

# Jacobian Ensembles Improve Robustness Trade-offs to Adversarial Attacks\*

Kenneth T. Co<sup>1,2</sup>[0000-0003-2766-7326], David  
Martinez-Rego<sup>2</sup>[0000-0003-1809-1169], Zhongyuan Hau<sup>1</sup>, and Emil C.  
Lupu<sup>1</sup>[0000-0002-2844-3917]

<sup>1</sup> Imperial College London, London SW7 2AZ, United Kingdom

{k.co, zy.hau17, e.c.lupu}@imperial.ac.uk

<sup>2</sup> DataSpartan, London EC2Y 9ST, United Kingdom

david@dataspartan.com

**Abstract.** Deep neural networks have become an integral part of our software infrastructure and are being deployed in many widely-used and safety-critical applications. However, their integration into many systems also brings with it the vulnerability to test time attacks in the form of Universal Adversarial Perturbations (UAPs). UAPs are a class of perturbations that when applied to *any input* causes model misclassification. Although there is an ongoing effort to defend models against these adversarial attacks, it is often difficult to reconcile the trade-offs in model accuracy and robustness to adversarial attacks. Jacobian regularization has been shown to improve the robustness of models against UAPs, whilst model ensembles have been widely adopted to improve both predictive performance and model robustness. In this work, we propose a novel approach, Jacobian Ensembles – a combination of Jacobian regularization and model ensembles to significantly increase the robustness against UAPs whilst maintaining or improving model accuracy. Our results show that Jacobian Ensembles achieves previously unseen levels of accuracy and robustness, greatly improving over previous methods that tend to skew towards only either accuracy or robustness.

**Keywords:** Adversarial machine learning · Computer vision · Jacobian regularization · Ensemble methods.

## 1 Introduction

Deep neural networks (DNNs) have achieved widespread use in many applications including image classification [15], real-time object detection [20], and speech recognition [12]. Despite these advances, there is an increasing recognition that DNNs are exposed to systemic vulnerabilities in the form of Universal Adversarial Perturbations (UAPs): where a single adversarial perturbation causes a model to misclassify a large set of inputs [19]. Thus, it is important to ensure that

---

\* Kenneth T. Co is supported in part by the DataSpartan research grant DSRD201801.

neural networks are robust to such devastating attacks whilst still maintaining their state-of-art accuracy on benign datasets.

UAP present a systemic risk, as they enable practical and physically realizable adversarial attacks. They have been demonstrated in many widely-used and safety-critical applications such as camera-based computer vision [7,8,2,18] and LiDAR-based object detection [10,11]. UAPs have also been shown to facilitate realistic attacks in both the physical [23] and digital [24] domains. In some cases, UAPs also enable very resource-efficient black-box attacks on DNNs [3,5].

An insufficiently studied aspect of existing defenses against UAPs is the trade-off between clean accuracy, or the model’s performance on a benign dataset, and its robustness to adversarial attacks. Indeed, a model with increased robustness to UAPs is desirable, but reduced clean accuracy could translate to reduced utility and additional financial or security costs depending on the application. Existing defenses primarily consider robustness to adversarial attacks, but neglect the cost it incurs on the model’s performance for the original task. For example, defenses like adversarial training or too much regularization improve robustness but greatly reduce clean accuracy [6,13,21].

*Jacobian regularization* (JR) has previously been shown to improve robustness against UAPs. However, JR can damage clean accuracy for large amounts of regularization [6,13,21]. *Model ensembles* on the other hand has widely been shown to achieve better classification performance and stability than a single (best) classifier [16,26]. Ensembles are created by combining the outputs of multiple base learners to generate an improved prediction.

In this work, we propose combining JR and model ensembles during training to create *Jacobian Ensembles*. JR is used to drastically improve the model’s robustness to UAPs while the ensemble methods stabilize the model’s predictions and improve its clean accuracy. JR and model ensembles individually each have been shown to improve UAP robustness but at some cost to clean accuracy [4,6]. We show that Jacobian Ensembles greatly improve on the accuracy-robustness trade-off when compared to either JR or model ensembles individually. First, we theoretically show that increasing the number of base learners in a model ensemble improves the expected robustness of classifiers. Then, we empirically verify our theoretical findings by applying JR with popular ensemble methods bagging [1], snapshot ensembles [14], soft voting [26] to DNNs trained on the popular benchmark datasets: MNIST [17], Fashion-MNIST [25] and then evaluating their robustness against UAPs.

To summarize, we make the following contributions:

- We derive theoretical formulations for robustness of ensemble methods and show that the robustness to UAPs increases monotonically with the number of base learners.
- We empirically verify our theoretical results and show that Jacobian Ensembles, a combination of Jacobian regularization and ensembles, achieves the best accuracy-robustness trade-off as measured by a combined weighted accuracy metric.

