Artificial Intelligence: A Blessing or a Curse. Discussing Security Concerns of Diagnostic Models in Radiological Assessment

Anonymous Full Paper Submission 10

OD1 Abstract

Radiology is increasingly adopting AI-based work-002 flows, which provide promise but also introduce new 003 security concerns. The goal of this research is to en-004 hance the security of these workflows by evaluating 005 the risks of data poisoning attacks using the Fast 006 Gradient Sign Method (FGSM) and Carlini-Wagner 007 (C&W) techniques. Detection methods commonly 008 employed in financial fraud are evaluated to assess 009 their effectiveness in this context. Knowledge dis-010 tillation is also explored as a defense mechanism 011 against data poisoning, offering a potential mitiga-012 tion strategy. By conducting these evaluations and 013 proposing defenses, this research aims to contribute 014 015 to the discussion of more robust deployment of AI systems in real-world radiology applications. 016

1 Introduction

Radiology has had a tendency to be an early adopter 018 of new technologies. Medical imaging is at the fore-019 front of research, since it is part of the initial pro-020 cesses for making a diagnosis [1]. The use of AI, 021 along with the discipline's eagerness, is part of a 022 tradition. There has been the X-ray, MRI, and 023 CT scan, each readily adopted into workflows. AI 024 marks a continuation of that trend as it becomes 025 more commonplace. However, AI fundamentally 026 differs from previous technological advancements. 027 Other imaging technologies simply provide more 028 data, but AI seeks to complement the radiologist. 029 It is also more mercurial by nature. Other techno-030 logical advancements are deterministic. If one takes 031 an X-ray, there is a high degree of certainty that 032 the image will yield the same result. This is not the 033 case with AI; it is stochastic. Given images, models 034 simply report what is most probable from their per-035 spective. It lacks the same promise of consistency. 036 Software engineers approach this issue by defining 037 key performance indicators and creating thresholds 038 for performance. They ensure that the model will 039 perform above a certain percent accuracy or with 040 a degree of specificity [2]. This measure is usually 041 based on the confusion matrix and Area Under the 042 Curve. 043

1.1 Data Drift

044

067

To ensure that performance metrics are maintained, 045 software engineers must regularly update and main-046 tain the models under their stewardship. The main 047 challenge is that the performance of models tends 048 to fluctuate and degrade over time [3]. This phe-049 nomenon is known as data drift. For example, con-050 sider a model trained to analyze X-rays of legs from 051 people in Scandinavia. If there were a significant 052 migration of shorter individuals into the population, 053 the model's performance would likely decline. The 054 AI relies on identifying certain characteristics in 055 specific locations. Due to the variation in height, 056 discrepancies could arise between what the AI has 057 learned and the data it encounters. Thus, leading to 058 degraded performance. This is data drift: the dis-059 crepancy between the environment in which the AI 060 operates and the data on which it was trained. To 061 address these issues, it is common practice to update 062 models with more representative data. While this 063 practice aims to create more robust models, it also 064 opens up opportunities for bad actors to manipulate 065 the results of models for malicious purposes. 066

1.2 Data Poisoning

One technique to keep models aligned with the cur-068 rent deployment context is to train them on data 069 that the AI processes during inference. Essentially, 070 the model may save the input provided by the user 071 along with the conclusion made, thereby reinforc-072 ing its predictive ability. However, this opens the 073 door for bad actors to craft payloads designed to 074 intentionally degrade the model, encouraging it to 075 learn incorrect associations [4]. This is known as 076 data poisoning. Consequently, while your model 077 updates to avoid data drift, it could inadvertently 078 poison itself. One critical aspect of data poisoning 079 to understand is that human inspection alone would 080 not detect anything amiss. The alterations made to 081 poison the data are not something the human eye 082 can easily discern, adding a layer of stealth to these 083 attacks if an adversary is determined to undermine. 084

Data Poisoning is an increasing significant concern through advancements in deployment of AI within clinical settings. For instance, Project MONAI is a framework that has become popular and focuses on clinical use-cases. The MONAI Deploy extension has already been utilized in the wild and implements continuous learning capabilities [5]. With increased
deployment, the risks associated with using AI in
clinical workflows will increase with time.

