

# Certifiable Robustness Against Patch Attacks Using an ERM Oracle

## Abstract

1 Consider patch attacks, where at test-time an adversary manipulates a test image with a patch in order  
2 to induce a targeted mis-classification. We consider a recent defense to patch attacks, Patch-Cleanser  
3 (Xiang et al., 2022). The Patch-Cleanser algorithm requires a prediction model to have a “two-mask  
4 correctness” property, meaning that the prediction model should correctly classify any image when  
5 any two blank masks replace portions of the image. To this end, Xiang et al. (2022) learn a prediction  
6 model to be robust to two-mask operations by augmenting the training set by adding pairs of masks  
7 at random locations of training images, and performing *empirical risk minimization* (ERM) on the  
8 augmented dataset. However, in the non-realizable setting when no predictor is perfectly correct on  
9 *all* two-mask operations on all images, we exhibit an example where ERM fails. To overcome this  
10 challenge, we propose a different algorithm that provably learns a predictor robust to *all* two-mask  
11 operations using an ERM oracle, based on prior work by Feige et al. (2015a).

## 12 1. Introduction

13 Patch attacks (Brown et al., 2017; Karmon et al., 2018; Yang et al., 2020) are an important threat  
14 model in the general field of test time evasion attacks (Goodfellow et al., 2014). In a patch attack, the  
15 adversary replaces a contiguous block of pixels with an adversarially crafted pattern. Patch attacks  
16 can realize physical world attacks to computer vision systems by printing and attaching a patch into  
17 an object. To secure performance of computer vision systems against patch-attacks, there has been  
18 an active line of research for providing certifiable robustness against them (see e.g., McCoyd et al.,  
19 2020; Xiang et al., 2020; Xiang and Mittal, 2021; Metzen and Yatsura, 2021; Zhang et al., 2020;  
20 Chiang et al., 2020).

21 Xiang et al. (2022) recently proposed a state-of-the-art algorithm called Patch-Cleanser that can  
22 provably defend against patch attacks. The high level idea of the Patch-Cleanser algorithm is to  
23 robustly remove all adversarial pixels of an input image in order to obtain accurate predictions. The  
24 main difficulty is that the patch location is unknown. One naive solution is to place a mask at all  
25 possible locations of an input image. As long as the masks are large enough, at least one of the masks  
26 would cover the patch and remove the adversarial effects of the patch so that the prediction model can  
27 induce a correct classification on the input image. However, it is challenging to distinguish between  
28 this correct prediction and the predictions on the other masked images. To overcome this challenge,  
29 they use a second mask. For each of the one-masked images produced in the first step, they add a  
30 second mask at all possible locations. For each one-masked image, if for all possible locations of the  
31 second mask, the prediction model outputs the same classification, it means that the first mask was  
32 removing the patch, and the agreed-upon prediction is correct. Also, any disagreements implies the  
33 contrary.

34 Crucially, the Patch-Cleanser algorithm relies on a *two-mask correctness* assumption of the prediction  
35 model that is defined as follows: for a given input  $(x, y)$ , if for any pair of masks applied to  $x$ , a  
36 prediction model  $F$  outputs the correct prediction  $y$ , then  $F$  has two-correctness property on  $(x, y)$   
37 (see Xiang et al. (2022, Definition 2)). They show as long as the two-mask correctness property holds,  
38 their double-masking algorithm guarantees robustness against patch attacks on the input image  $(x, y)$ .

39 In order to train a model with the two-mask correctness property, Xiang et al. (2022) use a heuristic  
40 data-augmentation approach as follows. They add pairs of masks at random locations to training  
41 images and learn a model that predicts correctly on the masked-images using *empirical risk*  
42 *minimization*.

43 However, we argue that in the non-realizable setting, when no predictor achieves zero robust loss, *this*  
44 *approach can fail*. Intuitively, an ERM oracle does not distinguish between distributing error over  
45 a few perturbations (i.e. masked-out variations) of many input images versus concentrating many  
46 mistakes on the perturbations of few input images. However, in the latter case, the robust loss can be

47 much higher than the former case. To wit, to obtain high robust loss the adversary only needs one  
 48 successful perturbation per clean image. If some input images have many perturbations that fool the  
 49 classifier, but most input images have none, then the adversary cannot obtain high robust loss. We  
 50 want to be in the second case. We have included a schematic demonstrating this failure mode in A.4.

51 Our main contribution is an algorithm to learn a predictor that is robust to a set of masking operations  
 52 (resulting from the two-mask), using an ERM oracle. The algorithm is based on prior work due to  
 53 Feige et al. (2015b), but the analysis and application are novel in this work. Combining our algorithm  
 54 with Patch-Cleanser yields a predictor that is provably robust to adversarial patch attacks.