## 2 Background

### 2.1 Universal Adversarial Perturbations

Let  $f : \mathcal{X} \subset \mathbb{R}^D \rightarrow \mathbb{R}^C$  be logits of a piece-wise linear classifier with input  $\mathbf{x} \in \mathcal{X}$ . We, define  $F(\mathbf{x}) = \arg \max(f(\mathbf{x}))$  to be the output of this classifier and write  $\tau(\mathbf{x})$  as the true class label of an input  $\mathbf{x}$ . **Universal Adversarial Perturbations (UAP)** are perturbations  $\delta \in \mathbb{R}^n$  to the data that satisfy  $F(\mathbf{x} + \delta) \neq \tau(\mathbf{x})$  for sufficiently many  $\mathbf{x} \in \mathcal{X}$  where  $\|\delta\|_p < \varepsilon$ . The latter condition  $\|\delta\|_p < \varepsilon$  constrains the magnitude of the perturbation and is often some  $\ell_p$ -norm and small  $\varepsilon > 0$  [19]. Given a classifier  $f$ , UAPs are generated by maximizing the loss  $\sum_i \mathcal{L}_f(\mathbf{x}_i + \delta)$  with an iterative stochastic gradient descent algorithm [4,22] where  $\mathcal{L}$  is the model’s training loss,  $\{\mathbf{x}_i\}$  are batches of inputs, and  $\delta$  are small perturbations that satisfy  $\|\delta\|_p < \varepsilon$ .

### 2.2 Model Ensembles

An ensemble consists of combining multiple classifiers (base learners) to obtain a resulting ensemble that has better accuracy or predictive performance on aggregate than any individual base learner. In practice, it is widely accepted that combining multiple classifiers can achieve better classification performance than a single “best” classifier [16,26]. Ensembles are typically generated in two ways: sequentially and in parallel. After generating the base learners the combination of their outputs is taken rather than choosing a single “best” learner [26].

In this work, we will analyze ensemble methods that aggregate their base learners in a convex combination. Formally, we define an ensemble  $\mathcal{F}$  as a convex combination of  $M$  base learners  $f_i$ :  $\mathcal{F}(x) = \sum_{i=1}^M c_i f_i(x)$  where  $\sum_{i=1}^M c_i = 1$  and  $0 < c_i < 1, \forall i$ . This is typical as many ensemble methods aggregate their methods via averaging or a similar form of weighted sum [26]. Note however that this will exclude some boosting algorithms such as AdaBoost that are typically not convex combinations [9].

Popular algorithms like *Bagging* [1] and newer methods like *Snapshot Ensembles* [14] take the average of their base learners. Other methods like *Soft Voting* [26] use a convex combination of their model outputs to vote. We refer the reader to each ensemble methods’ corresponding paper for further details on how the base learners are generated. For this work, we will consider Bagging, Snapshot Ensembles, and Soft Voting.

### 2.3 Jacobian Regularization

Let  $f(\mathbf{x})$  be the logit output of the classifier for input  $\mathbf{x}$ , we write  $\mathbf{J}_f(\mathbf{x})$  to denote the input-output Jacobian of  $f$  at  $\mathbf{x}$ . To train models with Jacobian regularization (JR) [13,6], the following joint loss is optimized:

$$\mathcal{L}_{\text{joint}}(\theta) = \mathcal{L}_{\text{train}}(\{\mathbf{x}_i, \mathbf{y}_i\}_i, \theta) + \frac{\lambda_{\text{JR}}}{2} \left( \frac{1}{B} \sum_i \|\mathbf{J}(\mathbf{x}_i)\|_F^2 \right) \quad (1)$$

where  $\theta$  represent the parameters of the model,  $\mathcal{L}_{\text{train}}$  is the standard cross-entropy training loss,  $\{\mathbf{x}_i, \mathbf{y}_i\}$  are input-output pairs from the mini-batch, and  $B$  is the mini-batch size. This optimization uses a regularization parameter  $\lambda_{\text{JR}}$ , which allows the adjustment between regularization and classification loss.

The primary idea is to reduce the Frobenius norm of the input-output Jacobian  $\|\mathbf{J}(\mathbf{x}_i)\|_F$  to decrease the model’s sensitivity to small perturbations such as UAPs. JR shows some promise in improving robustness to UAPs. However, it can often simultaneously decrease the model’s clean accuracy [6] especially for large values of  $\lambda_{\text{JR}}$ .

### 3 Bounds on UAP Effectiveness for Model Ensembles

In this section, we derive theoretical bounds for the expectation and variance of the Frobenius norm of the Jacobian of model ensembles. Note that we only consider ensembles that take a convex combination of their base learners. Similar to [6], we restrict the Frobenius norm of the Jacobian of a model to improve robustness against UAPs.