094 1.3 Contribution

There are several aims this paper seeks to target. 095 The first is to investigate adversarial attacks. Two 096 of the most common attacks are FGSM and C&W. 097 The literature suggests that luminosity is a factor in 098 the ability to conduct attacks. [6]. Thus, the data 099 from 2017 RSNA Pediatric Bone Age Challenge is 100 investigated [7]. The contrast between the brightness 101 of the bones and the black background make this 102 analysis within the domain of adversarial attacks 103 of medical images significant. The dataset is also 104 selected because the images are more sparse, whereas 105 typical scans usually contain significantly more noise. 106 By using X-rays, the manipulations become clearer 107 and more understandable, and the structures of the 108 bones are more defined. It is also multi-class. A 109 second dataset is used for only segmentation. It is 110 the Kvasir-SEG dataset [8]. This dataset is designed 111 for the segmentation of polyps. This is valuable as 112 a contrast to Bone Age. The images are color and 113 are of the GI tract; these are visually more complex. 114 By conducting this analysis, it provides guidance 115 how adversarial attacks affect segmentation and not 116 only classification. 117

The second aim is to better understand the applicability of Benford's Law to medical imaging. There has only been some initial investigation in the literature [9]. This requires consensus building to uncover the strengths and limitations of this approach.

The last aim is to explore how utilizing different loss functions within distillation affects their defensive capability and how different tasks are affected by adversarial attacks.

127 2 Background

Perturbations are the primary attack methods used 128 against computer vision models and pose a signifi-129 cant risk to models deployed in the healthcare space. 130 Therefore, it is essential to understand these attacks 131 and how they are used in context. Additionally, de-132 tection methods from other fields, such as Benford's 133 Law, are worth considering. Finally, it is important 134 to explore defenses against these attacks to learn 135 how developers can be both proactive in detection 136 and reactive by fortifying their models if poisoned 137 data slips through. 138

139 2.1 Computer Vision

YOLO (You Only Look Once) and U-Net are the
focus of this discussion as they represent the most
widely used models. YOLO is a productionized

object detection and classification model that is 143 widely used in industry. The first iteration of YOLO 144 was published in 2016 [10]. There have been upwards 145 of a eleven versions, and it has been foundational 146 in computer vision. YOLO will be primary focus of 147 this study. U-Net is a convolutional neural network 148 (CNN) like YOLO, but it is designed for medical 149 segmentation. It utilizes an encoder to compress 150 the inputs into the feature space and then uses a 151 decoder to up-sample to recover spatial information 152 to create the segmentation map [11]. A standard 153 U-Net is utilized in this work. 154

YOLO will be the focus of the work as it is gen-155 erally more popular. YOLO segments an image 156 based on labeled images with annotated bounding 157 boxes. It then predicts the presence of particular 158 classes within the image during inference. YOLO is 159 currently the benchmark for object detection tasks 160 and has been enhanced through a series of innova-161 tions. The first key innovation of YOLO is the use of 162 Cross-Stage Partial Networks [12], which minimizes 163 the computational load required for convolutions 164 by employing a partitioning strategy. The second 165 innovation involves Path Aggregation Networks [13], 166 a technique that enables the development of sub-167 networks for more robust feature pooling. Lastly, 168 YOLO produces three feature maps that are fused 169 to create a more informative output, which is then 170 further interpreted [12]. These innovations have so-171 lidified YOLO as the standard in object detection 172 and classification. 173

2.2 Attacks 174

There are two main attack methods evaluated in this 175 endeavor. The first is Fast Gradient Sign Method. 176 This method is far simpler and is more well suited 177 for general attacks. 178

$$x' = x + \epsilon \cdot \operatorname{sign}(\nabla_x J(\theta, x, y)) \tag{1}$$
 179

