

1 Universal adversarial attacks on deep neural networks 2 for medical image classification 3

4 **Hokuto Hirano¹, Akinori Minagi¹, Kazuhiro Takemoto^{1*}**

5 *1) Department of Bioscience and Bioinformatics, Kyushu Institute of Technology, Iizuka,*
6 *Fukuoka 820-8502, Japan*

7 **Corresponding author's e-mail: takemoto@bio.kyutech.ac.jp*
8

9 **Abstract**

10 **Background.** Deep neural networks (DNNs) are widely investigated in medical image
11 classification to achieve automated support for clinical diagnosis. It is necessary to evaluate
12 the robustness of medical DNN tasks against adversarial attacks, as high-stake decision
13 making will be made based on the diagnosis. Several previous studies have considered
14 simple adversarial attacks. However, the vulnerability of DNNs to more realistic and higher
15 risk attacks have not been evaluated yet, i.e., universal adversarial perturbation (UAP),
16 which is a single perturbation that can induce DNN failure in most classification tasks.

17 **Methods.** We focus on three representative DNN-based medical image classification tasks
18 (i.e., skin cancer, referable diabetic retinopathy, and pneumonia classifications) and
19 investigate their vulnerability of DNNs with various model architectures to UAPs.

20 **Results.** We demonstrate that the DNNs are vulnerable to both nontargeted UAPs, which
21 cause a task failure resulting in an input being assigned an incorrect class, and to targeted
22 UAPs, which cause the DNN to classify an input into a specific class. The almost
23 imperceptible UAPs achieved > 80% success rates for nontargeted and targeted attacks.
24 The vulnerability to UAPs barely depended on model architecture. Moreover, we
25 discovered that adversarial retraining, which is known to be an effective method for
26 adversarial defenses, increased the robustness of DNNs against UAPs in only limited cases.

27 **Conclusion.** Unlike previous assumptions, the results indicate that DNN-based clinical
28 diagnosis is easier to deceive because of adversarial attacks. Adversaries can result in failed
29 diagnoses at lower costs (e.g., without consideration of data distribution); moreover, they
30 can affect the diagnosis. The effects of adversarial defenses may be not limited. Our
31 findings emphasize that more careful consideration is required in developing DNNs for
32 medical imaging and their practical applications.

33 **Keywords:** deep neural networks, medical imaging, adversarial attacks, security and
34 privacy

35

36 Background

37 Deep neural networks (DNNs) are effective for image classification and are beginning to
38 be applied to medical image diagnosis to empower physicians and accelerate decision
39 making in clinical environments [1]. For example, DNNs have been used for classifying
40 skin cancer from photographic images [2], referable diabetic retinopathy from optical
41 coherence tomography (OCT) images of the retina [3], and pneumonia from chest X-ray
42 images [3]; they have demonstrated high diagnostic performances. A meta-analysis [4] has
43 indicated that the diagnostic performance of DNNs is equivalent to that of healthcare
44 professionals.

45 Despite the high performance of DNNs, practical applications of DNNs to disease
46 diagnosis are still debatable. High-stake decision-making will be made from disease
47 diagnosis. Complex classifiers, including DNNs, can potentially cause catastrophic harm
48 to society because they are often difficult to interpret [5]. More importantly, DNNs have a
49 number of security concerns [6]; specifically, DNNs are known to be vulnerable to
50 adversarial examples [7, 8], which are input images that cause misclassifications by DNNs
51 and typically generated by adding specific, imperceptible perturbations to original input
52 images that have been correctly classified using DNNs. The existence of adversarial
53 examples questions the generalization ability of DNNs, reduces model interpretability, and
54 limits the applications of deep learning in safety- and security-critical environments [9]; in
55 particular, the adversarial examples cause not only misdiagnosis, but also various social
56 disturbances [10]. The vulnerability of DNNs to adversarial attacks has been claimed in
57 skin cancer classification [10] and pneumonia classification based on chest X-ray images
58 [11].

59 Nevertheless, more focused investigations are required the vulnerability of DNNs to
60 adversarial attacks. Previous studies have only considered input-dependent adversarial
61 attacks (i.e., an individual adversarial perturbation is used such that each input image
62 misclassifies). Such adversarial attacks are difficult tasks because they require high
63 computational costs. More realistic adversarial attacks must be further considered. Notably,
64 a single small perturbation that can induce DNN failure in most image classification tasks
65 has been reported [12]. Such perturbations are called *universal adversarial perturbations*
66 (*UAPs*), as they are image agnostic. A previous study [12] has considered only UAPs for
67 nontargeted attacks, which cause misclassification (i.e., a task failure resulting in an input
68 image being assigned an incorrect class). However, we previously extended the algorithm
69 for generating UAPs to enable targeted attacks [13], which caused the DNN to classify an
70 input image into a specific class. UAPs are difficult to detect as such perturbations are
71 extremely small and hence do not significantly affect data distributions. UAP-based
72 adversarial attacks can be more straightforward to implement by adversaries in real-world
73 environments. UAPs are vulnerable to security threats in medical image diagnosis;
74 however, the vulnerability of UAPs is still poorly evaluated in DNN-based medical image
75 diagnosis to date. In addition, defense strategies against UAPs are still poorly investigated
76 although the vulnerability of DNNs to adversarial attacks indicates the need for strategies
77 to address security concerns (i.e., adversarial defense [8]). Specifically, adversarial
78 retraining is one of the few approaches that could not be defeated thus far [14].

79 The aim of this study to evaluate the vulnerability of DNNs to UAPs for medical image
80 classification. We focused on representative medical image classifications: skin cancer
81 classification from photographic images [2], referable diabetic retinopathy classification
82 from OCT images [3], and pneumonia classification from chest X-ray images [3]. We
83 obtained DNN models with various architectures for medical image diagnosis and
84 investigated the vulnerability to nontargeted and targeted attacks based on UAPs. Moreover,
85 adversarial defense was considered; in particular, we evaluated the increased robustness of
86 DNNs to nontargeted and targeted UAPs using adversarial retraining [12, 14–16], a
87 representative method for adversarial defenses.

88 **Methods**

89 *Medical image datasets*

90 We used three types of medical images: skin lesion images for skin cancer classification,
91 OCT images for referable diabetic retinopathy classification, and chest X-ray images for
92 pneumonia classification.