We show that using model ensembles result in tighter bounds on the Frobenius norm of the Jacobian, suggesting improved robustness and stability versus a single classifier. Our main result in **Theorem 1** shows that increasing the number of base learners decreases both the upper and lower bounds of the expectation and variance for the Frobenius norm of the ensemble’s Jacobian.

**Proposition 1** *Let  $x_i$  be independent random variables drawn from a Normal distribution  $x_i \sim \mathcal{N}(\mu, \sigma^2)$  for a fixed mean  $\mu$  and variance  $\sigma^2$ . Define  $\bar{x} = \sum_{i=1}^M c_i x_i$  where  $\sum_{i=1}^M c_i = 1$  and  $0 < c_i < 1, \forall i$ . We then have the following:*

$$\frac{\sigma^2}{M} \leq \sum_{i=1}^M c_i^2 \sigma^2 < \sigma^2 \quad (2)$$

*Proof.* By linearity of Normal distributions:  $\bar{x} \sim \mathcal{N}(\mu, \sum_{i=1}^M c_i^2 \sigma^2)$ . We then derive bounds for  $\sum_{i=1}^M c_i^2$  when  $M \geq 2$ :

$$\begin{aligned} \sum_{i=1}^M c_i^2 &= \sum_{i=1}^M c_i^2 + 2 \sum_{i=1}^M \sum_{j \neq i} c_i c_j - 2 \sum_{i=1}^M \sum_{j \neq i} c_i c_j \\ &= \left( \sum_{i=1}^M c_i \right)^2 - 2 \sum_{i=1}^M \sum_{j \neq i} c_i c_j \\ &= 1 - 2 \sum_{i=1}^M \sum_{j \neq i} c_i c_j \end{aligned}$$

Since  $c_i c_j > 0$  for all pairs  $i, j$ , then it follows that we have the upper bound  $\sum_{i=1}^M c_i^2 < 1$ . Note that equality:  $\sum_{i=1}^M c_i^2 = 1$  is only possible in the degenerate

case (when  $c_i = 1$  for exactly one  $i$  and  $c_j = 0$  for  $i \neq j$ ).  
 For the lower bound, we use Cauchy-Schwarz inequality to get:

$$\left( \sum_{i=1}^M (c_i \cdot 1) \right)^2 \leq \left( \sum_{i=1}^M c_i^2 \right) \left( \sum_{i=1}^M 1^2 \right) = \left( \sum_{i=1}^M c_i^2 \right) \cdot M$$

The left hand side reduces to 1, so it follows that  $\sum_{i=1}^M c_i^2 \geq \frac{1}{M}$  with equality when  $c_i = \frac{1}{M}$  for all  $i$ . Multiplying all sides with  $\sigma^2$  gives the desired bounds.  $\square$

Let  $\mathcal{F}$  be an ensemble of  $M$  base learners  $f_i: \mathcal{F}(x) = \sum_{i=1}^M c_i f_i(x)$  where  $\sum_{i=1}^M c_i = 1$  and  $0 < c_i < 1, \forall i$ . Let  $\mathbf{J}_{\mathcal{F}}$  denote the Jacobian of  $\mathcal{F}$  and  $\mathbf{J}_i$  the Jacobian of  $f_i$  for all  $i$ . It follows that  $\mathbf{J}_{\mathcal{F}} = \sum_{i=1}^M c_i \mathbf{J}_i$ .

**Theorem 1.** *Let each matrix  $\mathbf{J}_i \in \mathbb{R}^{C \times D}$  be comprised of the independent random variables  ${}_i a_{pq} \sim \mathcal{N}(\mu, \sigma^2)$ , where  ${}_i a_{pq}$  is the element on the  $p$ -th row and  $q$ -th column of matrix  $\mathbf{J}_i$ . It follows that their convex combination  $\mathbf{J}_{\mathcal{F}}$  satisfies:*

$$CD \left( \frac{\sigma^2}{M} + \mu^2 \right) \leq \mathbb{E}(\|\mathbf{J}_{\mathcal{F}}\|_F^2) < \mathbb{E}(\|\mathbf{J}_i\|_F^2) \quad (3)$$

$$CD \left( \frac{4\mu^2\sigma^2}{M} + \frac{2\sigma^4}{M^2} \right) \leq \text{Var}(\|\mathbf{J}_{\mathcal{F}}\|_F^2) < \text{Var}(\|\mathbf{J}_i\|_F^2) \quad (4)$$

