	0.289	59.0

212
213

214 **Characterizing the gaps between AI localization performance and expert**
 215 **localization performance**

216 In order to better understand under what circumstances the AI localization performance
 217 was closer to, or further from, the expert localization performance, we first conducted a

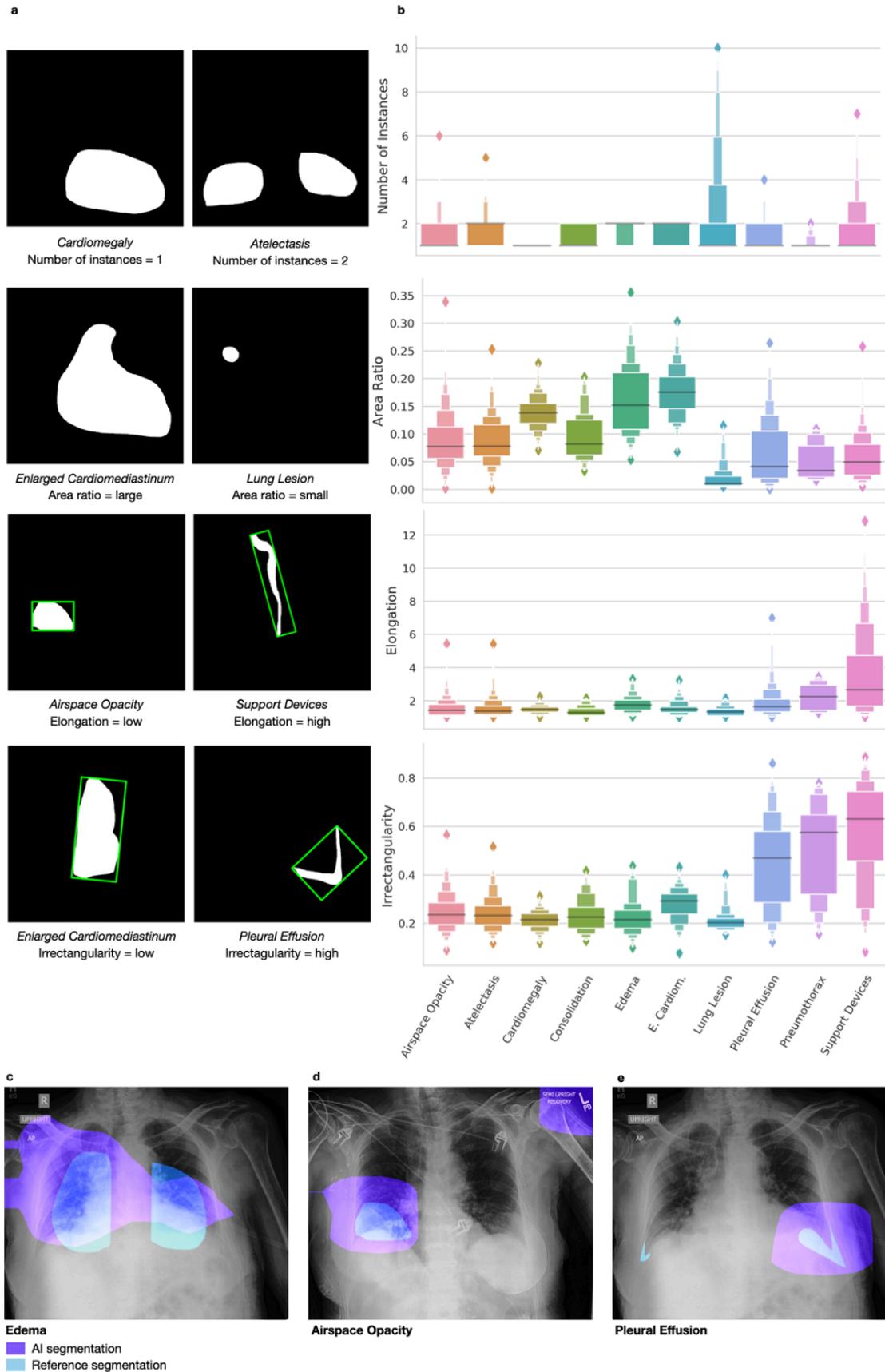
218 qualitative analysis by visually inspecting both the saliency method segmentations and
219 the benchmark radiologist segmentations with a radiologist. Then, to inform our qualitative
220 interpretation, we conducted a statistical analysis to quantify how both the saliency
221 method segmentations and the benchmark radiologist segmentations performed in the
222 presence of four pathological characteristics³¹: (1) number of instances for a given
223 pathology (for example, bilateral Pleural Effusion would have two instances, whereas
224 there is usually only one instance for Cardiomegaly), (2) area ratio (pathology area with
225 respect to the area of the whole CXR), (3) elongation, and (4) irrectangularity (the last two
226 features were meant to measure the complexity of the pathology's shape). See Fig. 3a
227 for example segmentations with the above 4 characteristics. See Fig. 3b for the
228 distribution of the four pathological characteristics across all 10 pathologies. Lung Lesion
229 had the largest number of instances on average, and the lowest mean area ratio among
230 all the pathologies. Support Devices and Pneumothorax had the two highest mean
231 elongation and irrectangularity values, respectively.