In FGSM, there is your input, x. In this context, 180 it is an image. ϵ is essentially the amount of manip-181 ulation the attack should inflict to the image. It is a 182 tuning parameter. The larger the value, the greater 183 the amount of perpetuation. However, this increases 184 the visibility of attack and increases the likelihood 185 that your attack will be detected. $\nabla_x J(\theta, x, y)$ is the 186 gradient of the loss function and is used by FGSM to 187 identify the directionality of perturbations to maxi-188 mize the model's prediction error [14]. This can be 189 better conceptualized when considering a prediction 190 landscape and pushing the poisoned data against 191 the gradient. 192

$$\min_{s} \|\delta\|_p + c \cdot f(x+\delta) \tag{2} 193$$

Carlini & Wagner is the second attack evaluated. 194 The C&W attack is an optimization-based approach. 195

262

The perturbation δ is measured with an L_2 norm in 196 most cases. This helps to constrain the space of the 197 perturbation while still allowing it to misclassify the 198 input. The constant c is a tuning parameter. It helps 199 to shrink δ to the point where it may still misclassify 200 inputs without it being larger than necessary, which 201 could draw attention and increase the chance of 202 detection [15]. 203

204 2.3 Detection

To detect these attacks, techniques from other do-205 mains are being explored. The investigators in [9] 206 apply methods grounded in Benford's Law to de-207 tect such attacks. Benford's Law suggests that the 208 209 leading digit distribution of natural datasets follows a logarithmic function. In their analysis, the goal 210 was to test whether the leading digit distribution 211 of pixel values deviates from that of natural images. 212 The approach involved treating the input image as a 213 vector, computing a gradient, transforming the data, 214 and then comparing the frequency of the first digits 215 using the Kolmogorov-Smirnov test [16]. The gradi-216 ent magnitude of the input image and the Discrete 217 Cosine Transformation typically adhere to Benford's 218 Law, and these were used in the comparison. The 219 Kolmogorov-Smirnov test was employed to deter-220 mine the divergence of these distributions. Using 221 these techniques, the investigators successfully iden-222 tified 94.7% of a Projected Gradient Descent attack 223

with infinity norm and 81.8% with L2 norm.

225 2.4 Defenses

Proactive measures can be taken to mitigate attacks, 226 rather than simply focusing on reducing response 227 times. In [17], the investigators introduce knowl-228 edge distillation, a strategy designed to make models 229 less susceptible to adversarial input by facilitating 230 knowledge transfer. Essentially, class probability 231 vectors are fed into a smaller secondary model that 232 produces a more discrete result. While the primary 233 goal is to enhance the model's robustness, knowl-234 edge distillation is designed to have minimal impact 235 on the model's architecture, maintain accuracy, and 236 preserve speed. The technique significantly reduced 237 the success rate of adversarial sample crafting from 238 95.89% to 0.45% on the MNIST dataset [17]. Ad-239 ditionally, the distilled model experienced only a 240 1.28% drop in accuracy, without requiring signifi-241 cant changes to the base model. 242

243 **3** Literature Review

U-Nets have been developed for medical segmentation and have been a focus of attacks along with
YOLO. In [6], the investigators applied FGSM to
U-Nets. They found image luminosity is a factor





(a) X-Ray of 1 Month Old

(b) X-Ray of 5 Month Old

Figure 1. Comparison of X-Ray Features

for attack success. This suggestion is in part moti- 248 vated using the Bone Age dataset in our analysis. 249 They used the Dice Score and found a 40% drop 250 after applying only FGSM to four varieties of U-251 Net. Furthering the investigation of FGSM attacks, 252 the investigators in [18] also applied C&W attacks 253 and attempted to defend against both. However, 254 this work was outside of the medical context. They 255 found that distillation was effective at the mitigation 256 of FGSM but not C&W. The loss function used in 257 their distillation process utilized a singular weight-258 ing parameter to balance loss of the student model 259 with the distillation loss, which can be thought as 260 agreement between the student and teacher model. 261

4 Detection Techniques

To guard against these attacks and to ensure quality checks for AI, qualitative and quantitative techniques have arisen. The first technique is more qualitative. It is inspecting saliency plots to understand which aspects of the input are most significant [19]. 267 Applications of Benford's Law are more quantitative measures for detecting manipulation. 269