93 In a previous study [2], skin lesion images (Red-Green-Blue color) were obtained from the
94 International Skin Imaging Collaboration (ISIC) 2018 dataset ([challenge2018.isic-
95 archive.com/task3/training/](http://challenge2018.isic-archive.com/task3/training/)), in which the images were classified into seven classes:
96 melanoma (MEL), melanocytic nevus (NV), basal cell carcinoma (BCC), actinic
97 keratosis/Bowens disease (intraepithelial carcinoma; AKIEC), benign keratosis (solar
98 lentigo/seborrheic keratosis/lichen planus-like keratosis; BKL), dermatofibroma (DF), and
99 vascular lesion (VASC). The dataset comprised 10,015 images. We randomly divided these
100 images into 7,000 training images (778 MEL images, 4689 NV images, 370 BCC images,
101 229 AKIEC images, 764 BKL images, 76 DF images, 94 VASC images) and 3,015 test
102 images (335 MEL images, 2016 NV images, 144 BCC images, 98 AKIEC images, 335
103 BKL images, 39 DF images, 48 VASC images).

104 The OCT images and chest X-ray images (grayscale) were obtained from a previous study
105 [3] (data.mendeley.com/datasets/rscbjbr9sj/3). The OCT images were classified into four
106 classes: choroidal neovascularization with neovascular membrane and associated
107 subretinal fluid (CNV), diabetic macular edema with retinal-thickening-associated
108 intraretinal fluid (DME), multiple drusen present in early age-related macular degeneration
109 (DRUSEN), and normal retina with preserved foveal contour and absence of any retinal
110 fluid/edema (NM). The original dataset comprised 37,455 CNV images, 11,598 DME
111 images, 8,866 DRUSEN images, and 51,390 NM images. We constructed a class-balanced
112 training image set and test image set by randomly selecting 1,960 images and 840 images
113 per class without duplicates, respectively. We finally obtained 7,840 training images and
114 3,360 test images.

115 The chest X-ray images were classified into binary classes: no pneumonia (NORMAL) and
116 viral or bacterial pneumonia (PNEUMONIA). The original dataset comprised 1,583
117 NORMAL images and 4,273 PNEUMONIA images. We constructed a class-balanced
118 training image set and test image set by randomly selecting 900 and 270 images per class
119 without duplicates, respectively. We finally obtained 1,800 training images and 540 test
120 images.

121 *Transfer learning methods*

122 Based on previous studies [2, 3], we obtained the DNN models using transfer learning; in
123 particular, we fine-tuned DNN models pretrained using the ImageNet dataset [17] with a
124 medical image dataset. We mainly considered the Inception V3 architecture [18], following
125 previous studies. To investigate the effect of model architecture on adversarial robustness,
126 we considered different model architectures: VGG16 [19], VGG19 [19], ResNet50 [20],
127 Inception ResNet V2 [21], DenseNet 121 [22], and DenseNet 169 [22]. For each model
128 architecture, we replaced the original last fully connected (FC) layer with a new FC layer
129 in which the output size is the number of classes. The images were resized to 299×299
130 pixels. All layer parameters were fine-tuned using the training images in a medical image
131 dataset. We used the stochastic gradient descent optimizer with a momentum of 0.9. The
132 batch size and number of epochs were set to 32 and 50, respectively. The learning rates
133 were scheduled based on the number of epochs: 0.001 for ≤ 40 epochs, $1e-4$ for 41–45
134 epochs, and $1e-5$ for > 45 epochs. To avoid overfitting, data augmentation was considered:
135 random image rotations with the angle ranging between -5° and 5° and random image
136 shifts with 5% of the height and 5% of width. For skin cancer classification, we adapted
137 oversampling to account for imbalances in the dataset. The transfer learning procedures
138 were performed using Keras (version 2.2.4; keras.io).

139 *Universal adversarial perturbations*

140 Simple iterative algorithms [12, 13] were used to generate UAPs for nontargeted and
141 targeted attacks. The details of the algorithms are described in [12, 13]. We used the
142 nontargeted UAP algorithm available in the Adversarial Robustness 360 Toolbox (ART)
143 [23] (version 1.0; github.com/Trusted-AI/adversarial-robustness-toolbox). The targeted
144 UAP algorithm was implemented by modifying the nontargeted UAP algorithm in the ART
145 in our previous study [13] (github.com/hkthirano/targeted_UAP_CIFAR10).

146 The algorithms consider a classifier and generate nontargeted (targeted) UPAs $\boldsymbol{\rho}$ from an
147 input image set \mathbf{X} , under the constraint that the L_p norm of the perturbation is equal to or
148 less than a small ξ value (i.e., $\|\boldsymbol{\rho}\|_p \leq \xi$). The algorithms start with $\boldsymbol{\rho} = \mathbf{0}$ (no
149 perturbation) and iteratively update $\boldsymbol{\rho}$ by additively obtaining an adversarial perturbation
150 for an input image \mathbf{x} , which is randomly selected from \mathbf{X} without replacement. These
151 iterative updates continue until the number of iterations reaches the maximum i_{\max} .

152 The fast gradient sign method (FGSM) [7] is used to obtain an adversarial perturbation for
153 the input image; meanwhile, the original UAP algorithm [12] uses the DeepFool method
154 [24]. This is because the FGSM is used for both nontargeted and targeted attacks, and
155 DeepFool requires a higher computational cost than the FGSM and only generates a
156 nontargeted adversarial example for the input image. The FGSM generates the adversarial
157 perturbation for \mathbf{x} based on the loss gradient [7], with attack strength parameter ϵ .

158 Nontargeted and targeted UAPs were generated using the training images in the dataset.
159 Parameter ϵ was set to 0.0024; cases where $p = 2$ and ∞ were considered. Parameter
160 ξ was determined based on ratio ζ of the L_p norm of the UAP to the average L_p norm
161 of an image in the dataset. For the ISIC 2018 (skin lesion image) dataset, the average L_∞

162 and L_2 norms were 237 and 85,662, respectively. For the OCT image dataset, the average
163 L_∞ and L_2 norms were 253 and 15,077, respectively. For the chest X-ray image dataset,
164 the average L_∞ and L_2 norms were 253 and 40,738, respectively. Parameter i_{\max} was
165 set to 15.

166 To compare the performances of the generated UAPs with those of the random controls,
167 we generated random vectors (random UAPs) sampled uniformly from the sphere of a
168 specified radius [12].