232
233 For each evaluation scheme (overlap and hit rate), we ran 12 simple linear regressions:
234 four with the AI evaluation metric (IoU or hit rate) as the response variable, four with the
235 expert evaluation metric as the response variable, and four with the difference between
236 the expert and AI evaluation metrics as the response variable. Each group of four
237 regressions used the above four pathological characteristics as the regression's single
238 attribute, respectively, and only CXRs with a positive label were included in each
239 regression (n=1534). Each regression coefficient can be interpreted as the effect of that

240 pathological characteristic on the evaluation metric at hand. See Table 2 for coefficients
241 from the regressions under the overlap and hit rate evaluation schemes.

242

243 Our qualitative analysis uncovered three patterns in the saliency method segmentations
244 that were associated with lower localization performance. First, we observed that when
245 multiple instances of a single pathology are present in a CXR, instead of highlighting each
246 distinct instance of the pathology separately, the saliency method segmentation often
247 highlights one large confluent area that encompasses all of the instances (see Fig. 3c).
248 Second, we found that saliency method segmentations tend to be significantly larger than
249 either the human benchmark or reference segmentations, and often fail to respect clear
250 anatomical boundaries (see Fig. 3d). Correspondingly, the AI overlap coefficient for area
251 ratio was 0.556 (95% CI [0.510, 0.601]), suggesting that as a pathology's area ratio
252 decreases, AI localization performance worsens. Furthermore, under the overlap
253 evaluation scheme, the gap between AI localization performance and expert localization
254 performance increases as a pathology's area ratio decreases: the area ratio coefficient
255 was -0.151 (95% CI [-0.233, -0.07]) when the difference between expert and AI overlap
256 was the response variable. Third, our qualitative analysis showed that, in segmenting
257 complex and elongated pathologies, when the AI segmentations include the pathology,
258 they also frequently enclose significant portions of the CXR where the pathology is not
259 present (see Fig. 3e). Similarly, our statistical analysis demonstrated that AI localization
260 performance worsens both as a pathology's elongation increases (overlap coefficient = -
261 0.375 95% CI [-0.453, -0.297]), and as a pathology's irrectangularity increases (overlap
262 coefficient = -0.205 95% CI [-0.245, -0.165]).



264 **Fig. 3 | Characterizing the gaps between AI localization performance and expert**
265 **localization performance.** **a**, Example segmentations with four pathological
266 characteristics: (1) number of instances (top row), (2) area ratio (second row), (3)
267 elongation (third row), and (4) irrectangularity (fourth row). Elongation and irrectangularity
268 were calculated by fitting a rectangle of minimum area enclosing the binary mask.
269 Elongation = $\text{maxLength}/\text{minLength}$. Irrectangularity = $1 - (\text{area of segmentation}/\text{area of}$
270 $\text{enclosing rectangle})$. **b**, Distribution of the four pathological characteristics across all 10
271 pathologies in letter value plot style. The black horizontal line in each box indicates the
272 median for that pathology, and from the middle each box represents the increasing
273 (middle to up) and decreasing (middle to bottom) quantile. **c**, Multiple instances of Edema
274 in this CXR are shown by the reference segmentations in blue. Instead of highlighting
275 each distinct instance of the pathology separately, the AI segmentation in purple
276 highlights one large confluent area that tries to encompass both instances. **d**, Airspace
277 Opacity in this CXR is shown by the reference segmentation in blue. Not only is the AI
278 segmentation larger than the Airspace Opacity, but it also fails to respect clear anatomical
279 boundaries by highlighting an area outside of the chest cavity. **e**, Bilateral Pleural Effusion
280 in this CXR is shown by the reference segmentations in blue. On the right, instead of
281 highlighting the distinct V-shape of the pathology, the AI segmentation highlights a large
282 portion of the CXR where the pathology is not present in trying to enclose the whole
283 pathology.

284

285 Our statistical analysis showed that the first and second trends listed above held not only
286 for the AI segmentations, but also for the expert segmentations: as the number of
287 instances of a pathology increases, expert localization performance worsens (overlap
288 coefficient = -0.178 95% CI [$-0.334, -0.021$]), and as the area ratio of a pathology
289 increases, expert localization performance improves (overlap coefficient = 0.404 95% CI
290 [$0.334, 0.475$]). However, our statistical analysis showed no evidence that the third trend
291 holds for the expert segmentations: expert localization performance does not worsen as

292 pathology complexity increases (irrectangularity overlap coefficient = 0.073 [0.016, 0.13];
 293 elongation overlap coefficient was not statistically significant).

294

295 The results of the above experiments using hit/miss as an evaluation metric were
 296 consistent with the results when using overlap as an evaluation metric. In both
 297 experiments, we reported the 95% confidence interval and the Bonferroni corrected p-
 298 values.