When reviewing these two plots, one can interpret 270 which features are most significant to the classifi-271 cation of hand bone ages. For the one-month-old. 272 the structure and definition of the carpal bones ap-273 pear to be most significant. For instance, one might 274 notice the roundness due to the lack of structure. 275 Additionally, there is increased highlighting in the 276 larger gaps between the phalanges. For the five-277 month-old, the carpal bones are more well-defined, 278 and there is notable highlighting of the radius and 279 ulna in the scan. 280

One issue to consider is the presence of the "L" 281 marker. This marker is intended to help radiologists 282 orient the scan by indicating the left side of the 283 person being scanned. The strongest highlighting is 284 found along the contours of the letters on the marker. 285 The concern is that the model might be associating 286 the marker with older hands. Younger hands are 287 smaller and often require zooming in, which can 288









Figure 2. Comparison in Digit Frequency

obscure the marker and implicitly train the model 289 to recognize the marker as a feature associated with 290 older hands. If the marker were superimposed on 291 younger hands, they might be incorrectly interpreted 292 as older due to this spurious relationship. Thus, 293 models may sometimes hamper their own learning 294 without external influence and may simply learn 295 problematic features. 296

297 4.1 Benford's Law

Benford's Law is an observation regarding the fre-298 quency of digits. For instance, the number "1" 299 should appear in the first digit of a number roughly 300 30% of the time. The number "2" should appear as 301 the first digit in a number 18% of the time. This 302 pattern in frequency decreases with the higher the 303 digit. The underlying process causing this distribu-304 tion is not well understood, but it provides a useful 305 benchmark to suggest if the data has been manip-306 ulated not. In [9], a cosine transform is applied to 307 the images which allows Benford's Law to hold. 308

A scan of a five month old is evaluated to determine if an attack by C&W could be determined by inspection for Bone Age. For Kvasir, C&W is also used on the image of a polyp.

When observing how these distribution differ from the expectation, it is suggestive that the has been manipulated. This is most noticeable when observing the frequency of "1" in both datasets.

A limitation to consider is that when applying the Kolmogorov-Smirnov test, the p-value tends to

produce too many false positives. It is likely that this 319 test should be applied to entire datasets rather than 320 individual images. If Benford's Law is to be applied 321 to specific images, it is more useful to inspect the 322 distribution through graphs or to use more advanced 323 versions of Benford's Law to increase the robustness 324 of the detection. For instance, it may be more robust 325 to jointly look at the first and second digits to avoid 326 false positives. 327

When evaluating Benford's Law on an image at-328 tacked by the FGSM, the technique did not perform 329 well. This holds for both datasets. The attack 330 caused the frequencies to more closely align with 331 the expected frequencies, leading to false negatives. 332 Observing entire datasets may help mitigate these 333 deviations in individual scans and determine whether 334 there has been a systematic attack on the dataset. 335

5 Results

336

YOLO is a model for object detection, and in this 337 formulation, the ages of bones were treated as classes 338 to segment. The accuracy measurement reflects how 339 well YOLO was able to determine a bounding box 340 for the correct age of the bone from the scan. After 341 training YOLO for 10 epochs, the model achieved 342 an accuracy of 64% in segmenting the images to 343 determine the correct class. This is a reasonable 344 baseline, with the highest accuracy being 78% for 345 1-month-old scans. However, it performed the least 346 well for 8-month-old scans, with an accuracy of only 347 28%.348

When applying the FGSM, the average accuracy 349 drops to 17% on the adversarial example set. For the 350 less representative classes, YOLO was completely 351 fooled. For instance, none of the 8-month-old scans 352 were correctly classified. However, for larger classes, 353 the model performed slightly better, with 5-month-354 old bones reaching 32% accuracy and 4-month-old 355 bones reaching 26 356

As for the C&W attack, the overall accuracy 357 drops to 12%. YOLO was again tricked by the 358 less representative classes. For newborns, none of 359 the poisoned test set scans were correctly classified. 360 Accuracy remained close to 0 for 7-month-old and 8month-old scans. The majority class, 5-month-olds, 362 saw a performance of 35%. 363

U-Net is a model suited for segmentation. It achieved a mean average precision of 42%. This drops to 22% under FGSM and 19% under C&W. 366 It was also trained for 10 epochs and adhered to the same general process that is used to evaluate YOLO. 369