169 *Vulnerability evaluation*

170 The fooling rate R_f and targeted attack success rate R_s were computed to evaluate the
171 vulnerability of the DNN models to nontargeted UAPs and targeted UAPs, respectively.
172 R_f for an image set is defined as adversarial images for which predicted labels are
173 inconsistent with the labels predicted from their associated clean images to all images in
174 the set (i.e., the probability that the labels predicted from clean images are inconsistent with
175 the label predicted from their adversarial images). R_s for an image set is the proportion of
176 adversarial images classified into the target class to all images in the set. It is noteworthy
177 that R_s has a baseline, defined as R_s observed without UAPs. Class (label) composition
178 in image data and prediction performance of DNNs affect the baseline. In this study, for
179 the OCT and chest X-ray image datasets, the baselines of R_s of targeted UAPs to a
180 specified class were $\sim 25\%$ and $\sim 50\%$, respectively. For the skin lesion dataset, the baselines
181 of R_s of targeted UAPs to MEL and NV were $\sim 10\%$ and $\sim 65\%$, respectively.

182 Additionally, we obtained the confusion matrices to evaluate the change in prediction due
183 to UAPs for each class. The confusion matrices were normalized to account for an
184 imbalanced dataset (ISIC 2018 dataset).

185 *Adversarial retraining*

186 Adversarial retraining was performed to increase the robustness of the DNN models to
187 UAPs [12, 15]; in particular, the models were fine-tuned with adversarial images. The
188 procedure was described in a previous study [12]. A brief description is as follows. i) Ten
189 UAPs against a DNN model were generated using the algorithm (for generating a
190 nontargeted or targeted UAP) with the (clean) training image set. ii) A modified training
191 image set was obtained by randomly selecting half of the training images and combining
192 them with the remaining, in which each image was perturbed by a UAP randomly selected
193 from 10 UAPs. iii) The model was fine-tuned by performing five additional epochs of
194 training on the modified training image set. iv) A new UAP against the fine-tuned model
195 was generated using the algorithm with the training image set. v) The R_f and R_s of the
196 UAP for the test images were computed. Steps i)–v) were repeated five times.

197 **Results**

198 *Performance for medical image classification*

199 The test and training accuracies of seven DNN models for three medical image datasets

200 are summarized in Table S1 in Addition file 1. The DNN models achieved good
201 accuracies. For the skin lesion image, OCT image, chest X-ray image datasets, the test
202 accuracies averaged over seven models were 87.3%, 95.8%, and 98.4%, respectively;
203 specifically, the test accuracies of Inception V3 models, which were frequently used in
204 previous studies on medical image diagnosis (e.g., [2, 3]), were 87.7%, 95.5%, and
205 97.6%, respectively. The normalized confusion matrices for the Inception V3 models on
206 the test images are shown in Fig. S1 in Addition file 1.

207 *Nontargeted universal adversarial attacks*

208 Meanwhile, the DNN models showed vulnerability to nontargeted UAPs. The fooling rates
209 R_f for both the training and test images increased rapidly with the perturbation magnitude
210 ζ and reached a high R_f , despite a low ζ (2%–6%). Figure 1 shows the case of
211 nontargeted UAPs with $p = 2$ against the Inception V3 models. The UAPs with $\zeta = 4\%$
212 achieved R_f of $> 80\%$ for the skin lesion (Fig. 1A) and chest X-ray image datasets (Fig.
213 1C), whereas slightly larger UAPs (with $\zeta = 6\%$) were required to achieve R_f of $\sim 70\%$
214 for the OCT image dataset (Fig. 1B). The R_f of the UAPs was significantly higher than
215 those of random UAPs. The confusion matrices on test images show that the models
216 classified most images into several specific classes (i.e., dominant classes) due to the UAPs
217 for the skin lesion and OCT image datasets. Specifically, most skin lesion images tended
218 to be classified into AKIEC or DF (Fig. 1D); moreover, most OCT images were classified
219 into CNV (Fig. 1E). For the chest X-ray image dataset, the model wrongly predicted the
220 true labels (Fig. 1F). A high R_f at a low ζ and dominant labels were observed in the case
221 of UAP with $p = \infty$ against the Inception V3 models for all medical image datasets (Fig.
222 S2 in Addition file 1). However, the skin lesion images tended to be classified into more
223 broad classes: BCC, AKIEC, BKL, or DF (Fig. S2D in Addition file 1).

224 The vulnerability to nontargeted UAPs is typically observed in other models. Table 1 shows
225 the R_f of the UAPs against the DNN models for the test images in the medical image
226 datasets. Overall, the small UAPs ($\zeta = 4\%$ for the skin lesion and chest X-ray image
227 datasets, and $\zeta = 6\%$ for the OCT image dataset) achieved a high R_f (70%–90%). The
228 R_f of the UAPs were significantly higher than those of random UAPs. R_f may depend
229 on model architectures; specifically, the R_f of the UAPs against the VGG16 and VGG19
230 models were $\sim 50\%$ for the chest X-ray image dataset, whereas those of the UAPs against
231 the other models were almost between 70% and 80%. This indicates that the models
232 classified images into either NORMAL or PNEUMONIA. In the case of UAPs with $p =$
233 2, the VGG16 and VGG19 models classified most test images into PNEUMONIA and
234 NORMAL, respectively (Fig. S3 in Addition file 1). In the case of UAPs with $p = \infty$, both
235 VGG16 and VGG19 models predicted most test images as NORMAL. This indicates that
236 the confusion matrix patterns (dominant classes) might change according to the model
237 architecture and p . Additionally, the change in confusion matrix patterns (on test images)
238 was observed in the skin lesion and OCT image datasets. For example, the VGG16 model
239 classified most skin lesion images into BKL owing to the UAP with $\zeta = 4\%$ and $p = 2$
240 (Figure S4A), whereas the Inception V3 models classified them into AKIEC or DF (Fig.
241 1D). The ResNet 50 model classified most OCT images into DME owing to the UAP with
242 $\zeta = 6\%$ and $p = 2$ (Fig. S4B in Addition file 1), whereas Inception V3 models classified

243 them into CNV (Fig. 1E).