55 **Setup and Notation** Let  $\mathcal{X}$  denote the instance space and  $\mathcal{Y}$  denote the label space. Our main  
 56 objective is to be robust against adversarial patches  $\mathcal{A} : \mathcal{X} \rightarrow 2^{\mathcal{X}}$ , where  $\mathcal{A}(x)$  represents the  
 57 (potentially infinite) set of adversarially patched images that an adversary might attack with at  
 58 test-time. Xiang et al. (2022) showed that even though the space of adversarial patches  $\mathcal{A}$  can be  
 59 exponential or infinite, one can consider a “covering” set  $\mathcal{U} : \mathcal{X} \rightarrow 2^{\mathcal{X}}$  of masking operations on  
 60 images where  $|\mathcal{U}(x)|$  is polynomially finite.

61 Thus for the remainder of the paper, we focus on the task of learning a predictor robust to a perturbation  
 62 set  $\mathcal{U} : \mathcal{X} \rightarrow 2^{\mathcal{X}}$ , where  $\mathcal{U}(x) \subseteq \mathcal{X}$  is the set of allowed masking operations that can be performed  
 63 on  $x$ . We assume that  $\mathcal{U}(x)$  is finite where  $|\mathcal{U}(x)| \leq k$ . Let  $\mathcal{H} \subseteq \mathcal{Y}^{\mathcal{X}}$  be a hypothesis class, and  
 64 denote by  $\text{vc}(\mathcal{H})$  its VC dimension. Let  $\text{ERM}_{\mathcal{H}}$  be an ERM oracle for  $\mathcal{H}$ . For any set arbitrary set  
 65  $W$ , denote by  $\Delta(W)$  the set of distributions over  $W$ .  $\text{OPT}_{\mathcal{H}}$  is defined as follows:

$$\text{OPT}_{\mathcal{H}} \triangleq \min_{h \in \mathcal{H}} \mathbb{E}_{(x,y) \sim \mathcal{D}} \max_{z \in \mathcal{U}(x)} \mathbb{1}[h(z) \neq y]. \quad (1)$$

## 66 2. Main Result: Minimizing Robust Loss Using an ERM Oracle

67 In this section, we present our main contribution: an algorithm to learn a predictor that is simulta-  
 68 neously robust to a set of (polynomially many) masking operations, using an  $\text{ERM}_{\mathcal{H}}$  oracle. The  
 69 algorithm is based on prior work due to Feige et al. (2015b), but the analysis and application are novel  
 70 in this work. The main interesting feature of this algorithm is that it achieves stronger robustness  
 71 guarantees in the non-realizable regime when  $\text{OPT}_{\mathcal{H}} \gg 0$ , where the approach of Xiang et al. (2022)  
 72 — as we highlighted in the introduction — of calling  $\text{ERM}_{\mathcal{H}}$  on the inflated dataset: original training  
 73 points plus all possible perturbations resulting from the allowed masking operations, can provably  
 74 fail (see e.g., A.4).

---

**Algorithm 1:** Feige, Mansour, and Schapire (2015b)

---

**Input:** weight update parameter  $\eta > 0$ , number of rounds  $T$ , and training dataset

$$S = \{(x_1, y_1), \dots, (x_m, y_m)\}.$$

1 Set  $w_1(z, (x, y)) = 1$ , for each  $(x, y) \in S, z \in \mathcal{U}(x)$ .

2 Set  $P^1(z, (x, y)) = \frac{w_1(z, (x, y))}{\sum_{z' \in \mathcal{U}(x)} w_1(z', (x, y))}$ , for each  $(x, y) \in S, z \in \mathcal{U}(x)$ .

3 **for each**  $t \leftarrow 1$  **to**  $T$  **do**

75 4 Call ERM on the empirical weighted distribution:

$$h_t = \operatorname{argmin}_{h \in \mathcal{H}} \sum_{(x,y) \in S} \sum_{z \in \mathcal{U}(x)} \frac{1}{m} P^t(z, (x, y)) \mathbb{1}[h_t(z) \neq y].$$

5 **for each**  $(x, y) \in S$  **and**  $z \in \mathcal{U}(x)$  **do**

6  $w_{t+1}(z, (x, y)) = (1 + \eta \mathbb{1}[h_t(z) \neq y]) \cdot w_t(z, (x, y)).$

7  $P^{t+1}(z, (x, y)) = \frac{w_t(z, (x, y))}{\sum_{z' \in \mathcal{U}(x)} w_t(z', (x, y))}.$

