

# Adversarial Vulnerability of Randomized Ensembles

**Hassan Dbouk & Naresh Shanbhag**

*Department of Electrical and Computer Engineering*

*University of Illinois at Urbana-Champaign*

**I ILLINOIS**

Electrical & Computer Engineering

COLLEGE OF ENGINEERING

# Robust and Efficient Inference

deep nets are vulnerable

original sample



decision: 'panda'

+ .007 ×



=

adversarial sample



decision: 'gibbon'

# Robust and Efficient Inference

deep nets are vulnerable

original sample



decision: 'panda'

+ .007 ×



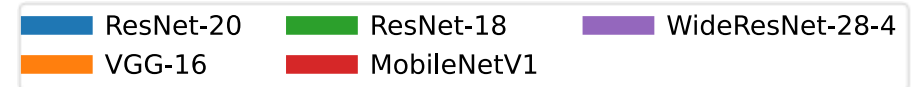
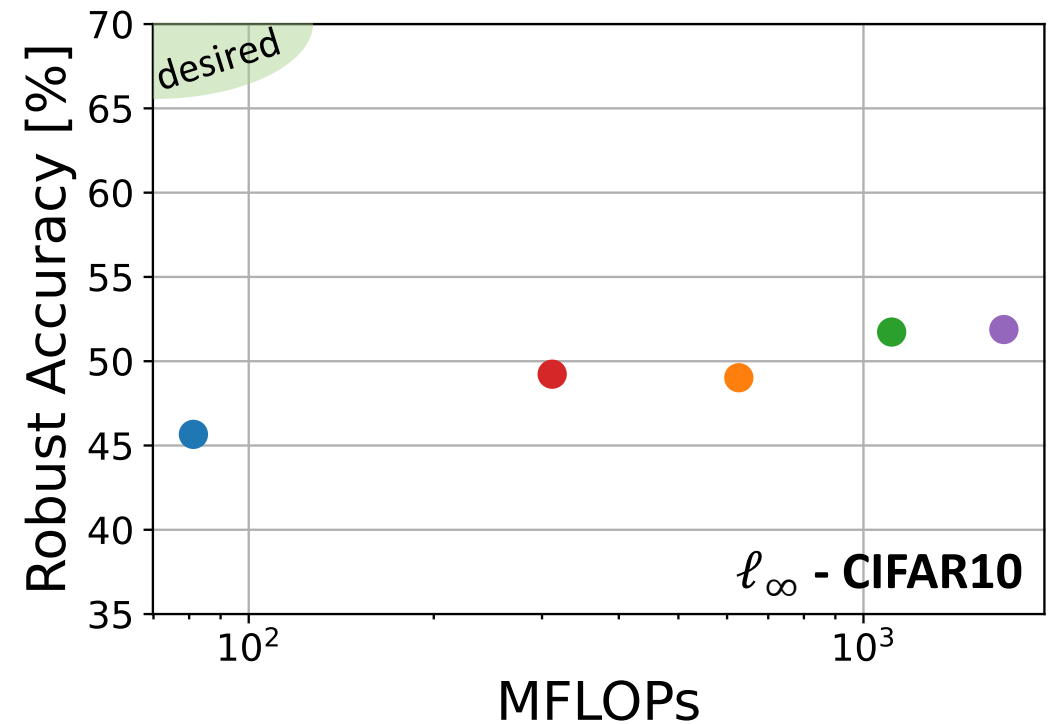
=

adversarial sample



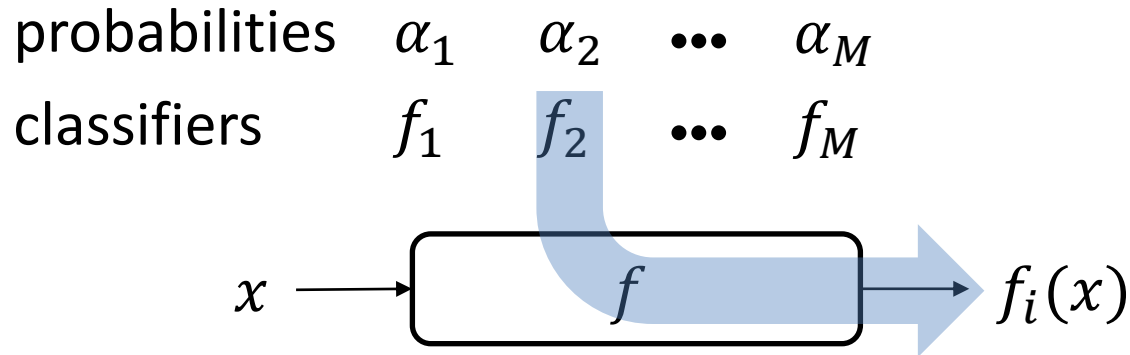
decision: 'gibbon'