244 The nontargeted UAPs (with $\zeta = 4\%$ for the skin lesion and chest X-ray image datasets,
245 and with $\zeta = 6\%$ for the OCT image dataset) were almost imperceptible. Figure 2 shows
246 the nontargeted UAPs with $p = 2$ against the Inception V3 models and examples of
247 adversarial images for the medical image datasets. The models predicted the original
248 images as their actual classes; however, they classified the adversarial images into any
249 incorrect class due to the nontargeted UAPs. The UAPs with $p = \infty$ and against the other
250 DNN models were also almost imperceptible for the skin lesion (Fig. S5 in Addition file
251 1), OCT (Fig. S6 in Addition file 1), and chest X-ray image datasets (Fig. S7 in Addition
252 file 1). The UAP patterns tended to depend on the model architecture for each medical
253 image dataset (Figs. S5–S7 in Addition file 1). The transferability of UAPs, which indicates
254 that UAPs generated based on DNNs with one model architecture can be used to deceive
255 DNNs with another model architecture, was not confirmed for the OCT (Table S3 in
256 Addition file 1) and chest X-ray image datasets (Table S4 in Addition file 1); however, a
257 weak transferability of UAPs was observed in the skin lesion image dataset (Table S5 in
258 Addition file 1). Specifically, the nontargeted UAPs with $p = 2$ generated based on the
259 Inception V3 models achieved R_f of $\sim 45\%$, $\sim 2\%$, and $\sim 10\%$ on average against the DNNs
260 with another model architecture for the skin lesion, OCT, and chest X-ray image datasets,
261 respectively.

262 *Targeted universal adversarial attacks*

263 Vulnerability to targeted UAPs was confirmed. Table 2 shows the targeted attack success
264 rates R_s of the UAPs with $p = 2$ against the DNN models for the test images in the
265 medical image datasets. As representative examples, we considered targeted attacks to the
266 most significant case and the control in each medical image dataset. For the skin lesion
267 image dataset, targeted attacks to MEL and NV were considered for the skin lesion image
268 dataset. For the OCT image dataset, targeted attacks to CNV and NM were considered. For
269 the chest X-ray image dataset, targeted attacks to PNEUMONIA and NORMAL were
270 considered. Overall, a high ($> 85\%$) R_s was observed regardless of the model architecture
271 despite small UAPs (with $\zeta = 2\%$ for the skin lesion and chest X-ray image datasets, and
272 with $\zeta = 6\%$ for the OCT image dataset). Furthermore, the confusion matrices (Fig.3)
273 indicate that the UAP-based targeted attacks were successful: most ($R_s\%$ of) test images
274 were classified into the targeted class owing to the UAPs. However, a smaller R_s was
275 partly observed according to the model architectures and datasets. For the skin lesion image
276 dataset, the R_s of the UAPs against the VGG16 ($\sim 40\%$) and VGG19 ($\sim 65\%$) models were
277 lower than those ($\sim 90\%$) of the UAPs against the other models. For the targeted attacks to
278 NM in the OCT image dataset, the R_s (30%–40%) of the UAPs against the VGG and
279 DensNet models were lower than those ($\sim 85\%$) of the UAPs against the other models. The
280 R_s of random UAPs were almost equivalent to those of the baselines. The R_s of the UAPs
281 were significantly higher than those of the random UAPs. Furthermore, a high R_s for a
282 small ζ were observed for the targeted UAPs with $p = \infty$ (Table S2 in Addition file 1).
283 However, the R_s for targeted attacks to MEL were lower overall, compared with the R_s
284 of the UAPs with $p = 2$. For example, the R_s of the UAPs with $p = 2$ and $p = \infty$
285 against the Inception V3 model were $\sim 95\%$ and $\sim 75\%$, respectively.

286 The targeted UAPs (with $\zeta = 2\%$ for the skin lesion and chest X-ray image datasets, and
287 with $\zeta = 6\%$ for the OCT image dataset) were also almost imperceptible. Figure 3 shows
288 the targeted UAPs with $p = 2$ against the Inception V3 models and examples of
289 adversarial images for the medical image datasets. The models predicted the original
290 images as their actual classes; however, they classified the adversarial images into the
291 targeted class owing to the UAPs. The UAPs with $p = \infty$ and against the other DNN
292 models were also almost imperceptible. For the skin lesion image dataset, Figures S8 and
293 S9 show the targeted attacks to NV and MEL, respectively. For the OCT image dataset,
294 Figures S10 and S11 in Addition file 1 show the targeted attacks to NM and CNV,
295 respectively. For the chest X-ray image dataset, Figures S12 and S13 in Addition file 1
296 show the targeted attacks to NORMAL and PNEUMONIA, respectively. The UAP patterns
297 tended to depend on the model architecture for each medical image dataset (Figs. S8–S13
298 in Addition file 1). The transferability of UAPs was not confirmed for the skin lesion (Table
299 S6), OCT (Table S7 in Addition file 1), and chest X-ray image datasets (Table S8 in
300 Addition file 1); specifically, the R_s observed when the targeted UAPs with $p = 2$
301 generated based on the Inception V3 model that attacked the DNN models with another
302 architecture was almost equivalent to their baselines of R_s ($\sim 10\%$, $\sim 25\%$, and $\sim 50\%$, for
303 the skin lesion, OCT, and chest X-ray image datasets, respectively).

304 *Adversarial retraining*

305 We analyzed the usefulness of adversarial retraining against universal adversarial attacks.
306 Figure 5 shows the effect of adversarial retraining on the R_f of nontargeted UAPs with
307 $p = 2$ against the Inception V3 models for the skin lesion, OCT, and chest X-ray image
308 datasets. $\zeta = 4\%$ for the skin lesion and chest X-ray image datasets; $\zeta = 6\%$ for the
309 OCT image dataset. Adversarial retraining did not affect the test accuracy. For the OCT
310 image dataset, R_f decreased with the iterations of adversarial retraining; specifically, R_f
311 decreased from 70.2% to 13.1% after five iterations (Fig. 5B); however, $\sim 40\%$ of the NM
312 images were still classified into an incorrect class (DME; Fig. 5E). The effect of adversarial
313 retraining on R_f was limited for the skin lesion (Fig. 5A) and chest X-ray image datasets
314 (Fig. 5B). For the chest X-ray image dataset, R_f decreased from 81.7% to 46.7%. A R_f
315 of $\sim 50\%$ indicates that the model classified most images into either one of two classes;
316 specifically, most images were classified into NORMAL at the fifth iteration (Fig. 5F). For
317 the skin lesion image dataset, any remarkable decrease in R_f due to adversarial retraining
318 was not confirmed; specifically, R_f decreased from 92.2% to 82.1% (Fig. 5A). Most
319 images were classified into MEL at the fifth iteration (Fig. 5C). However, the dominant
320 classes changed for each iteration. For example, the dominant classes were AKIEC and
321 BKL at the third and fourth iterations, respectively (Fig. S14 in Addition file 1).