**Output:** The majority-vote predictor  $\text{MAJ}(h_1, \dots, h_T)$ .

---

76 **Theorem 1.** Set  $T(\epsilon) = \frac{32 \ln k}{\epsilon^2}$  and  $m(\epsilon, \delta) = O\left(\frac{\text{vc}(\mathcal{H})(\ln k)^2}{\epsilon^4} \ln\left(\frac{\ln k}{\epsilon^2}\right) + \frac{\ln(1/\delta)}{\epsilon^2}\right)$ . Then, for any  
 77 distribution  $\mathcal{D}$  over  $\mathcal{X} \times \mathcal{Y}$ , with probability at least  $1 - \delta$  over  $S \sim \mathcal{D}^{m(\epsilon, \delta)}$ , running Algorithm 1  
 78 for  $T(\epsilon)$  rounds produces  $h_1, \dots, h_{T(\epsilon)}$  satisfying:

$$\mathbb{E}_{(x,y) \sim \mathcal{D}} \left[ \max_{z \in \mathcal{U}(x)} \mathbb{1}[\text{MAJ}(h_1, \dots, h_{T(\epsilon)})(z) \neq y] \right] \leq 2\text{OPT}_{\mathcal{H}} + \epsilon.$$

79 **Comparison with prior related work** As presented, Feige et al. (2015b) only considered *finite*  
80 hypothesis classes  $\mathcal{H}$  and provided generalization guarantees depending on  $\log |\mathcal{H}|$ . On the other  
81 hand, we consider here infinite classes  $\mathcal{H}$  with bounded VC dimension and provide tighter robust  
82 generalization bounds. The robust learning guarantee (Attias et al., 2022, Theorem 2) assumes access  
83 to a *robust* ERM oracle, which minimizes the robust loss on the training dataset. On the other hand,  
84 at the expense of higher sample complexity, we provide a robust learning guarantee using only an  
85 ERM oracle in the challenging *non-realizable* setting. Prior work due to Montasser et al. (2020)  
86 considered using an ERM oracle for robust learning but only in the simpler realizable setting (when  
87  $\text{OPT}_{\mathcal{H}} = 0$ ).

88 Before proceeding with the proof Theorem 1, we describe now at a high-level the proof strategy. The  
89 main insight is to solve a finite zero-sum game. In particular, our goal is to find a mixed-strategy over  
90 the hypothesis class that is approximately close to the value of the game:

$$\text{OPT}_{S, \mathcal{H}} \triangleq \min_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m \max_{z_i \in \mathcal{U}(x_i)} \mathbb{1}[h(z_i) \neq y_i].$$

91 We observe that Algorithm 1 due to (Feige et al., 2015b) solves a similar finite zero-sum game  
92 (see Lemma 3), and then we relate it to the value of the game we are interested in (see Lemma  
93 2). Combined together, this only establishes that we can minimize the robust loss on the empirical  
94 dataset using an ERM oracle. We then appeal to uniform convergence guarantees for the robust loss  
95 in Lemma 4 to show that, with large enough training data, our output predictor achieves robust risk  
96 that is close to the value of the game.

97 **Lemma 2.** For any data set  $S = \{(x_1, y_1), \dots, (x_m, y_m)\} \in (\mathcal{X} \times \mathcal{Y})^m$ ,

$$\begin{aligned} \text{OPT}_{S, \mathcal{H}} &= \min_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m \max_{z_i \in \mathcal{U}(x_i)} \mathbb{1}[h(z_i) \neq y_i] \geq \\ &\quad \min_{Q \in \Delta(\mathcal{H})} \max_{P_1 \in \Delta(\mathcal{U}(x_1)), \dots, P_m \in \Delta(\mathcal{U}(x_m))} \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{z_i \sim P_i} \mathbb{E}_{h \sim Q} \mathbb{1}[h(z_i) \neq y_i]. \end{aligned}$$

98 **Lemma 3** (Feige, Mansour, and Schapire (2015b)). For any data set  $S =$   
99  $\{(x_1, y_1), \dots, (x_m, y_m)\} \in (\mathcal{X} \times \mathcal{Y})^m$ , running Algorithm 1 for  $T$  rounds produces a mixed-strategy  
100  $\hat{Q} = \frac{1}{T} \sum_{t=1}^T h_t \in \Delta(\mathcal{H})$  satisfying:

$$\begin{aligned} &\max_{P_1 \in \Delta(\mathcal{U}(x_1)), \dots, P_m \in \Delta(\mathcal{U}(x_m))} \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{z_i \sim P_i} \frac{1}{T} \sum_{t=1}^T \mathbb{1}[h_t(z_i) \neq y_i] \leq \\ &\quad \min_{Q \in \Delta(\mathcal{H})} \max_{P_1 \in \Delta(\mathcal{U}(x_1)), \dots, P_m \in \Delta(\mathcal{U}(x_m))} \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{z_i \sim P_i} \mathbb{E}_{h \sim Q} \mathbb{1}[h(z_i) \neq y_i] + 2\sqrt{\frac{\ln k}{T}}. \end{aligned}$$