Table 2 | Coefficients from regressions on pathological characteristics

Pathological characteristic (independent variable)	Coefficient using AI localization as response	Coefficient using expert localization as response	Coefficient using localization difference (expert - AI)
IoU			
Number of instances	0.003 (-0.111, 0.116)	-0.178 (-0.334, -0.021) *	-0.18 (-0.355, -0.006) *
Area ratio	0.556 (0.510, 0.601) ***	0.404 (0.334, 0.475) ***	-0.151 (-0.233, -0.07) ***
Elongation	-0.375 (-0.453, -0.297) ***	0.067 (-0.043, 0.178)	0.442 (0.321, 0.563) ***
Irrectangularity	-0.205 (-0.245, -0.165) ***	0.073 (0.016, 0.13) **	0.278 (0.216, 0.341) ***
Hit/Miss			
Number of instances	-0.101 (-1.307, -0.714) ***	0.403 (0.175, 0.632) **	1.413 (1.059, 1.768) ***
Area ratio	1.111 (0.982, 1.24) ***	0.313 (0.207, 0.42) ***	-0.78 (-0.962, -0.633) ***
Elongation	-0.842 (-1.05, -0.634) ***	0.267 (0.105, 0.428) **	1.108 (0.86, 1.357) ***
Irrectangularity	-0.536 (-0.642, -0.429) ***	0.169 (0.0852, 0.252) ***	0.704 (0.577, 0.832) ***

* p-value < 0.05, ** p-value < 0.01, *** p-value < 0.001

299

300

301 **Effect of model confidence on AI localization performance**

302 Since the saliency method is highly dependent on the DNN's architecture, we conducted
 303 statistical analyses to determine whether there was any correlation between the model's
 304 confidence in its prediction and AI localization performance. We first ran a simple

305 regression for each pathology using the model's probability output for the pathology as
306 the single independent variable and using the IoU of the AI segmentation with reference
307 segmentation as the response variable. We then performed a simple regression that uses
308 the same approach as above, but that includes all 10 pathologies. For each of the 11
309 regressions, we excluded true negative cases in order to calculate the IoU score for the
310 expert segmentations. In addition to the linear regression coefficients, we also computed
311 the Spearman correlation coefficients to capture any potential non-linear associations
312 (see Table 3).

313
314 We found that for all pathologies, model confidence was positively correlated with AI
315 localization performance. The p-values for all of the coefficients were below 0.001 except
316 for the coefficients for Pneumothorax (n=11) and Lung Lesion (n=50), the two pathologies
317 for which we had the fewest positive examples. Of all the pathologies, model confidence
318 for positive predictions of Enlarged Cardiomediatinum had the largest linear regression
319 coefficient with AI localization performance (1.974, p-value = 2.523e-24). Model
320 confidence for positive predictions of Pneumothorax had the largest Spearman correlation
321 coefficient with AI localization performance (0.734, p-value = 0.01), although the
322 coefficient was not as statistically significant as the Spearman correlation coefficient for
323 Pleural Effusion (0.69, p-value=8.08e-24), the second largest of all the pathologies.
324 Combining all of the pathologies (n=2365), the linear regression coefficient was 0.109
325 (95% CI [0.083, 0.135]), and the Spearman correlation coefficient was 0.285 (95%CI
326 [0.239, 0.331]). We also performed analogous experiments using hit rate as the response
327 variable and found comparable results (see Supplementary Table 1).

Pathology	CXRs (n)	Linear regression coefficient	Spearman correlation coefficient
Airspace Opacity	381	0.714 (0.601, 0.826) ***	0.577 (0.542, 0.61) ***
Atelectasis	296	0.489 (0.333, 0.645) ***	0.348 (0.303, 0.391) ***
Cardiomegaly	229	0.679 (0.535, 0.823) ***	0.592 (0.559, 0.624) ***
Consolidation	120	1.155 (0.674, 1.635) ***	0.384 (0.341, 0.426) ***
Edema	124	0.642 (0.459, 0.826) ***	0.548 (0.512, 0.582) ***
Enlarged Cardiomeastinum	668	1.974 (1.608, 2.340) ***	0.428 (0.386, 0.468) ***
Lung Lesion	50	0.218 (0.087, 0.349) **	0.509 (0.47, 0.545) ***
Pleural Effusion	159	0.632 (0.489, 0.776) ***	0.69 (0.663, 0.715) ***
Pneumothorax	11	0.446 (0.108, 0.783) *	0.734 (0.710, 0.756) **
Support Devices	327	0.211 (0.172, 0.25) ***	0.468 (0.428, 0.506) ***
All pathologies	2365	0.109 (0.083, 0.135) ***	0.285 (0.239, 0.331) ***

* p-value < 0.05, ** p-value < 0.01, *** p-value < 0.001

328

329

330 Discussion

331 The purpose of this work was to evaluate the performance of saliency methods, which
332 are widely used in clinical practice for DNN prediction explainability. We demonstrate that
333 saliency maps are consistently worse than expert radiologists at localizing a variety of
334 pathologies on CXRs. We use qualitative and quantitative analyses to establish that AI
335 localization performance is furthest from expert localization performance in the face of
336 pathologies that have multiple instances, are smaller in size, and have shapes that are
337 more complex, suggesting that deep learning explainability as a clinical interface may be
338 less reliable and less useful when used for pathologies with those characteristics. We
339 show that model assurance is positively correlated with AI localization performance,

340 which could indicate that the saliency methods are safer to use as a decision aid to
341 clinicians when the model has made a positive prediction with high confidence. Finally,
342 since IoU computes the overlap of two segmentations but pointing game hit rate better
343 captures diagnostic attention, we suggest using both metrics to evaluate both AI and
344 expert localization performance.