322 Figure 6 shows the effect of adversarial retraining on the R_s of targeted UAPs with $p =$
323 2 against the Inception V3 models for the skin lesion, OCT, and chest X-ray image datasets.
324 As representative examples, we considered targeted attacks to the most significant cases:
325 MEL, CNV, and PNEUMONIA for the skin lesion, OCT, and chest X-ray image datasets,
326 respectively. $\zeta = 2\%$ for the skin lesion and chest X-ray image datasets; $\zeta = 6\%$ for the
327 OCT image dataset. Adversarial retraining did not affect the test accuracy and reduced R_s
328 for all medical image datasets (Figs. 6A–6C). For the OCT and chest X-ray image dataset,

329 R_s decreased from ~95% to the baseline R_s (~25% and ~50%, respectively) after five
330 iterations. For the skin lesion image dataset, R_s decreased from ~95% to ~30%; however,
331 R_s at the fifth iteration was higher than the baseline (~10%). The confusion matrices (Figs.
332 6D–6F) indicated that adversarial retraining was useful against UAP-based targeted
333 attacks: most images were correctly classified into the original classes despite the
334 adversarial attacks. However, the effect of adversarial retraining was partially limited for
335 the skin lesion image dataset. For example, 30% of the NV images were still classified into
336 the target class (MEL) despite five iterations of adversarial retraining (Fig. 6C).
337 Furthermore, ~20% of BKL and VASC images were still classified into the target class.

338 Discussion

339 We showed the vulnerability of the DNN models for medical image classification to small
340 UAPs. Furthermore, previous studies [10, 11] have indicated the vulnerability to
341 adversarial attacks toward medical DNNs; however, they were limited to input image-
342 dependent perturbations. In this study, we demonstrated that almost imperceptible UAPs
343 caused DNN misclassifications. Unlike previous assumptions, the results indicate that a
344 DNN-based medical image diagnosis is easier to deceive. Adversaries can result in failed
345 DNN-based medical image diagnoses at lower costs (i.e., using a single perturbation);
346 specifically, they do not need to consider the distribution and diversity of input images
347 when attacking DNNs using UAPs, as UAPs are image agnostic. The vulnerability to UAPs
348 was confirmed in various model architectures. The vulnerability to UAPs may be a
349 universal feature in DNNs.

350 We demonstrated that nontargeted attacks based on UAPs were possible (Figs. 1 and 2;
351 Table 1). Most images were classified into a few specific classes for the skin lesion and
352 OCT image (multiclass) datasets. This result is consistent with the existence of dominant
353 classes in UAP-based nontargeted attacks [12]. For the skin lesion image dataset, the
354 dominant classes of AKIEC and DF observed in this study may be due to the imbalanced
355 dataset. The numbers of AKIEC and DF images are relatively fewer than those of other
356 class images. As the algorithm considers maximizing the R_f , a relatively large R_f is
357 achieved when all inputs are classified into AKIEC and DF due to UAPs. The use of
358 imbalanced datasets may be one of the causes of vulnerability to UAPs. For the OCT image
359 (binary-class) dataset, the DNN models wrongly predicted the actual labels because of
360 R_f maximization; however, the existence of dominant classes was partly confirmed
361 according to the model architecture. These misclassifications result in false positives and
362 false negatives in medical diagnosis. False positives may cause unwanted mental stress to
363 patients, whereas false negatives may result in significant misdiagnoses involving human
364 lives; specifically, they fail to perform early detection and render therapeutic strategies
365 difficult. Moreover, they can cause the social credibility of medical doctors and medical
366 organizations to be undermined.

367 The transferability of nontargeted UAPs across model architectures was limited (Tables
368 S3–S5 in Addition file 1). This indicates that UAPs are architecture specific, inconsistent
369 with a previous study [12]. This discrepancy might be due to differences in the image
370 datasets. Specifically, the number of classes (2–7) in the medical image datasets was lower
371 than that (1,000) in the dataset in the previous study. This study partly considered grayscale

372 images, whereas the previous study used colored images. Transferability may be observed
373 in datasets comprising colored images with more classes. In fact, a weak transferability
374 was observed for the skin lesion image dataset (Table S5 in Addition file 1).

375 Furthermore, we showed that targeted attacks based on UAPs were possible in medical
376 image diagnosis (Figs. 3 and 4; Table 2), although the UAPs were not transferable across
377 model architectures (Tables S6–S8 in Addition file 1). The results imply that adversaries
378 can control DNN-based medical image diagnoses. As targeted attacks are more realistic,
379 they may result in more significant security concerns compared with nontargeted attacks.
380 In particular, adversaries can obtain any diagnosis; specifically, they can intentionally cause
381 not only problems resulting from misdiagnosis, but also various social disturbances. As
382 mentioned in a previous study [10], adversarial attacks can be used for insurance fraud as
383 well as drug and device approval adjustments, thereby fraudulently providing and
384 obtaining high-quality care when DNNs are used for decision making.

385 We considered adversarial retraining, which is known to be an effective method for
386 adversarial defenses [14], to reduce the vulnerability to UAPs. However, the effect of
387 adversarial retraining was limited for nontargeted UAPs (Fig. 5). For targeted attacks,
388 adversarial retraining reduced the vulnerability to UAPs significantly but did not
389 completely avoid it (particularly for the skin lesion image dataset; Fig. 6). In addition,
390 adversarial retraining requires high computational costs, as it is an iterative fine-tuning
391 method. Alternative simple methods such as dimensionality reduction (e.g., principle
392 component analysis), distributional detection (e.g., maximum mean discrepancy), and
393 normalization detection (e.g., dropout randomization) are available; however, they are
394 known to be easily detected as adversarial examples [15]. Despite the recent development
395 in adversarial defenses such as regularized surrogate loss optimization [25] and the use of
396 a discontinuous activation function [26], many promising defense methods have failed [27].
397 Defending against adversarial attacks is a cat-and-mouse game [10]. Furthermore,
398 properties inherent to image processing may cause misclassifications. For instance, DNN-
399 based image reconstructions are often performed for purifying adversarial examples [28];
400 however, it causes image artifacts, resulting in misclassifications by DNNs [29]. It may be
401 difficult to completely avoid security concerns caused by adversarial attacks.