101 **Lemma 4** (VC Dimension for the Robust Loss (Attias et al., 2022)). For any class  $\mathcal{H}$  and any  $\mathcal{U}$   
102 such that  $\sup_{x \in \mathcal{X}} |\mathcal{U}(x)| \leq k$ , denote the robust loss class of  $\mathcal{H}$  with respect to  $\mathcal{U}$  by

$$\mathcal{L}_{\mathcal{H}}^{\mathcal{U}} = \left\{ (x, y) \mapsto \max_{z \in \mathcal{U}(x)} \mathbb{1}[h(z) \neq y] : h \in \mathcal{H} \right\}.$$

103 Then, it holds that  $\text{vc}(\mathcal{L}_{\mathcal{H}}^{\mathcal{U}}) \leq O(\text{vc}(\mathcal{H}) \log(k))$ .

104 We are now ready to proceed with the proof of Theorem 1.

105 **Proof** Let  $S \sim \mathcal{D}^m$  be an iid sample from  $\mathcal{D}$ , where the size of the sample  $m$  will be determined  
106 later. By invoking Lemma 3 and Lemma 2, we observe that running Algorithm 1 on  $S$  for  $T$  rounds,  
107 produces  $h_1, \dots, h_T$  satisfying

$$\max_{P_1 \in \Delta(\mathcal{U}(x_1)), \dots, P_m \in \Delta(\mathcal{U}(x_m))} \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{z_i \sim P_i} \frac{1}{T} \sum_{t=1}^T \mathbb{1}[h_t(z_i) \neq y_i] \leq \text{OPT}_{S, \mathcal{H}} + \frac{\epsilon}{4}.$$

108 Next, the average robust loss for the majority-vote predictor  $\text{MAJ}(h_1, \dots, h_T)$  can be bounded from  
 109 above as follows

$$\begin{aligned}
 & \frac{1}{m} \sum_{i=1}^m \max_{z_i \in \mathcal{U}(x_i)} \mathbb{1} [\text{MAJ}(h_1, \dots, h_T)(z_i) \neq y_i] \\
 & \leq \frac{1}{m} \sum_{i=1}^m \max_{z_i \in \mathcal{U}(x_i)} 2 \mathbb{E}_{t \sim [T]} \mathbb{1} [h_t(z_i) \neq y_i] \\
 & = 2 \frac{1}{m} \sum_{i=1}^m \max_{z_i \in \mathcal{U}(x_i)} \frac{1}{T} \sum_{t=1}^T \mathbb{1} [h_t(z_i) \neq y_i] \\
 & \leq 2 \max_{P_1 \in \Delta(\mathcal{U}(x_1)), \dots, P_m \in \Delta(\mathcal{U}(x_m))} \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{z_i \sim P_i} \frac{1}{T} \sum_{t=1}^T \mathbb{1} [h_t(z_i) \neq y_i] \\
 & \leq 2 \text{OPT}_{S, \mathcal{H}} + \frac{\epsilon}{2}.
 \end{aligned}$$

110 Next, we invoke Lemma 4 to obtain a uniform convergence guarantee on the robust loss. In particular,  
 111 we apply Lemma 4 on the *convex-hull* of  $\mathcal{H}$ :  $\mathcal{H}^T = \{\text{MAJ}(h_1, \dots, h_T) : h_1, \dots, h_T \in \mathcal{H}\}$ . By a  
 112 classic result due to Blumer, Ehrenfeucht, Haussler, and Warmuth (1989), it holds that  $\text{vc}(\mathcal{H}^T) =$   
 113  $O(\text{vc}(\mathcal{H})T \ln T)$ . Combining this with Lemma 4 and plugging-in the value of  $T = \frac{32 \ln k}{\epsilon^2}$ , we get  
 114 that the VC dimension of the robust loss class of  $\mathcal{H}^T$  is bounded from above by

$$\text{vc}(\mathcal{L}_{\mathcal{H}^T}^{\mathcal{U}}) \leq O\left(\frac{\text{vc}(\mathcal{H})(\ln k)^2}{\epsilon^2} \ln\left(\frac{\ln k}{\epsilon^2}\right)\right).$$