345

346 Our work has several potential implications for patient care. Heat maps generated using
347 saliency methods are advocated as clinical decision support in the hope that the heat
348 maps not only improve clinical decision-making, but also encourage clinicians to trust
349 model predictions^{32–34}. However, we found that AI localization performance, on balance,
350 performed worse than expert localization across multiple analyses. This is consistent with
351 recent work focused on localizing a single pathology, Pneumothorax, in CXRs³⁵. Our work
352 expanded this exploration and found that common saliency methods may underperform
353 for many important pathologies on CXRs. If used in clinical practice, heat maps that
354 incorrectly highlight medical images may exacerbate well documented biases, chiefly
355 automation bias, and erode trust in model predictions (even when model output is
356 correct), limiting clinical translation²⁵.

357

358 Determining when saliency methods are more likely to succeed or fail in localizing
359 pathologies could have further implications for patient care. That knowledge could inform
360 not only under what clinical conditions saliency methods might be safer to use, but also
361 how we might improve saliency methods in the future. We found that AI localization
362 performance worsens in the presence of pathologies that have multiple instances. We

363 also found that AI localization performance worsens in the presence of pathologies that
364 are smaller in size compared with the CXR image. This result explains why, under both
365 the overlap and hit rate schemes, the gap between AI and expert localization performance
366 was the largest for Lung Lesion, whose mean area ratio was the smallest of the 10
367 pathologies we explored. Moreover, AI localization performance under both evaluation
368 schemes was best for Enlarged Cardiomeastinum and Cardiomegaly, which had the
369 largest and third largest area ratios, respectively, of the 10 pathologies, suggesting that
370 saliency methods might be safer to use in the context of these two pathologies, or
371 pathologies with similar characteristics. Grad-CAM segmentations often fail to respect
372 clear anatomical boundaries, and we hypothesize that this is an algorithmic artifact of
373 Grad-CAM, whose feature map sized (14 x 14) heatmap is interpolated to the original
374 image dimension (usually 2000 x 2000), resulting in coarse resolution. We also found that
375 AI localization performance worsens in the presence of pathologies whose shapes are
376 more complex. AI localization for Pneumothorax and Support Devices, both of which were
377 more elongated and complex than any of the other conditions, underperformed compared
378 to expert localization performance; however, this performance gap must also be
379 considered in the context of the model training data prevalence, and future work may
380 explore the impact of training data prevalence on the localization performance of saliency
381 methods.

382

383 While IoU is a commonly used metric for evaluating semantic segmentation outputs, there
384 are inherent limitations to the metric in the pathological context. This is indicated by our
385 finding that even the expert segmentations had relatively low overlap with the reference

386 segmentations (the highest expert mIoU was 0.712 for Cardiomegaly). One potential
387 explanation for this consistent underperformance is that pathologies can be hard to
388 distinguish, especially without clinical context. Furthermore, whereas many people might
389 agree on how to segment, say, a cat or a stop sign in traditional computer vision tasks,
390 radiologists use a certain amount of clinical discretion when defining the boundaries of a
391 pathology on a CXR. There can also be institutional and geographic differences in how
392 radiologists are taught to recognize pathologies, and studies have shown that there can
393 be high interobserver variability in the interpretation of CXRs^{36–38}. We sought to address
394 this with the hit rate evaluation metric, which highlights when two radiologists share the
395 same diagnostic intention, even if it is less exact than IoU in comparing segmentations
396 directly. Expert performance using hit rate was above 0.95 for four pathologies
397 (Pneumothorax, Cardiomegaly, Support Devices, and Enlarged Cardiomediatinum);
398 these are pathologies for which there is often little disagreement between radiologists
399 about where the pathologies are located, even if the expert segmentations are noisy. The
400 only pathology for which AI localization performance was better than expert localization
401 performance under the hit rate scheme was Consolidation. However, because the hit rate
402 scheme required the benchmark radiologists to select only one point on the CXR, even if
403 there were multiple instances of the pathology present (as is often the case with
404 Consolidation), it is likely that the hit rate setup unfairly penalized expert performance in
405 this case and that it is not the best evaluation metric to use for Consolidation. Further
406 work is needed to validate this hypothesis and to demonstrate which segmentation
407 evaluation metrics, even beyond overlap and hit rate, are more appropriate for which
408 pathologies when evaluating saliency methods for the clinical setting.

409

410 Our work builds upon several studies investigating the validity of saliency maps in
411 localization^{39,40} and upon some early work on trustworthiness of saliency methods to
412 explain DNNs in medical imaging⁴¹. We substantially extend the body of literature by
413 doing a comprehensive analysis on a multi-label classification task using the most popular
414 saliency method in the medical context, Grad-CAM. Our work analyzes 10 different
415 commonly occurring clinical pathologies as opposed to only Pneumothorax and
416 Pneumonia. In addition to demonstrating that Grad-CAM is not yet ready for clinical use,
417 we also highlight the strengths and weaknesses of Grad-CAM by doing quantitative
418 statistical analysis, thus opening future avenues of research to improve Grad-CAM in
419 particular and saliency methods in general. Finally, we establish new ground truth and
420 benchmark segmentations on 10 different CXR observations facilitating future research
421 on attribution methods.