441

³⁷⁰ 6 Knowledge Distillation

Knowledge Distillation is a technique designed to 371 guard against adversarial examples. The effective-372 ness of the previous attacks is due to the inherent 373 fragility of CNNs, where small perturbations can 374 significantly manipulate inference. Knowledge dis-375 tillation involves a student-teacher model approach. 376 In this instance, the teacher model is YOLO. The 377 confidence score of the predicted class is extracted 378 from YOLO and used as a new input. This score, 379 along with the class (in this case, the age of the in-380 fant being scanned), serves as input for a secondary 381 ResNet model. The ResNet is then evaluated both 382 383 with and without distillation.

For the U-Net, instead of feeding the output directly into the student model, it is added to the loss function. For segmentation tasks, basing the loss on ground truth with the teacher model's soft targets supports student model in learning segmentation boundaries from the teacher. The loss used is

 $loss = \alpha \cdot label_loss + (1 - \alpha) \cdot distillation_loss (3)$

This formulation allows for a more precise weighting to be applied to learning the ground truth and agreement with the teacher model.

To test YOLO, 500 adversarial examples were 394 generated using the C&W attack, and 10,000 ad-395 versarial examples were generated using FGSM to 396 create a poisoned version of the original dataset. 397 C&W is too computationally intensive to warrant so 398 many examples, since the adversarial examples are 399 more potent. As for Kvasir, 1,000 examples were 400 created under FGSM with 100 with C&W with the 401 402 same reasoning.

When reviewing the example model without dis-403 tillation, its ability to classify images is very limited. 404 The model is essentially only capable of correctly 405 classifying bones from 4-month-olds and 5-month-406 olds, while misclassifying all other age groups. This 407 suggests that the model is not effectively learning 408 the distinguishing characteristics of the other classes 409 and is instead relying on the prominence of certain 410 classes in making its predictions. 411

In contrast, when applying distillation, the model 412 demonstrates a greater ability to differentiate be-413 tween characteristics across classes. This improve-414 ment is evident in the confusion matrix, which shows 415 a more diverse distribution of predictions. The dis-416 tilled model is attempting to generalize and learn 417 from the data. However, the limitations of this 418 model may indicate that the teacher model (YOLO) 419 has not been sufficiently trained. The YOLO model 420 was only trained for 10 epochs, yet there is still a 421 significant difference in performance between a stan-422 dard ResNet and a ResNet that has undergone dis-423 tillation. However, the dataset is large and is using 424



Figure 3. Comparison of Attacks Under Distillation for Kvasir

a foundational model to conduct transfer learning, 425 so the concern may be overstated. Also, the difference in accuracy is the finding and not the absolute 427 accuracy. 428

When evaluating the effectiveness of distillation 429 in mitigating C&W, the results are not particularly 430 promising. Without distillation, the model struggles to generalize effectively. Although the distilled 432 model performs marginally better, such as predicting 7-month-olds more accurately, it largely fails to 434 cope with the attack. 435

As for the segmentation results, distillation was 436 effective at improving model robustness under C&W. 437 However, it becomes less robust to FGSM and informs a trade-off in how loss functions should be formulated for targeting specific attacks. 440

7 Conclusion

In this study, two attack types were applied to the 442 RSNA Bone Age dataset and Kvasir using FGSM 443 and C&W. FGSM, while simple and effective, is 444 harder to detect, with Benford's Law failing to catch 445 its perturbations. C&W, being more complex, is 446 more damaging but also more noticeable. Knowl-447 edge distillation showed some effectiveness, particu-448 larly against FGSM when considering the attack on 449 YOLO. However, it is ineffective when considering 450 U-Net and FGSM. This is likely attributable to the 451 loss function that is able to more effectively protect 452 against C&W. 453

It is recommended to apply Benford's Law for detecting severe attacks, while knowledge distillation can help mitigate simpler ones. Future work should focus on understanding how loss function formulations factor into distillation effectiveness and the robustness of Benford's Law in the medical image domain to further support robustness initiatives.

461 A Appendix A



Figure A.1. Comparison of ResNets when Attacked by FGSM $\,$



Figure A.2. Comparison of ResNets when Attacked by C&W