402 A simple solution for avoiding adversarial attacks is to render DNNs closed source and
403 publicly unavailable; however, this hinders the accelerated development of medical DNNs
404 and practical applications of DNNs to automated support for clinical diagnosis. Because
405 the amount of medical image data is limited, collaboration among multiple institutions is
406 required to achieve high diagnosis performance [30]. For similar reasons, medical DNNs
407 are often developed by fine-tuning existing DNNs such as VGG, ResNet, and Inception
408 pretrained using the ImageNet dataset (i.e., via transfer learning), although a previous study
409 [30] debated the effect of transfer learning on the improvement in prediction performance
410 for medical imaging; consequently, model architectures and model weights may be
411 important. Furthermore, DNNs are aimed for real-world usage (e.g., automated support for
412 clinical diagnosis). The assumption that DNNs are closed source and publicly unavailable
413 may be unrealistic. Even if DNNs are black-box (e.g., model architectures and weights are
414 unknown and loss gradient is not accessible), adversarial attacks on DNNs may be possible.
415 Several methods for adversarial attacks on black-box DNNs, which estimate adversarial

416 perturbations using only model outputs (e.g., confidence scores), have been proposed [31–
417 33]. The development and operation of secure, privacy-preserving, and federated DNNs
418 are required in medical imaging [6].

419 **Conclusion**

420 We demonstrated the vulnerability of DNNs for medical image classification to
421 nontargeted and targeted UAPs. Our findings emphasized that careful consideration is
422 required in developing DNNs for medical imaging and their practical applications. Our
423 study enhances our understanding of the vulnerabilities of DNNs to adversarial attacks and
424 may help increase the security of DNNs. UAPs are useful for reliability evaluation and for
425 designing the operation strategy of medical DNNs.

426 **List of abbreviations**

427 **AKIEC:** actinic keratosis/Bowens disease (intraepithelial carcinoma)

428 **BCC:** basal cell carcinoma

429 **BKL:** benign keratosis (solar lentigo/seborrheic keratosis/lichen planus-like keratosis)

430 **CNV:** neovascular membrane and associated subretinal fluid

431 **DF:** dermatofibroma

432 **DME:** diabetic macular edema with retinal-thickening-associated intraretinal fluid

433 **DNN:** deep neural network

434 **DRUSEN:** multiple drusen present in early age-related macular degeneration

435 **DensNet:** dense convolutional network

436 **FC:** fully connected

437 **FGSM:** fast gradient sign method

438 **ISIC:** International Skin Imaging Collaboration

439 **MEL:** melanoma

440 **NM:** normal retina with preserved foveal contour and absence of any retinal fluid/edema

441 **NV:** melanocytic nevus

442 **OCT:** optical coherence tomography

443 **ResNet:** residual network

444 **UAP:** universal adversarial perturbation

445 **VASC:** vascular lesion

446 **VGG:** visual geometry group

447 **Declarations**

448 *Ethics approval and consent to participate*

449 Not applicable.

450 *Consent for publication*

451 Not applicable.

452 *Availability of data and material*

453 All data generated and analyzed during this study are included in this published article
454 and its supplementary information files. The code and data used in this study are
455 available from our GitHub repository: github.com/hkthirano/MedicalAI-UAP.

456 *Competing interests*

457 The authors declare that they have no competing interests.

458 *Funding*

459 No funding was received.

460 *Authors' contributions*

461 KT conceived and designed the study. HH and AM prepared the data and models. HH
462 coded and performed experimental evaluation. HH and KT interpreted the results. HH
463 and KT wrote the manuscript. All authors provided the final approval for publication.

464 *Acknowledgments*

465 We would like to thank Editage (www.editage.jp) for English language editing.

466 **References**

467 1. Litjens G, Kooi T, Bejnordi BE, Setio AAA, Ciompi F, Ghafoorian M, et al. A survey
468 on deep learning in medical image analysis. *Med Image Anal.* 2017;42 December
469 2012:60–88. doi:10.1016/j.media.2017.07.005.

470 2. Esteva A, Kuprel B, Novoa RA, Ko J, Swetter SM, Blau HM, et al. Dermatologist-
471 level classification of skin cancer with deep neural networks. *Nature.* 2017;542:115–8.
472 doi:10.1038/nature21056.

- 473 3. Kermany DS, Goldbaum M, Cai W, Valentim CCS, Liang H, Baxter SL, et al.
474 Identifying Medical Diagnoses and Treatable Diseases by Image-Based Deep Learning.
475 Cell. 2018;172:1122-1131.e9. doi:10.1016/j.cell.2018.02.010.
- 476 4. Liu X, Faes L, Kale AU, Wagner SK, Fu DJ, Bruynseels A, et al. A comparison of
477 deep learning performance against health-care professionals in detecting diseases from
478 medical imaging: a systematic review and meta-analysis. Lancet Digit Heal.
479 2019;1:e271–97. doi:10.1016/S2589-7500(19)30123-2.
- 480 5. Rudin C. Stop explaining black box machine learning models for high stakes decisions
481 and use interpretable models instead. Nat Mach Intell. 2019;1:206–15.
482 doi:10.1038/s42256-019-0048-x.
- 483 6. Kaissis GA, Makowski MR, Rückert D, Braren RF. Secure, privacy-preserving and
484 federated machine learning in medical imaging. Nat Mach Intell. 2020;2:305–11.
485 doi:10.1038/s42256-020-0186-1.
- 486 7. Goodfellow IJ, Shlens J, Szegedy C. Explaining and harnessing adversarial examples.
487 2014. <http://arxiv.org/abs/1412.6572>.
- 488 8. Yuan X, He P, Zhu Q, Li X. Adversarial examples: attacks and defenses for deep
489 learning. IEEE Trans Neural Networks Learn Syst. 2019;30:2805–24.
490 doi:10.1109/TNNLS.2018.2886017.
- 491 9. Matyasko A, Chau L-P. Improved network robustness with adversary critic. 2018.
492 <http://arxiv.org/abs/1810.12576>.
- 493 10. Finlayson SG, Bowers JD, Ito J, Zittrain JL, Beam AL, Kohane IS. Adversarial
494 attacks on medical machine learning. Science (80-). 2019;363:1287–9.
495 doi:10.1126/science.aaw4399.
- 496 11. Asgari Taghanaki S, Das A, Hamarneh G. Vulnerability Analysis of Chest X-Ray
497 Image Classification Against Adversarial Attacks. In: Understanding and Interpreting
498 Machine Learning in Medical Image Computing Applications. 2018. p. 87–94.
499 doi:10.1007/978-3-030-02628-8_10.
- 500 12. Moosavi-Dezfooli SM, Fawzi A, Fawzi O, Frossard P. Universal adversarial
501 perturbations. Proc - 30th IEEE Conf Comput Vis Pattern Recognition, CVPR 2017.
502 2017;2017-Janua:86–94.
- 503 13. Hirano H, Takemoto K. Simple iterative method for generating targeted universal
504 adversarial perturbations. In: Proceedings of 25th International Symposium on Artificial
505 Life and Robotics. 2020. p. 426–30. <http://arxiv.org/abs/1911.06502>.
- 506 14. Madry A, Makelov A, Schmidt L, Tsipras D, Vladu A. Towards Deep Learning
507 Models Resistant to Adversarial Attacks. In: International Conference on Learning
508 Representations. 2018. <https://openreview.net/forum?id=rJzIBfZAb>.
- 509 15. Carlini N, Wagner D. Adversarial examples are not easily detected. In: Proceedings