422

423 There are several limitations of our work. We chose to evaluate Grad-CAM, since it has
424 become one of the most popular explainability methods for CXRs, but future work may
425 also evaluate other saliency methods, including Integrated Gradients⁴², WILDCAT⁴³, and
426 SmoothGrad⁴⁴. While we used DenseNet121 as the underlying model architecture for the
427 saliency method, since it was shown to produce the best classification results on the
428 CheXpert dataset, future work should explore and compare different model architectures
429 for CXR explainability. Our dataset had only 11 CXRs with Pneumothorax and 50 CXRs
430 with Lung Lesion. Future work should investigate the impact of pathology prevalence in
431 the training data on AI localization performance. Some pathologies, such as effusions

432 and cardiomegaly, are always in the same place in a frontal view CXR; others, such as
433 lesions and opacities, can occur in different locations on a CXR. Future work could
434 investigate how a pathology's location on a CXR, and the consistency of that location,
435 affect AI localization performance. Finally, we compared Grad-CAM-generated pixel-level
436 segmentations to human expert pixel-level segmentations, but future work might explore
437 how AI localization performance changes when comparing bounding-box annotations,
438 instead of pixel-level segmentations.

439
440 In conclusion, we demonstrate that not only should we not yet rely on saliency methods
441 for deep learning explainability in CXRs, but saliency methods are particularly brittle in
442 the presence of pathologies that have multiple instances, are smaller in size, and have
443 shapes that are more complex. Although our findings suggest that care should be taken
444 when deploying saliency methods into clinical practice, we demonstrate scenarios in
445 which heat maps generated by saliency methods are most consistent with human expert
446 annotations. Our work serves as the foundation for future work that rigorously evaluates
447 a range of saliency methods using a variety of evaluation metrics before deep learning
448 explainability techniques are integrated into the medical workflow.

449

450 **Methods**

451 **Ethical and information governance approvals.**

452 This study does not involve human subject participants.

453

454 **Dataset and clinical taxonomy.** *Dataset description.* The localization experiments were
455 performed using CheXpert, a large public dataset for chest X-ray interpretation. The
456 CheXpert dataset contains 224,316 chest X-rays for 65,240 patients labeled for the
457 presence of 14 observations (13 pathologies and an observation of “No Finding”) as
458 positive, negative, or uncertain. The CheXpert validation set consists of 234 chest X-rays
459 from 200 patients randomly sampled from the full dataset and was labeled according to
460 the consensus of three board-certified radiologists. The test set consists of 668 chest X-
461 rays from 500 patients not included in the training or validation sets and was labeled
462 according to the consensus of five board-certified radiologists. See Supplementary Table
463 2 for dataset summary statistics.

464
465 *Reference segmentation.* The chest X-rays in our validation set and test set were
466 manually segmented by two board-certified radiologists with 18 and 27 years of
467 experience, using the annotation software tool MD.ai⁴⁵ (see Supplementary Figs. 12
468 through 14). The radiologists were asked to contour the region of interest for all
469 observations in the chest X-rays for which there was a positive ground truth label in the
470 CheXpert dataset. For a pathology with multiple instances, all the instances were
471 contoured. For Support Devices, radiologists were asked to contour any implanted or
472 invasive devices including pacemakers, PICC/central catheters, chest tubes,
473 endotracheal tubes, feeding tubes and stents and ignore ECG lead wires or external
474 stickers visible in the chest X-ray. Finally, of the 14 observations labeled in the CheXpert
475 dataset, Fracture, Pleural Other, and Pneumonia were not segmented because they
476 either had low prevalence and/or ill-defined boundaries unfit for segmentation.

477

478 *Evaluating the expert performance using benchmark segmentation.* To evaluate the
479 expert performance on the test set using the IoU evaluation method, three radiologists,
480 certified in Vietnam with 9, 10, and 18 years of experience, were asked to segment the
481 regions of interest for all observations in the chest X-rays for which there was a positive
482 ground truth label in the CheXpert dataset. These radiologists were also provided the
483 same instructions for contouring as were provided to the radiologists drawing the
484 reference segmentations. To extract the “maximally activated” point from the benchmark
485 segmentations, we asked the same radiologists to locate each pathology present on each
486 CXR using only a single most representative point for that pathology on the CXR (see
487 Supplementary Figs. 1 through 11 for the detailed instructions given to the radiologists).
488 There was no overlap between these three radiologists and the two who drew the
489 reference segmentations.

490

491 **Classification network architecture and training protocol.** *Multi-label classification*
492 *model.* The model takes as input a single-view chest X-ray and outputs the probability for
493 each of the 14 observations. In case of availability of more than one view, the models
494 output the maximum probability of the observations across the views. Each chest X-ray
495 was resized to 320×320 pixels and normalized before it was fed into the network. The
496 DenseNet121 model architecture⁴⁶ was used. Cross-entropy loss was used to train the
497 model. The Adam optimizer⁴⁷ was used with default β -parameters of $\beta_1 = 0.9$ and $\beta_2 =$
498 0.999 , and the learning rate was fixed at 1×10^{-4} for the duration of the training. Batches
499 were sampled using a fixed batch size of 16 images.