*Proof.* Taking the square of Frobenius norm, we have the following for the Jacobian of a single model:

$$\|\mathbf{J}_i\|_F^2 = \sum_{p=1}^C \sum_{q=1}^D |{}_i a_{pq}|^2$$

The moments of  $\|\mathbf{J}_i\|_F^2$  are proportional to the moments of the random variables  ${}_i a_{pq}^2$ . These follow a chi-squared distribution with 1 degree of freedom, and have the expectation and variance:

$$\begin{aligned} \mathbb{E}({}_i a_{pq}^2) &= \text{Var}({}_i a_{pq}) + [\mathbb{E}({}_i a_{pq})]^2 \\ &= \sigma^2 + \mu^2 \\ \text{Var}({}_i a_{pq}^2) &= \mathbb{E}({}_i a_{pq}^4) - [\mathbb{E}({}_i a_{pq}^2)]^2 \\ &= \mu^4 + 6\mu^2\sigma^2 + 3\sigma^4 - \mu^4 - 2\mu^2\sigma^2 - \sigma^4 \\ &= 4\mu^2\sigma^2 + 2\sigma^4 \end{aligned}$$

Define  $\bar{a}_{pq} = \sum_{i=1}^M c_i {}_i a_{pq}$ , the elements of  $\mathbf{J}_{\mathcal{F}}$ . Note that  $\bar{a}_{pq} \sim \mathcal{N}(\mu, \sum_{i=1}^M c_i^2 \sigma^2)$ . Thus for the ensemble model's Jacobian, we have:

$$\|\mathbf{J}_{\mathcal{F}}\|_F^2 = \sum_{p=1}^C \sum_{q=1}^D |\bar{a}_{pq}|^2$$

It is clear that the moments of  $\|\mathbf{J}_{\mathcal{F}}\|_F$  are proportional to that of  $\bar{a}_{pq}^2$ . These random variables follow a chi-squared distribution with  $M$  degrees of freedom, and have the following expectation and variance:

$$\begin{aligned} \mathbb{E}(\bar{a}_{pq}^2) &= \sum_{i=1}^M c_i^2 \sigma^2 + \mu^2 \\ \text{Var}(\bar{a}_{pq}^2) &= 4\mu^2 \sum_{i=1}^M c_i^2 \sigma^2 + 2 \left( \sum_{i=1}^M c_i^2 \sigma^2 \right)^2 \end{aligned}$$

Applying **Proposition 1**, we then have the following bounds for the expectation and variance for these random variables:

$$\begin{aligned} \frac{\sigma^2}{M} + \mu^2 &\leq \mathbb{E}(\bar{a}_{pq}^2) < \mathbb{E}(i a_{pq}^2) \\ \frac{4\mu^2 \sigma^2}{M} + \frac{2\sigma^4}{M^2} &\leq \text{Var}(\bar{a}_{pq}^2) < \text{Var}(i a_{pq}^2) \end{aligned}$$

As the random variables are independently drawn, our desired result follows:

$$\begin{aligned} CD \left( \frac{\sigma^2}{M} + \mu^2 \right) &\leq \mathbb{E}(\|\mathbf{J}_{\mathcal{F}}\|_F^2) < \mathbb{E}(\|\mathbf{J}_i\|_F^2) \\ CD \left( \frac{4\mu^2 \sigma^2}{M} + \frac{2\sigma^4}{M^2} \right) &\leq \text{Var}(\|\mathbf{J}_{\mathcal{F}}\|_F^2) < \text{Var}(\|\mathbf{J}_i\|_F^2) \quad \square \end{aligned}$$

**Conclusion.** These are proportional to the expectation and variance of the Frobenius norms of our Jacobian matrices, so we can derive the following conclusions in this scenario. Ensembles decrease both the expected value and variance of the Jacobian's Frobenius norms when compared to that of a single model's. As  $M$  increases, the lower bounds of both the expectation and variance decreases.

Averaging is one of the most common methods for aggregating base learner outputs in an ensemble [16,26], so it is important to consider this case. When the ensemble is done via averaging, i.e.  $c_i = \frac{1}{M}$  for all  $i$ , this achieves the equality condition for the lower bounds of both  $\mathbb{E}(\bar{a}_{pq}^2)$  and  $\text{Var}(\bar{a}_{pq}^2)$ . Therefore, increasing the number of models in the ensemble *strictly decreases* the expectation and variance of the ensemble's Jacobian's norm  $\|\mathbf{J}_{\mathcal{F}}\|_F$ .

This theoretical result gives us the motivation on how model ensembles also improve the stability of models and thus their robustness to UAPs. We show this by deriving the above bounds on the Frobenius norm on their Jacobian. As model ensembles have also been shown to have improved performance over a single classifier, this makes it an ideal candidate for improving both model accuracy and robustness. In the next section, we explore the robustness of model ensembles and verify our theory with empirical results.