115 Finally, using Vapnik’s “General Learning” uniform convergence (Vapnik, 1982), with probability at  
 116 least  $1 - \delta$  over  $S \sim \mathcal{D}^m$  where  $m = O\left(\frac{\text{vc}(\mathcal{H})(\ln k)^2}{\epsilon^4} \ln\left(\frac{\ln k}{\epsilon^2}\right) + \frac{\ln(1/\delta)}{\epsilon^2}\right)$ , it holds that

$$\forall f \in \mathcal{H}^T : \mathbb{E}_{(x, y) \sim \mathcal{D}} \left[ \max_{z \in \mathcal{U}(x)} \mathbb{1} [f(z) \neq y] \right] \leq \frac{1}{m} \sum_{i=1}^m \max_{z_i \in \mathcal{U}(x_i)} \mathbb{1} [f(z_i) \neq y_i] + \frac{\epsilon}{4}.$$

117 This also applies to the particular output  $\text{MAJ}(h_1, \dots, h_T)$  of Algorithm 1, and thus

$$\begin{aligned}
 \mathbb{E}_{(x, y) \sim \mathcal{D}} \left[ \max_{z \in \mathcal{U}(x)} \mathbb{1} [\text{MAJ}(h_1, \dots, h_{T(\epsilon)})(z) \neq y] \right] & \leq \frac{1}{m} \sum_{i=1}^m \max_{z_i \in \mathcal{U}(x_i)} \mathbb{1} [\text{MAJ}(h_1, \dots, h_T)(z_i) \neq y_i] + \frac{\epsilon}{4} \\
 & \leq 2 \text{OPT}_{S, \mathcal{H}} + \frac{\epsilon}{2} + \frac{\epsilon}{4}.
 \end{aligned}$$

118 Finally, by applying a standard Chernoff-Hoeffding concentration inequality, we get that  
 119  $\text{OPT}_{S, \mathcal{H}} \leq \text{OPT}_{\mathcal{H}} + \frac{\epsilon}{8}$ . Combining this with the above inequality concludes the proof.  $\blacksquare$

120

### 121 3. Conclusion

122 Per the call for papers, we discuss the scalability of our method, which depends on multiple factors.  
 123 As is argued/demonstrated empirically in the original Patch-Cleanser paper, their original defense  
 124 scales to high resolution images. An additional strength of this research direction initiated by Patch-  
 125 Cleanser and maintained by our approach is that it is *agnostic* to the network/structure of the model,  
 126 and can be applied as a module on top of any state-of-the-art model. The complexity of Algorithm 1  
 127 has two components, the complexity of the ERM Oracle and the number of iterations  $T$ . As noted  
 128 in Section 2, our algorithm makes  $T = \Omega\left(\frac{\ln k}{\epsilon^2}\right)$  oracle calls where  $\ln k$  is the bit-complexity of the  
 129 perturbations and thus we are oracle-efficient.

130 In order to modify Xiang et al. (2022) to handle the case where two-mask correctness *is not realizable*,  
 131 we exhibit polynomial time algorithms for learning a classifier that satisfies the two-mask property  
 132 and analyze the provable robustness of this approach, based upon prior work by Feige et al. (2015a).  
 133 The key future work that we intend for the full version of this work includes an empirical evaluation  
 134 of this method and extensions to a new multi-group robustness notion.

135 **Appendix A. Missing Proofs**

136 **A.1 Proof of Lemma 2**

137 **Proof** By definition of  $\text{OPT}_{S, \mathcal{H}}$ , it follows that

$$\begin{aligned}
\text{OPT}_{S, \mathcal{H}} &= \min_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m \max_{z_i \in \mathcal{U}(x_i)} \mathbb{1}[h(z_i) \neq y_i] \\
&\geq \min_{h \in \mathcal{H}} \max_{z_1 \in \mathcal{U}(x_1), \dots, z_m \in \mathcal{U}(x_m)} \frac{1}{m} \sum_{i=1}^m \mathbb{1}[h(z_i) \neq y_i] \\
&\geq \min_{Q \in \Delta(\mathcal{H})} \max_{z_1 \in \mathcal{U}(x_1), \dots, z_m \in \mathcal{U}(x_m)} \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{h \sim Q} \mathbb{1}[h(z_i) \neq y_i] \\
&\geq \min_{Q \in \Delta(\mathcal{H})} \max_{P_1 \in \Delta(\mathcal{U}(x_1)), \dots, P_m \in \Delta(\mathcal{U}(x_m))} \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{z_i \sim P_i} \mathbb{E}_{h \sim Q} \mathbb{1}[h(z_i) \neq y_i].
\end{aligned}$$