500

501 *Ensembling.* An ensemble of 30 DenseNet121 checkpoints was created to improve the
502 performance of the model. The 30 checkpoints were generated by training the model for
503 3 epochs and selecting the 10 checkpoints from each epoch with the highest average
504 AUC across 5 observations selected for their clinical importance and prevalence in the
505 validation set: Atelectasis, Cardiomegaly, Consolidation, Edema, and Pleural Effusion.
506 See Supplementary Table 3 for the performance of the model on each of the pathologies.

507

508 **DNN interpretation strategy.** *Saliency method.* Grad-CAM was used to visualize the
509 decision made by the classification network. Grad-CAM uses the gradients of the target
510 flowing into the final convolutional layer to produce a saliency map that highlights the
511 regions on which the model focuses while making the decision. The saliency map
512 outputted by Grad-CAM was resized to the original image dimension. It was then
513 normalized using max-min normalization and then converted into a binary segmentation
514 using binary thresholding (Otsu's method⁴⁸). To further ensure that the final binary
515 segmentation is consistent with model probability output, another layer of thresholding
516 was applied such that the segmentation mask produced all zeros if the predicted
517 probability was below a chosen level. The probability threshold is searched on the interval
518 of [0,0.8] with steps of 0.1. The exact value is determined per pathology by maximizing
519 the mIoU on validation set.

520

521 *Segmentation evaluation metrics.* Localization performance of each segmentation was
522 evaluated using Intersection over Union (IoU) score. The (Intersection over Union) IoU is

523 the ratio between the area of overlap and the area of union between the ground truth and
524 the predicted areas, ranging from 0–1 with 0 signifying no overlap and 1 signifying
525 perfectly overlapping segmentation. We then compared the mean Intersection over Union
526 (mIoU) of Grad-CAM and radiologist benchmark on each pathology. The mIoU is the
527 average IoU of all the images in the test dataset. True negatives where both
528 segmentations were labeled as all 0s are excluded in the mean calculation. Confidence
529 intervals are calculated using bootstrapping with 1000 bootstrap samples. The variance
530 in the width of CI across pathologies can be explained by difference in sample sizes.

531

532 **Statistical analysis.**

533 *Pathology Characteristics.* We used four features to characterize the pathologies. 1.
534 Number of instances is defined as the number of disjoint components in the
535 segmentation. 2. Area ratio area is the area of the pathology divided by the total image
536 area. 3.4. Elongation and irrectangularity are geometric features that measure shape
537 complexities. They were designed to quantify what radiologists qualitatively described as
538 focal or diffused. To calculate the metrics, a rectangle of minimum area enclosing the
539 contour is fitted to each pathology. Elongation is defined as the ratio of the rectangle's
540 longer side to short side. Irrectangularity = $1 - \frac{\text{area of segmentation}}{\text{area of enclosing}}$
541 rectangle, with values ranging from 0 to 1 with 1 being very irrectangular. When there are
542 multiple instances within one pathology, we used the characteristics of the dominant
543 instance (largest in perimeter).

544

545 *Model Confidence.* We used the probability output of the DNN architecture for model
546 confidence. The probabilities were normalized using max-min normalization per
547 pathology before aggregation.

548

549 *Linear Regression.* For each evaluation scheme (overlap and hit rate), we ran three
550 groups of simple linear regressions, with expert and AI evaluation metrics and their
551 differences as the response variables. Each group has four regressions using the above
552 four pathological characteristics as the regression's single attribute, respectively, and only
553 CXRs with a positive label were included in each regression (n=1534). All features are
554 normalized using min-max normalization so that they are comparable on scales of
555 magnitudes. We report the 95% confidence interval and p-value of the regression
556 coefficients.

557

558 **Data Availability**

559 CheXpert data is available at <https://stanfordmlgroup.github.io/competitions/chexpert/>.

560 The validation set and corresponding benchmark radiologist annotations will be
561 available online for the purpose of extending the study.

562

563 **Code Availability**

564 All code used to produce the results of the paper will be in a public repository for the
565 purpose of reproducing the study. The link to the code will be added to the text of the
566 paper for the camera-ready version.

567

568 **Competing Interests**

569 There are no competing interests.

570

571 **Author Contributions**

572 (1) Conceptualization: P.R. and A.P., (2) Design: P.R., A.P., A.S., X.G. and A.A., (3) Data
573 analysis and interpretation: A.S., X.G., A.A., P.R., A.P., S.T., C.N., V.N., J.S., and F.B.,
574 (4) Drafting of the manuscript: A.S., X.G., A.A., and P.R., (5) Critical revision of the
575 manuscript for important intellectual content: A.P, S.T., C.N., V.N., J.S., F.B, A.N., and
576 M.L., (6) Supervision: A.N., M.L, and P.R.