## 4 Experiments with Jacobian Ensembles

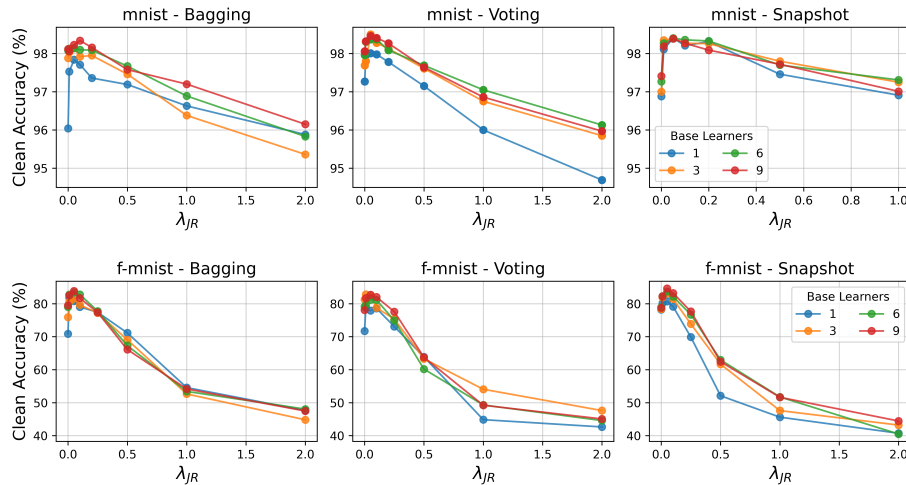
### 4.1 Experimental Setup

**Jacobian Ensembles.** To apply Jacobian Ensembles, we only need to include the Jacobian regularization as described in Eq. 1 to the joint loss of standard ensemble methods. The Jacobian regularization parameter  $\lambda_{JR}$  is tested for values between 0 and 2: where the resulting models manage to maintain good accuracy as informed by previous work [6].

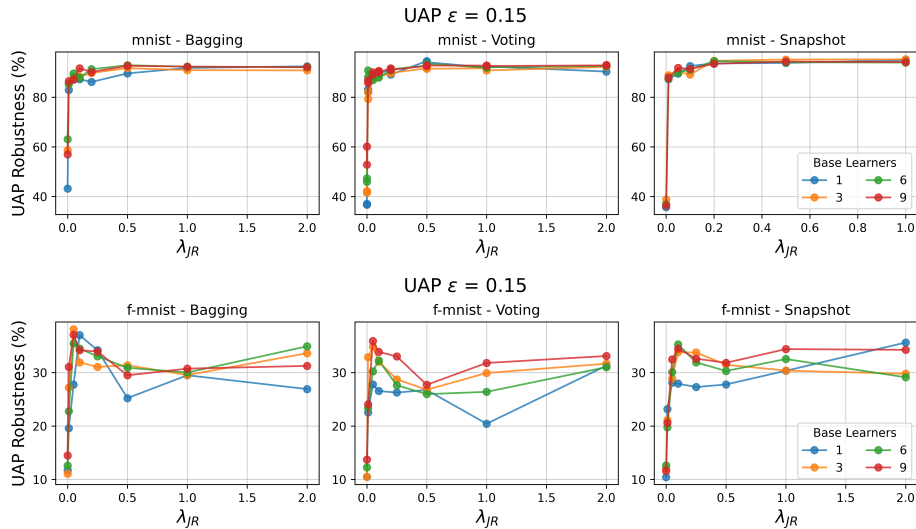
We evaluate the following ensemble methods: *Bagging* [1], *Snapshot Ensembles* [14], and *Soft Voting* [26]. For these, we evaluate all experiments with ensembles trained on 1, 3, 6, and 9 base learners. Effectively, one base learner is similar to using no ensemble method at all.

**Models & Datasets.** We use the MNIST [17] and Fashion-MNIST [25] datasets. These are popular image classification benchmarks, each with 10 classes, and 28 by 28 pixel images whose their pixel values range from 0 to 1. For the DNN architecture, we use a version of LeNet-5 [17,13], which we refer to as LeNet.

**UAP Attacks.** We evaluate the robustness of these models to UAPs generated via iterative stochastic gradient descent with 100 iterations and a batch size of 200. Perturbations are applied under  $\ell_\infty$ -norm constraints. The  $\varepsilon$  we consider in our attacks for this norm are from 0.10 to 0.25, this perturbation magnitude is equivalent to 10%-25% of the maximum total possible change in pixel values. UAPs are generated over 50 different random seeds, and we report UAPs with the highest attack success rate, as this would represent the worst-case scenario.



**Fig. 1.** Clean accuracy of LeNet on MNIST (top) and Fashion-MNIST (bottom) for various Jacobian regularization strengths  $\lambda_{JR}$  and with varying number of base learners per ensemble method.



**Fig. 2.** Model robustness against UAP with  $\varepsilon = 0.15$  of LeNet on MNIST (top) and Fashion-MNIST (bottom) for various Jacobian regularization strengths  $\lambda_{JR}$  and with varying number of base learners per ensemble method.