138  
139

140 **A.2 Proof of Lemma 3**

141 **Proof** By the minimax theorem and (Feige, Mansour, and Schapire, 2015b, Equation 3 and 9 in  
142 proof of Theorem 1), we have that

$$\begin{aligned}
\max_{P_1 \in \Delta(\mathcal{U}(x_1)), \dots, P_m \in \Delta(\mathcal{U}(x_m))} \sum_{i=1}^m \mathbb{E}_{z_i \sim P_i} \frac{1}{T} \sum_{t=1}^T \mathbb{1}[h_t(z_i) \neq y_i] \leq \\
\min_{Q \in \Delta(\mathcal{H})} \max_{P_1 \in \Delta(\mathcal{U}(x_1)), \dots, P_m \in \Delta(\mathcal{U}(x_m))} \mathbb{E}_{z_i \sim P_i} \mathbb{E}_{h \sim Q} \mathbb{1}[h(z_i) \neq y_i] + 2 \frac{\sqrt{\mathcal{L}^* m \ln k}}{T},
\end{aligned}$$

143 where  $\mathcal{L}^* = \sum_{i=1}^m \max_{z \in \mathcal{U}(x_i)} \sum_{t=1}^T \mathbb{1}[h_t(z) \neq y]$ . By observing that  $\mathcal{L}^* \leq mT$  and dividing both  
144 sides of the inequality above by  $m$ , we arrive at the inequality stated in the lemma.  $\blacksquare$

145

146 **A.3 Proof of Lemma 4**

**Proof** By finiteness of  $\mathcal{U}$ , observe that for any dataset  $S \in (\mathcal{X} \times \mathcal{Y})^m$ , each robust loss vector in the set of robust loss behaviors:

$$\begin{aligned}
\Pi_{\mathcal{L}_{\mathcal{H}}^{\mathcal{U}}}(S) &= \{(f(x_1, y_1), \dots, f(x_m, y_m)) : f \in \mathcal{L}_{\mathcal{H}}^{\mathcal{U}}\} \\
\text{maps to a 0-1 loss vector on the inflated set } S_{\mathcal{U}} &= \\
&= \{(z_1^1, y_1), \dots, (z_1^k, y_1), (z_2^1, y_2), \dots, (z_2^k, y_2), \dots, (z_m^1, y_m), \dots, (z_m^k, y_m)\}, \\
\Pi_{\mathcal{H}}(S_{\mathcal{U}}) &= \{(h(z_1^1), \dots, h(z_1^k), h(z_2^1), \dots, h(z_2^k), \dots, h(z_m^1), \dots, h(z_m^k)) : h \in \mathcal{H}\}.
\end{aligned}$$

147 Therefore, it follows that  $|\Pi_{\mathcal{L}_{\mathcal{H}}^{\mathcal{U}}}(S)| \leq |\Pi_{\mathcal{H}}(S_{\mathcal{U}})|$ . Then, by applying the Sauer-Shelah lemma, it  
148 follows that  $|\Pi_{\mathcal{H}}(S_{\mathcal{U}})| \leq O((mk)^{\text{vc}(\mathcal{H})})$ . Then, by solving for  $m$  such that  $O((mk)^{\text{vc}(\mathcal{H})}) \leq 2^m$ ,  
149 we get that  $\text{vc}(\mathcal{L}_{\mathcal{H}}^{\mathcal{U}}) \leq O(\text{vc}(\mathcal{H}) \log(k))$ .  $\blacksquare$

150

151 **A.4 ERM failure Example**

152 **Appendix B. Risk Analysis**

153 Per the call for papers for this workshop, in this section we will include our risk analysis, which is a  
154 novel contribution for the authors. Some of this analysis is general to theory papers in robustness and

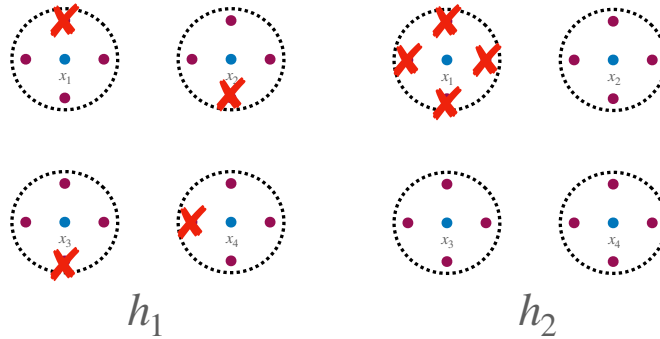


Figure 1:  $\{x_i\}$  are natural data points in blue and each is surrounded by their perturbed points in purple. Red means mis-classified points. The approach in this section is to solve the ERM on an inflated data-set consisting of natural points and their perturbations. Observe that both  $h_1$  and  $h_2$  have the same 0 – 1 loss on this toy data-set with four points but  $h_2$  has much lower robust loss since it can correctly classify 3/4 of the original examples no matter what the adversary does, while for  $h_1$  the adversary can perturb any point to induce a mis-classification

155 some of it is specific to our work. Some of our risk analysis is based on discussion in Hendrycks and  
 156 Mazeika (2022). This work attempts to mitigate existing risks due to patch attacks.