- 510 of the 10th ACM Workshop on Artificial Intelligence and Security - AISEC '17. New
511 York, New York, USA: ACM Press; 2017. p. 3–14. doi:10.1145/3128572.3140444.
- 512 16. Wong E, Rice L, Kolter JZ. Fast is better than free: Revisiting adversarial training. In:
513 International Conference on Learning Representations. 2020.
514 <https://openreview.net/forum?id=BJx040EFvH>.
- 515 17. Russakovsky O, Deng J, Su H, Krause J, Satheesh S, Ma S, et al. ImageNet Large
516 Scale Visual Recognition Challenge. *Int J Comput Vis*. 2015.
- 517 18. Szegedy C, Vanhoucke V, Ioffe S, Shlens J, Wojna Z. Rethinking the Inception
518 Architecture for Computer Vision. In: 2016 IEEE Conference on Computer Vision and
519 Pattern Recognition (CVPR). IEEE; 2016. p. 2818–26. doi:10.1109/CVPR.2016.308.
- 520 19. Simonyan K, Zisserman A. Very deep convolutional networks for large-scale image
521 recognition. In: 3rd International Conference on Learning Representations, ICLR 2015 -
522 Conference Track Proceedings. 2015.
- 523 20. He K, Zhang X, Ren S, Sun J. Deep Residual Learning for Image Recognition. In:
524 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). IEEE;
525 2016. p. 770–8. doi:10.1109/CVPR.2016.90.
- 526 21. Szegedy C, Ioffe S, Vanhoucke V, Alemi AA. Inception-v4, inception-ResNet and
527 the impact of residual connections on learning. In: 31st AAAI Conference on Artificial
528 Intelligence, AAAI 2017. 2017.
- 529 22. Huang G, Liu Z, Van Der Maaten L, Weinberger KQ. Densely connected
530 convolutional networks. In: Proceedings - 30th IEEE Conference on Computer Vision
531 and Pattern Recognition, CVPR 2017. 2017.
- 532 23. Nicolae M-I, Sinn M, Tran MN, Buesser B, Rawat A, Wistuba M, et al. Adversarial
533 Robustness Toolbox v1.0.0. 2018. <http://arxiv.org/abs/1807.01069>.
- 534 24. Moosavi-Dezfooli S-M, Fawzi A, Frossard P. DeepFool: a simple and accurate
535 method to fool deep neural networks. In: 2016 IEEE Conference on Computer Vision and
536 Pattern Recognition (CVPR). IEEE; 2016. p. 2574–82. doi:10.1109/CVPR.2016.282.
- 537 25. Zhang H, Yu Y, Jiao J, Xing E, Ghaoui L El, Jordan M. Theoretically Principled
538 Trade-off between Robustness and Accuracy. In: Chaudhuri K, Salakhutdinov R, editors.
539 Proceedings of the 36th International Conference on Machine Learning. Long Beach,
540 California, USA: PMLR; 2019. p. 7472–82.
541 <http://proceedings.mlr.press/v97/zhang19p.html>.
- 542 26. Xiao C, Zhong P, Zheng C. Enhancing adversarial defense by k-winners-take-all.
543 Proc 8th Int Conf Learn Represent. 2020. <http://arxiv.org/abs/1905.10510>.
- 544 27. Croce F, Hein M. Reliable evaluation of adversarial robustness with an ensemble of
545 diverse parameter-free attacks. Proc 37th Int Conf Mach Learn. 2020. doi:2003.01690.

- 546 28. Hwang U, Park J, Jang H, Yoon S, Cho NI. PuVAE: a variational autoencoder to
547 purify adversarial examples. *IEEE Access*. 2019;7:126582–93.
548 doi:10.1109/ACCESS.2019.2939352.
- 549 29. Antun V, Renna F, Poon C, Adcock B, Hansen AC. On instabilities of deep learning
550 in image reconstruction and the potential costs of AI. *Proc Natl Acad Sci*.
551 2020;:201907377. doi:10.1073/pnas.1907377117.
- 552 30. Chang K, Balachandar N, Lam C, Yi D, Brown J, Beers A, et al. Distributed deep
553 learning networks among institutions for medical imaging. *J Am Med Informatics Assoc*.
554 2018;25:945–54.
- 555 31. Chen J, Su M, Shen S, Xiong H, Zheng H. POBA-GA: Perturbation optimized black-
556 box adversarial attacks via genetic algorithm. *Comput Secur*. 2019;85:89–106.
557 doi:10.1016/j.cose.2019.04.014.
- 558 32. Guo C, Gardner JR, You Y, Wilson AG, Weinberger KQ. Simple black-box
559 adversarial attacks. *Proc 36th Int Conf Mach Learn*. 2019;:2484–93.
560 <http://arxiv.org/abs/1905.07121>.
- 561 33. Co KT, Muñoz-González L, de Maupeou S, Lupu EC. Procedural noise adversarial
562 examples for black-box attacks on deep convolutional networks. In: *Proceedings of the*
563 *2019 ACM SIGSAC Conference on Computer and Communications Security*. New York,
564 NY, USA: ACM; 2019. p. 275–89. doi:10.1145/3319535.3345660.
- 565
- 566

567 **Tables**

568 **Table 1:** Fooling rates R_f (%) of nontargeted UAPs against various DNN models for test
 569 images of skin lesion, OCT, and chest X-ray image datasets. $\zeta = 4\%$ for skin lesion and
 570 chest X-ray image datasets. $\zeta = 6\%$ for OCT image dataset. Values in brackets are R_f
 571 of random UAPs (random controls).