577

578 **References**

- 579 1. Rajpurkar, P. *et al.* CheXNet: Radiologist-Level Pneumonia Detection on Chest X-
580 Rays with Deep Learning. *ArXiv171105225 Cs Stat* (2017).
- 581 2. Baselli, G., Codari, M. & Sardanelli, F. Opening the black box of machine learning in
582 radiology: can the proximity of annotated cases be a way? *Eur. Radiol. Exp.* **4**, 30
583 (2020).
- 584 3. Litjens, G. *et al.* A survey on deep learning in medical image analysis. *Med. Image*
585 *Anal.* **42**, 60–88 (2017).
- 586 4. Wang, F., Kaushal, R. & Khullar, D. Should Health Care Demand Interpretable
587 Artificial Intelligence or Accept “Black Box” Medicine? *Ann. Intern. Med.* **172**, 59–60
588 (2019).

- 589 5. Reyes, M., Meier, R., Pereira, S., Silva, C. A. & Dahlweid, F.-M. On the
590 Interpretability of Artificial Intelligence in Radiology: Challenges and Opportunities. **2**,
591 12.
- 592 6. Pasa, F., Golkov, V., Pfeiffer, F., Cremers, D. & Pfeiffer, D. Efficient Deep Network
593 Architectures for Fast Chest X-Ray Tuberculosis Screening and Visualization. *Sci.*
594 *Rep.* **9**, 6268 (2019).
- 595 7. Simonyan, K., Vedaldi, A. & Zisserman, A. Deep Inside Convolutional Networks:
596 Visualising Image Classification Models and Saliency Maps. *ArXiv13126034 Cs*
597 (2014).
- 598 8. Aggarwal, M. *et al.* Towards Trainable Saliency Maps in Medical Imaging.
599 *ArXiv201107482 Cs Eess* (2020).
- 600 9. Tjoa, E. & Guan, C. Quantifying Explainability of Saliency Methods in Deep Neural
601 Networks. *ArXiv200902899 Cs* (2020).
- 602 10. Badgeley, M. A. *et al.* Deep learning predicts hip fracture using confounding
603 patient and healthcare variables. *Npj Digit. Med.* **2**, 1–10 (2019).
- 604 11. Zech, J. R. *et al.* Variable generalization performance of a deep learning model
605 to detect pneumonia in chest radiographs: A cross-sectional study. *PLOS Med.* **15**,
606 e1002683 (2018).
- 607 12. DeGrave, A. J., Janizek, J. D. & Lee, S.-I. AI for radiographic COVID-19
608 detection selects shortcuts over signal. *medRxiv* (2020)
609 doi:10.1101/2020.09.13.20193565.
- 610 13. Rudin, C. Stop explaining black box machine learning models for high stakes
611 decisions and use interpretable models instead. *Nat. Mach. Intell.* **1**, 206–215 (2019).

- 612 14. Selvaraju, R. R. *et al.* Grad-CAM: Visual Explanations from Deep Networks via
613 Gradient-based Localization. *Int. J. Comput. Vis.* **128**, 336–359 (2020).
- 614 15. Rizwan I Haque, I. & Neubert, J. Deep learning approaches to biomedical image
615 segmentation. *Inform. Med. Unlocked* **18**, 100297 (2020).
- 616 16. Makimoto, H. *et al.* Performance of a convolutional neural network derived from
617 an ECG database in recognizing myocardial infarction. *Sci. Rep.* **10**, 8445 (2020).
- 618 17. Raghunath, S. *et al.* Prediction of mortality from 12-lead electrocardiogram
619 voltage data using a deep neural network. *Nat. Med.* **26**, 886–891 (2020).
- 620 18. Porumb, M., Stranges, S., Pescapè, A. & Pecchia, L. Precision Medicine and
621 Artificial Intelligence: A Pilot Study on Deep Learning for Hypoglycemic Events
622 Detection based on ECG. *Sci. Rep.* **10**, 1–16 (2020).
- 623 19. Tham, Y.-C. *et al.* Referral for disease-related visual impairment using retinal
624 photograph-based deep learning: a proof-of-concept, model development study.
625 *Lancet Digit. Health* **3**, e29–e40 (2021).
- 626 20. Varadarajan, A. V. *et al.* Deep Learning for Predicting Refractive Error From
627 Retinal Fundus Images. *Invest. Ophthalmol. Vis. Sci.* **59**, 2861–2868 (2018).
- 628 21. Mitani, A. *et al.* Detection of anaemia from retinal fundus images via deep
629 learning. *Nat. Biomed. Eng.* **4**, 18–27 (2020).
- 630 22. Deep Learning to Assess Long-term Mortality From Chest Radiographs |
631 Pulmonary Medicine | JAMA Network Open | JAMA Network. [https://jamanetwork-](https://jamanetwork-com.stanford.idm.oclc.org/journals/jamanetworkopen/fullarticle/2738349)
632 [com.stanford.idm.oclc.org/journals/jamanetworkopen/fullarticle/2738349](https://jamanetwork-com.stanford.idm.oclc.org/journals/jamanetworkopen/fullarticle/2738349).
- 633 23. Rajpurkar, P. *et al.* CheXaid: deep learning assistance for physician diagnosis of
634 tuberculosis using chest x-rays in patients with HIV. *Npj Digit. Med.* **3**, 1–8 (2020).