157 **B.1 Short Term Risk of Patch Attacks**

158 First, some discussion of the short term vulnerability of learned systems. Adversarial attacks of this  
 159 nature lend themselves immediately to targeted attacks by malicious actors. For instance, if adversarial  
 160 patch attacks remain a systemic flaw of vision models, and self driving cars with vulnerable vision  
 161 systems are widely deployed, malefactors could dangerously target specific vehicles. Alternatively,  
 162 if software to design universal adversarial patches continues to proliferate, then lone wolves could  
 163 spread patches widely without a specific target and pose an acute and hard to mitigate risk to any  
 164 driver.

165 In addition to acute harms caused by attackers using these systems, they could also delay or prevent  
 166 the beneficial use of AI systems. This type of vulnerability could limit the reliance of automakers  
 167 on vision systems or delay the implementation of self driving technology. While some of the safety  
 168 benefits of self driving technology remain conjectural, something on the order of 50,000 Americans  
 169 die per year in automotive accidents (NHTSA) and on the order of 1 million people annually WHO.  
 170 There is a plausible argument that self driving technology can mitigate these risks by achieving  
 171 super-human performance and consistent driving behavior.

172 Our work mitigates some of this patch risk by exhibiting an algorithm that can learn a classifier that  
 173 competes with the global optima for this problem.

174 The most important limitation currently is we have not yet implemented an empirical evaluation. This  
 175 is intended for a future version of this work.

176 **B.2 Long Term Risk of Patch Attacks**

177 We observe a key long term concern, the risk of catastrophic failures due to AI/ML based controllers  
 178 subject to patch attacks or other perturbations. For instance while the author was writing this, they  
 179 observed an advertisement in an undisclosed airport for ‘AI for Air Traffic Control’. Safety critical

180 systems increasingly have possibly vulnerable AI sub-systems. There are some efforts to integrate  
181 AI/ML techniques into nuclear command and control systems Lowther and McGiffin (2019).

182 Some of the hypothetical benefits of AI include possibly simplifying decision making for human actors  
183 by reducing information overload and giving them the time to make thoughtful choices Lowther and  
184 McGiffin (2019). If AI systems are too vulnerable Klare (2020), these benefits would be unrealized  
185 and some decision makers may remain stuck with sub-optimal choices. If patch attacks/adversarial  
186 attacks remain a credible threat and these control systems are deployed, that could have extreme  
187 consequences. For instance, if an early launch warning system had a satellite based vision component  
188 focused on missile silos, a patch attack could prevent early detection.

189 Alternatively, there is also a risk to continuing to use legacy and non-AI systems in that we may be  
190 stuck with poor human decision making or static systems.

191 Moving out in terms of generality, there are also questions raised by adversarial robustness about  
192 whether or not models can be relied upon to perform consistently, when subject to natural perturbations  
193 of distribution shift, and our work is progress in this direction.