**Metrics.** The following metrics are evaluated on the entire 10,000 sample test sets for each dataset. *Clean Accuracy* is the accuracy of the model on the test set. *Model Robustness* measures the accuracy of the model on the test set when the corresponding worst-case UAP is applied or present. We then average this model robustness over all the UAP attack scenarios that we consider,  $\ell_\infty$ -norm of  $\varepsilon = 0.10, 0.15, 0.2, 0.25$ , to get an overall *mean UAP Accuracy*.

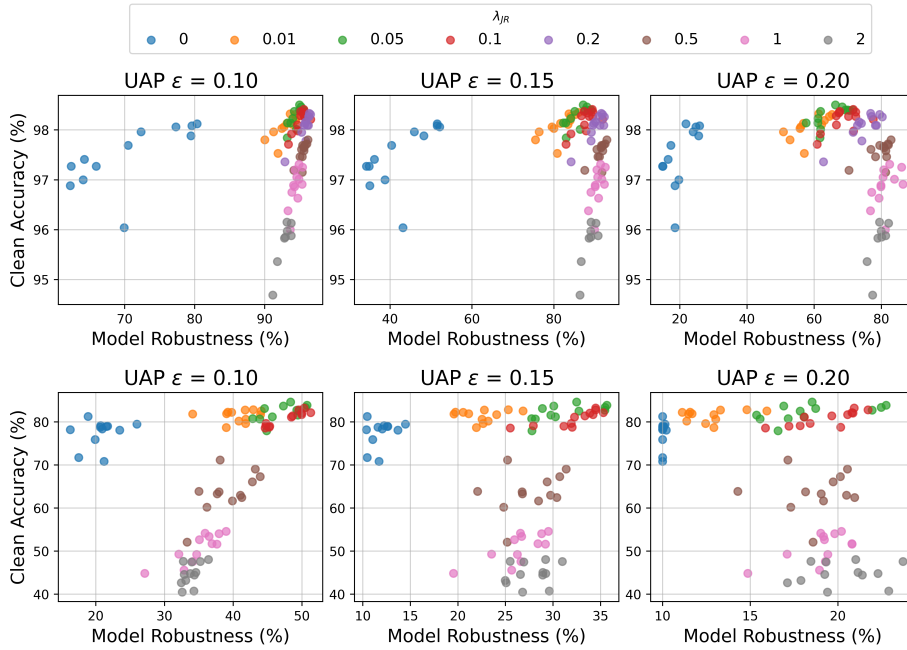
## 4.2 Improvements with Jacobian Ensembles

**Clean Accuracy.** In Fig. 1 there is a rapid degradation in clean accuracy when  $\lambda_{JR}$  is large. This is when JR is more heavily weighted, and this is consistent with previous work [6,13] as too much regularization damages accuracy on the test set. However, having a small amount of JR is still noticeably more beneficial than no JR as indicated by the clean accuracy when considering  $\lambda_{JR}$  in the range of 0 to 0.1 for both datasets across all settings.

We also see benefits of using ensembles: models with more than one base learner have noticeably better clean accuracy across all settings in Fig. 1. Overall, Jacobian Ensembles, which is a combination of ensemble methods and JR, achieves the best performance in our experiments.

**Model Robustness.** In the interest of space, we only present the robustness results for a particular UAP attack strength  $\varepsilon = 0.15$  in Fig. 2. In terms of robustness, trends for both datasets are slightly different since Fashion-MNIST is a more difficult dataset than MNIST: a regularly trained LeNet on MNIST





**Fig. 3.** Clean accuracy versus robustness trade-off of LeNet on MNIST (top) and Fashion-MNIST (bottom) labeled by  $\lambda_{JR}$  for various UAP attack strengths  $\epsilon$ .

can be expected to have 98-99% clean accuracy whereas it is 90-91% for the same model on Fashion-MNIST. Thus, it is expected that models on Fashion-MNIST are considerably less robust. For most settings, the general trends in Fig. 2 show that ensembles have a robustness benefit. JR monotonically increases robustness for MNIST and for Fashion-MNIST has a range between  $0 < \lambda_{JR} < 0.5$  that achieves the best robustness.

In conclusion, ensemble methods with more than one base learner introduce measurable advantages in both accuracy and robustness. These advantages become even more pronounced when combined with JR. Next we analyze the general trade-off between accuracy and robustness.

### 4.3 Accuracy-Robustness Trade-Off

In Fig. 3, we plot the accuracy versus robustness of the various models we trained under the different configurations accounting for various ensemble methods, number of base learners, and  $\lambda_{JR}$  values. The best models achieve both high accuracy and robustness, so these will be on the top right side of the graph.