Model architecture	Skin lesion		OCT		Chest X-ray	
	$p = 2$	$p = \infty$	$p = 2$	$p = \infty$	$p = 2$	$p = \infty$
Inception V3	92.2 (14.1)	90.0 (11.8)	70.2 (1.0)	73.9 (3.4)	81.7 (2.4)	79.8 (3.0)
VGG16	87.6 (4.9)	86.4 (3.5)	72.4 (0.2)	74.9 (1.8)	49.8 (2.2)	50.0 (2.2)
VGG19	89.2 (5.2)	87.0 (3.7)	72.8 (0.4)	74.7 (2.1)	49.3 (3.9)	49.3 (4.4)
ResNet50	91.9 (11.6)	87.9 (10.1)	71.2 (1.1)	74.8 (5.4)	72.6 (7.2)	73.0 (7.4)
Inception ResNet V2	94.5 (16.7)	90.3 (15.2)	69.6 (1.4)	74.0 (3.2)	78.0 (2.6)	77.0 (3.3)
DenseNet 121	93.8 (12.0)	82.9 (10.2)	68.8 (1.3)	73.0 (3.6)	69.8 (3.9)	71.7 (4.1)
DenseNet 169	93.8 (11.7)	84.2 (9.1)	50.3 (1.3)	72.3 (4.0)	67.6 (2.8)	71.3 (3.7)

572

573

574 **Table 2:** Targeted attack success rates R_s (%) of targeted UAPs with $p = 2$ against
 575 various DNN models to each target class. R_s were for test images. $\zeta = 2\%$ for skin
 576 lesion and chest X-ray image datasets. $\zeta = 6\%$ for OCT image dataset. Values in brackets
 577 are R_s of random UAPs (random controls).

Model architecture / Target class	Skin lesion		OCT		Chest X-ray	
	NV	MEL	NM	CNV	NORMAL	PNEUMONIA
Inception V3	93.3 (65.6)	94.4 (12.2)	84.1 (25.7)	95.9 (24.8)	96.1 (52.8)	93.3 (47.2)
VGG16	89.6 (71.7)	40.4 (8.3)	32.4 (25.4)	97.7 (24.9)	95.6 (50.2)	95.0 (49.8)
VGG19	91.6 (72.1)	64.6 (8.7)	41.2 (25.9)	97.5 (24.9)	97.6 (51.7)	95.2 (48.3)
ResNet50	97.9 (66.5)	92.4 (11.8)	84.9 (25.8)	98.5 (24.5)	95.7 (53.5)	95.2 (46.5)
Inception ResNet V2	92.4 (61.0)	97.3 (16.1)	84.5 (25.6)	96.2 (24.7)	98.3 (53.1)	93.9 (46.9)
DenseNet 121	92.1 (65.2)	90.5 (13.4)	41.8 (25.3)	88.1 (24.7)	94.8 (51.9)	92.0 (48.1)
DenseNet 169	92.9 (65.8)	92.9 (12.2)	41.7 (25.0)	92.7 (24.2)	95.7 (52.0)	93.1 (48.0)

578

579

580 Figure captions

581 **Fig. 1:** Vulnerability to nontargeted UAPs with $p = 2$. Line plots of the fooling rate R_f
582 against Inception V3 model versus perturbation magnitude ζ for skin lesion (A), OCT (B),
583 and chest X-ray (C) image datasets. Legend label indicates image set used for computing
584 R_f . Additional argument “(random)” indicates that random UAPs were used instead of
585 UAPs. Normalized confusion matrices for Inception V3 models attacked using UAPs on
586 test images of skin lesion (D), OCT (E), and chest X-ray (F) image datasets are also shown.
587 $\zeta = 4\%$ in (D) and (F). $\zeta = 6\%$ in (E).

588 **Fig. 2:** Nontargeted UAPs with $p = 2$ against Inception V3 models and their adversarial
589 images for skin lesion (A), OCT (B), and chest X-ray image datasets (C). $\zeta = 4\%$ in (A)
590 and (C). $\zeta = 6\%$ in (B). Labels in brackets beside the images are the predicted classes.
591 The original (clean) images are correctly classified into their actual labels. UAPs are
592 emphatically displayed for clarity; in particular, each UAP is scaled by a maximum of 1
593 and minimum of 0.

594 **Fig. 3:** Normalized confusion matrices for Inception V3 models attacked with targeted
595 UPAs with $p = 2$ on test images in skin lesion (left panels), OCT (middle panels), and
596 chest X-ray image datasets (right panels). $\zeta = 2\%$ for skin lesion and chest X-ray image
597 datasets. $\zeta = 6\%$ for OCT image dataset.

598 **Fig. 4:** Targeted UAPs with $p = 2$ against Inception V3 models and their adversarial
599 images for skin lesion (A), OCT (B), and chest X-ray image datasets. $\zeta = 2\%$ in (A) and
600 (C). $\zeta = 6\%$ in (B). Labels in brackets beside the images are predicted classes. Original
601 (clean) images were correctly classified into their actual labels. Adversarial images were
602 classified into the target classes. UAPs are emphatically displayed for clarity; in particular,
603 each UAP is scaled by a maximum of 1 and minimum of 0.

604 **Fig. 5:** Effect of adversarial retraining on robustness of nontargeted UAPs with $p = 2$
605 against Inception V3 models for skin lesion, OCT, and chest X-ray image datasets. $\zeta =$
606 4% for the skin lesion and chest X-ray image datasets. $\zeta = 6\%$ for OCT image dataset.
607 Top panels indicate scatter plots of fooling rate R_f (%) of UAPs versus number of
608 iterations for adversarial retraining. Bottom panels indicate normalized confusion matrices
609 for fine-tuned models obtained after five iterations of adversarial retraining. These
610 confusion matrices are on adversarial test images.

611 **Fig. 6:** Effect of adversarial retraining on robustness of targeted UAPs with $p = 2$ against
612 Inception V3 models for skin lesion, OCT, and chest X-ray image dataset. $\zeta = 2\%$ for
613 skin lesion and chest X-ray image datasets. $\zeta = 6\%$ for OCT image dataset. Top panels
614 indicate scatter plots of targeted attack success rate R_s (%) of UAPs versus number of
615 iterations for adversarial retraining. Bottom panels indicate normalized confusion matrices
616 for fine-tuned models obtained after five iterations of adversarial retraining. These
617 confusion matrices are on adversarial test images.

618 **Additional files**

619

620 **Additional file 1:** Supplementary tables and figures. (PDF)

621

622