## 194 **References**

- 195 Idan Attias, Aryeh Kontorovich, and Yishay Mansour. Improved generalization bounds for adversari-  
196 ally robust learning. *Journal of Machine Learning Research*, 23(175):1–31, 2022.
- 197 A. Blumer, A. Ehrenfeucht, D. Haussler, and M. Warmuth. Learnability and the Vapnik-Chervonenkis  
198 dimension. *Journal of the Association for Computing Machinery*, 36(4):929–965, 1989.
- 199 Tom B Brown, Dandelion Mané, Aurko Roy, Martín Abadi, and Justin Gilmer. Adversarial patch.  
200 *arXiv preprint arXiv:1712.09665*, 2017.
- 201 Ping-yeh Chiang, Renkun Ni, Ahmed Abdelkader, Chen Zhu, Christoph Studer, and Tom Goldstein.  
202 Certified defenses for adversarial patches. *CoRR*, abs/2003.06693, 2020. URL <https://arxiv.org/abs/2003.06693>.
- 204 Uriel Feige, Yishay Mansour, and Robert Schapire. Learning and inference in the presence of  
205 corrupted inputs. In *Conference on Learning Theory*, pages 637–657. PMLR, 2015a.
- 206 Uriel Feige, Yishay Mansour, and Robert E. Schapire. Learning and inference in the presence of  
207 corrupted inputs. In Peter Grünwald, Elad Hazan, and Satyen Kale, editors, *Proceedings of The  
208 28th Conference on Learning Theory, COLT 2015, Paris, France, July 3-6, 2015*, volume 40  
209 of *JMLR Workshop and Conference Proceedings*, pages 637–657. JMLR.org, 2015b. URL  
210 <http://proceedings.mlr.press/v40/Feige15.html>.
- 211 Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial  
212 examples. *arXiv preprint arXiv:1412.6572*, 2014.
- 213 Dan Hendrycks and Mantas Mazeika. X-risk analysis for ai research, 2022. URL <https://arxiv.org/abs/2206.05862>.
- 215 Danny Karmon, Daniel Zoran, and Yoav Goldberg. Lavan: Localized and visible adversarial noise.  
216 In *International Conference on Machine Learning*, pages 2507–2515. PMLR, 2018.
- 217 Michael T Klare. ‘skynet’ revisited: The dangerous allure of nuclear command  
218 automation, 2020. URL [https://www.armscontrol.org/act/2020-04/features/  
219 skynet-revisited-dangerous-allure-nuclear-command-automation](https://www.armscontrol.org/act/2020-04/features/skynet-revisited-dangerous-allure-nuclear-command-automation).
- 220 Adam Lowther and Curtis McGiffin. MS Windows NT kernel description, 2019. URL <https://warontherocks.com/2019/08/america-needs-a-dead-hand/>.
- 222 Michael McCoyd, Won Park, Steven Chen, Neil Shah, Ryan Roggenkemper, Minjune Hwang,  
223 Jason Xinyu Liu, and David A. Wagner. Minority reports defense: Defending against adversarial  
224 patches. *CoRR*, abs/2004.13799, 2020. URL <https://arxiv.org/abs/2004.13799>.
- 225 Jan Hendrik Metzen and Maksym Yatsura. Efficient certified defenses against patch attacks on image  
226 classifiers. *CoRR*, abs/2102.04154, 2021. URL <https://arxiv.org/abs/2102.04154>.

- 227 Omar Montasser, Steve Hanneke, and Nati Srebro. Reducing adversarially robust learning to non-  
228 robust PAC learning. In Hugo Larochelle, Marc’ Aurelio Ranzato, Raia Hadsell, Maria-Florina  
229 Balcan, and Hsuan-Tien Lin, editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. URL <https://proceedings.neurips.cc/paper/2020/hash/a822554e5403b1d370db84cfbc530503-Abstract.html>.
- 233 NHTSA. "newly released estimates show traffic fatalities reached a 16-  
234 year high in 2021". URL <https://www.nhtsa.gov/press-releases/early-estimate-2021-traffic-fatalities>.
- 236 V. Vapnik. *Estimation of Dependencies Based on Empirical Data*. Springer-Verlag, New York, 1982.
- 237 WHO. "road traffic injures". URL <https://www.who.int/news-room/fact-sheets/detail/road-traffic-injuries>.
- 239 Chong Xiang and Prateek Mittal. Patchguard++: Efficient provable attack detection against adversarial  
240 patches. *CoRR*, abs/2104.12609, 2021. URL <https://arxiv.org/abs/2104.12609>.
- 241 Chong Xiang, Arjun Nitin Bhagoji, Vikash Sehwal, and Prateek Mittal. Patchguard: Provable  
242 defense against adversarial patches using masks on small receptive fields. *CoRR*, abs/2005.10884,  
243 2020. URL <https://arxiv.org/abs/2005.10884>.
- 244 Chong Xiang, Saeed Mahloujifar, and Prateek Mittal. Patchcleanser: Certifiably robust defense  
245 against adversarial patches for any image classifier. In *31st USENIX Security Symposium (USENIX Security 22)*, pages 2065–2082, 2022.
- 247 Chenglin Yang, Adam Kortylewski, Cihang Xie, Yinzhi Cao, and Alan Yuille. Patchattack: A  
248 black-box texture-based attack with reinforcement learning. In *European Conference on Computer Vision*, pages 681–698. Springer, 2020.
- 250 Zhanyuan Zhang, Benson Yuan, Michael McCoyd, and David Wagner. Clipped bagnet: Defending  
251 against sticker attacks with clipped bag-of-features. In *2020 IEEE Security and Privacy Workshops (SPW)*, pages 55–61. IEEE, 2020.