In a scenario when robustness is not accounted for (i.e.  $\lambda_{JR} = 0$ ), models would appear in the top left. They are trained to have good accuracy, but remain extremely vulnerable to adversarial attacks like UAPs. On the other hand,

**Table 1.** Weighted accuracy (in %) of LeNet for MNIST (top) and Fashion-MNIST (bottom). The first 3 rows show the models with highest weighted accuracy. The bottom 3 rows show the *best* weighted accuracy of models trained with only JR, only ensemble, and “standard training” (neither JR nor ensemble).

Model			MNIST Accuracy (%)		
Ensemble	Learners	$\lambda_{JR}$	Clean	Avg. UAP	<b>Weighted</b>
Snapshot	3	0.50	97.8	82.8	<b>90.3</b>
Snapshot	9	0.50	97.7	82.3	90.0
Snapshot	3	0.20	98.3	81.5	89.9
JR Only	1	0.05	98.4	74.2	86.3
Bagging	6	0	98.1	42.8	70.4
Standard	1	0	99.1	31.9	65.5

Model			Fashion-MNIST Accuracy (%)		
Ensemble	Learners	$\lambda_{JR}$	Clean	Avg. UAP	<b>Weighted</b>
Bagging	9	0.05	83.9	43.2	<b>63.5</b>
Bagging	6	0.05	83.4	43.0	63.2
Snapshot	9	0.10	83.2	42.4	62.8
JR Only	1	0.10	79.0	31.2	55.1
Bagging	9	0	91.2	12.7	51.2
Standard	1	0	90.8	12.1	50.9

models that overcompensate for robustness such as those with very high regularization (e.g.  $\lambda_{JR} = 2$ ), will appear on the bottom right. As these models sacrifice a significant amount of clean accuracy for improved robustness, especially against UAP attacks with larger strength  $\varepsilon$ . These delineations become clear when labeling the models according to their  $\lambda_{JR}$  value as in Fig. 3. Thus, the role of  $\lambda_{JR}$  is evident in improving overall robustness.

To better capture the model performance on both benign and adversarial inputs, we compute the *Weighted Accuracy* by averaging the clean accuracy and mean UAP accuracy. In practice, the defender can adjust the weighting of each accuracy metric in their final assessment to better match their application and risk profile. We choose the mean as the base setting.

Next, we perform an ablation study on the effect of only JR and only ensemble models compared against the top 3 Jacobian Ensembles with the best weighted accuracy in Table 1. We find that Jacobian Ensembles achieve the best weighted accuracy. To compare, we also show in the bottom 3 rows for each table the best models with only JR, only ensembles, and neither JR nor ensembles. The best Jacobian Ensembles have a clear advantage with average UAP accuracy over the the non-Jacobian Ensembles whilst maintaining very close clean accuracy. This difference is even more pronounced on the Fashion-MNIST dataset.

Differences between the two datasets MNIST and Fashion-MNIST also show in Table 1. Since MNIST is an easier dataset, performance degradation by large  $\lambda_{JR} = 0$  are not as prominent, so larger  $\lambda_{JR}$  are favored by the combined score. For Fashion-MNIST, a large  $\lambda_{JR}$  is detrimental as the base model begin with

a relatively low clean accuracy ( $< 91\%$ ). In both cases, ensembles demonstrate a noticeably large boost in weighted accuracy, and further tuning is likely to improve their performance.

## 5 Conclusion

In this work, we propose Jacobian Ensembles to significantly increase model robustness against UAPs whilst maintaining the clean accuracy of models. Our results show that Jacobian Ensembles takes the advantages of both Jacobian regularization and model ensembles to achieve superior accuracy and robustness than each of these methods on their own, as measured by our weighted metric.

In addition, we derive theoretical upper and lower bounds on the robustness to UAPs for model ensembles, showing that increasing the number of base classifiers in the models' ensembles reduces the expected Frobenius norm of their Jacobian and thus improves stability. We then empirically verify our results and show that a combination of both JR and ensembles achieve the best performance.

These results give us confidence in recommending Jacobian Ensembles as a general methodology when training models as UAPs present a great threat to model adoption and safety. Our results show that it is indeed possible to maintain great test accuracy whilst achieving significant UAP robustness in previously unseen levels of accuracy-robustness trade-off. Thus, it is indeed possible to get the best of both worlds.