- 635 24. Tschandl, P. *et al.* Human–computer collaboration for skin cancer recognition.
636 *Nat. Med.* **26**, 1229–1234 (2020).
- 637 25. Rajpurkar, P. *et al.* AppendiXNet: Deep Learning for Diagnosis of Appendicitis
638 from A Small Dataset of CT Exams Using Video Pretraining. *Sci. Rep.* **10**, 3958
639 (2020).
- 640 26. Bien, N. *et al.* Deep-learning-assisted diagnosis for knee magnetic resonance
641 imaging: Development and retrospective validation of MRNet. *PLOS Med.* **15**,
642 e1002699 (2018).
- 643 27. Irvin, J. *et al.* CheXpert: A Large Chest Radiograph Dataset with Uncertainty
644 Labels and Expert Comparison. *Proc. AAAI Conf. Artif. Intell.* **33**, 590–597 (2019).
- 645 28. Zhang, J., Lin, Z., Brandt, J., Shen, X. & Sclaroff, S. Top-down Neural Attention
646 by Excitation Backprop. *ArXiv160800507 Cs* (2016).
- 647 29. Kim, H.-E. *et al.* Changes in cancer detection and false-positive recall in
648 mammography using artificial intelligence: a retrospective, multireader study. *Lancet*
649 *Digit. Health* **2**, e138–e148 (2020).
- 650 30. Vrabac, D. *et al.* DLBCL-Morph: Morphological features computed using deep
651 learning for an annotated digital DLBCL image set. *ArXiv200908123 Cs* (2020).
- 652 31. Steiner, D. F. *et al.* Impact of Deep Learning Assistance on the Histopathologic
653 Review of Lymph Nodes for Metastatic Breast Cancer. *Am. J. Surg. Pathol.* **42**,
654 1636–1646 (2018).
- 655 32. Uyumazturk, B. *et al.* Deep Learning for the Digital Pathologic Diagnosis of
656 Cholangiocarcinoma and Hepatocellular Carcinoma: Evaluating the Impact of a Web-
657 based Diagnostic Assistant. *ArXiv191107372 Eess* (2019).

- 658 33. Park, A. *et al.* Deep Learning–Assisted Diagnosis of Cerebral Aneurysms Using
659 the HeadXNet Model. *JAMA Netw. Open* **2**, e195600 (2019).
- 660 34. Crosby, J., Chen, S., Li, F., MacMahon, H. & Giger, M. Network output
661 visualization to uncover limitations of deep learning detection of pneumothorax. in
662 *Medical Imaging 2020: Image Perception, Observer Performance, and Technology*
663 *Assessment* vol. 11316 113160O (International Society for Optics and Photonics,
664 2020).
- 665 35. Melbye, H. & Dale, K. Interobserver Variability in the Radiographic Diagnosis of
666 Adult Outpatient Pneumonia. *Acta Radiol.* **33**, 79–81 (1992).
- 667 36. Herman, P. G. *et al.* Disagreements in Chest Roentgen Interpretation. *CHEST*
668 **68**, 278–282 (1975).
- 669 37. Albaum, M. N. *et al.* Interobserver Reliability of the Chest Radiograph in
670 Community-Acquired Pneumonia. *CHEST* **110**, 343–350 (1996).
- 671 38. Arun, N. T. *et al.* Assessing the validity of saliency maps for abnormality
672 localization in medical imaging. *ArXiv200600063 Cs* (2020).
- 673 39. Graziani, M., Lompech, T., Müller, H. & Andrearczyk, V. *Evaluation and*
674 *Comparison of CNN Visual Explanations for Histopathology.* (2020).
- 675 40. Arun, N. *et al.* Assessing the (Un)Trustworthiness of Saliency Maps for Localizing
676 Abnormalities in Medical Imaging. *ArXiv200802766 Cs* (2020).
- 677 41. Sundararajan, M., Taly, A. & Yan, Q. Axiomatic Attribution for Deep Networks.
678 *ArXiv170301365 Cs* (2017).
- 679 42. Durand, T., Mordan, T., Thome, N. & Cord, M. WILDCAT: Weakly Supervised
680 Learning of Deep ConvNets for Image Classification, Pointwise Localization and

- 681 Segmentation. in *2017 IEEE Conference on Computer Vision and Pattern*
682 *Recognition (CVPR)* 5957–5966 (IEEE, 2017). doi:10.1109/CVPR.2017.631.
- 683 43. Smilkov, D., Thorat, N., Kim, B., Viégas, F. & Wattenberg, M. SmoothGrad:
684 removing noise by adding noise. *ArXiv170603825 Cs Stat* (2017).
- 685 44. MD.ai. <https://www.md.ai/>.
- 686 45. Huang, G., Liu, Z., van der Maaten, L. & Weinberger, K. Q. Densely Connected
687 Convolutional Networks. *ArXiv160806993 Cs* (2018).
- 688 46. Kingma, D. P. & Ba, J. Adam: A Method for Stochastic Optimization.
689 *ArXiv14126980 Cs* (2017).
- 690 47. Otsu, N. A Threshold Selection Method from Gray-Level Histograms. *IEEE Trans*
691 *Syst. Man Cybern.* 62–66 (1979).
- 692