## References

1. Breiman, L.: Bagging predictors. *Machine learning* **24**(2), 123–140 (1996)
2. Brown, T.B., Mané, D.: Adversarial patch. arXiv preprint arXiv:1712.09665 (2017)
3. Co, K.T., Muñoz González, L., de Maupeou, S., Lupu, E.C.: Procedural noise adversarial examples for black-box attacks on deep convolutional networks. In: *Proceedings of the 2019 ACM SIGSAC Conf. on Computer and Communications Security*. pp. 275–289. CCS '19 (2019). <https://doi.org/10.1145/3319535.3345660>
4. Co, K.T., Muñoz-González, L., Kanthan, L., Glocker, B., Lupu, E.C.: Universal adversarial robustness of texture and shape-biased models. arXiv preprint arXiv:1911.10364 (2019)
5. Co, K.T., Muñoz-González, L., Lupu, E.C.: Sensitivity of deep convolutional networks to gabor noise. arXiv preprint arXiv:1906.03455 (2019)
6. Co, K.T., Rego, D.M., Lupu, E.C.: Jacobian regularization for mitigating universal adversarial perturbations. In: *International Conference on Artificial Neural Networks*. pp. 202–213. Springer (2021)
7. Eykholt, K., Evtimov, I., Fernandes, E., Li, B., Rahmati, A., Tramèr, F., Prakash, A., Kohno, T., Song, D.: Physical adversarial examples for object detectors. In: *12th USENIX Workshop on Offensive Technologies (WOOT 18)* (2018)
8. Eykholt, K., Evtimov, I., Fernandes, E., Li, B., Rahmati, A., Xiao, C., Prakash, A., Kohno, T., Song, D.: Robust physical-world attacks on deep learning visual classification. In: *Proceedings of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. pp. 1625–1634 (2018)

9. Freund, Y., Schapire, R., Abe, N.: A short introduction to boosting. *Journal-Japanese Society For Artificial Intelligence* **14**(771-780), 1612 (1999)
10. Hau, Z., Co, K.T., Demetriou, S., Lupu, E.C.: Object removal attacks on lidar-based 3d object detectors. *arXiv preprint arXiv:2102.03722* (2021)
11. Hau, Z., Demetriou, S., Muñoz-González, L., Lupu, E.C.: Shadow-catcher: Looking into shadows to detect ghost objects in autonomous vehicle 3d sensing. In: *Euro. Symposium on Research in Computer Security*. pp. 691–711. Springer (2021)
12. Hinton, G., Deng, L., Yu, D., Dahl, G.E., Mohamed, A.r., Jaitly, N., Senior, A., Vanhoucke, V., Nguyen, P., Sainath, T.N., et al.: Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups. *IEEE Signal Processing Magazine* **29**(6), 82–97 (2012)
13. Hoffman, J., Roberts, D.A., Yaida, S.: Robust learning with jacobian regularization. *arXiv preprint arXiv:1908.02729* (2019)
14. Huang, G., Li, Y., Pleiss, G., Liu, Z., Hopcroft, J.E., Weinberger, K.Q.: Snapshot ensembles: Train 1, get m for free. In: *Intl. Conf. on Learning Rep.* (2017)
15. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. In: *Advances in Neural Information Processing Systems (NeurIPS)*. pp. 1097–1105 (2012)
16. Kuncheva, L.I.: *Combining pattern classifiers: methods and algorithms*. John Wiley & Sons (2014)
17. LeCun, Y., Bottou, L., Bengio, Y., Haffner, P.: Gradient-based learning applied to document recognition. *Proceedings of the IEEE* **86**(11), 2278–2324 (1998)
18. Matachana, A.G., Co, K.T., Muñoz-González, L., Martínez, D., Lupu, E.C.: Robustness and transferability of universal attacks on compressed models. *arXiv preprint arXiv:2012.06024* (2020)
19. Moosavi-Dezfooli, S.M., Fawzi, A., Fawzi, O., Frossard, P.: Universal adversarial perturbations. In: *Proceedings of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. pp. 1765–1773 (2017)
20. Redmon, J., Divvala, S., Girshick, R., Farhadi, A.: You only look once: Unified, real-time object detection. In: *Proceedings of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. pp. 779–788 (2016)
21. Roth, K., Kilcher, Y., Hofmann, T.: Adversarial training is a form of data-dependent operator norm regularization. In: *Advances in Neural Information Processing Systems (NeurIPS)* (2020)
22. Shafahi, A., Najibi, M., Xu, Z., Dickerson, J., Davis, L.S., Goldstein, T.: Universal adversarial training. *arXiv preprint arXiv:1811.11304* (2018)
23. Thys, S., Van Ranst, W., Goedemé, T.: Fooling automated surveillance cameras: adversarial patches to attack person detection. In: *CVPRW: Workshop on The Bright and Dark Sides of Computer Vision: Challenges and Opportunities for Privacy and Security* (2019)
24. Tramèr, F., Dupré, P., Rusak, G., Pellegrino, G., Boneh, D.: Adversarial: Perceptual ad blocking meets adversarial machine learning. In: *Proceedings of the 2019 ACM SIGSAC Conf. on Computer and Communications Security*. p. 2005–2021. *CCS '19* (2019). <https://doi.org/10.1145/3319535.3354222>
25. Xiao, H., Rasul, K., Vollgraf, R.: Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747* (2017)
26. Zhou, Z.H.: *Ensemble methods: foundations and algorithms*. CRC press (2012)