robustness is expensive



# Robustness via Randomized Ensembles

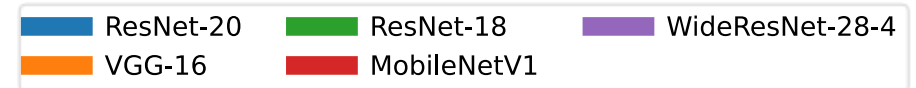
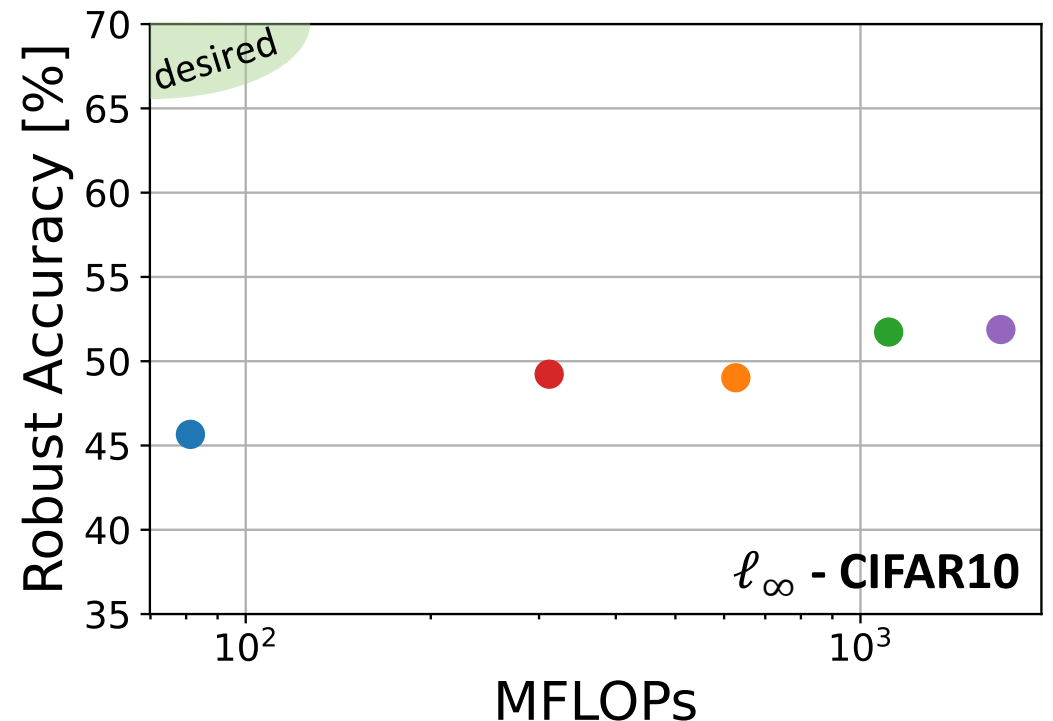
**multiple** classifiers  $f_1, \dots, f_M$



inference: pick **one** at random

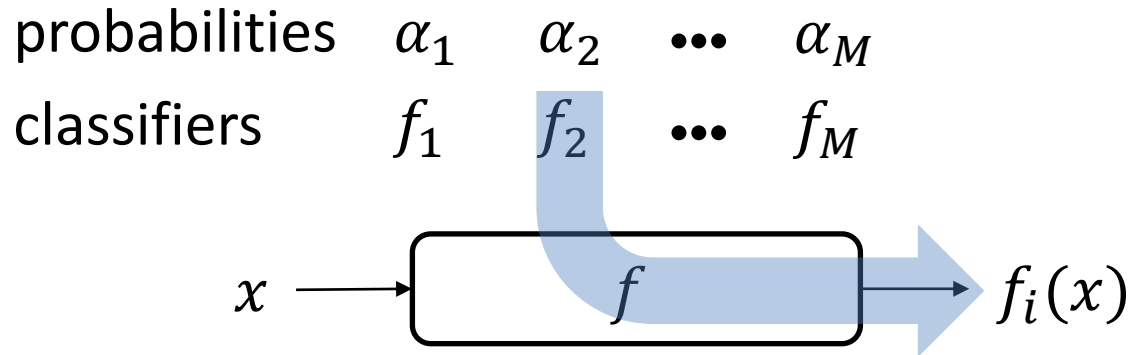
**no** increase in # of FLOPs

robustness is expensive



# Robustness via Randomized Ensembles

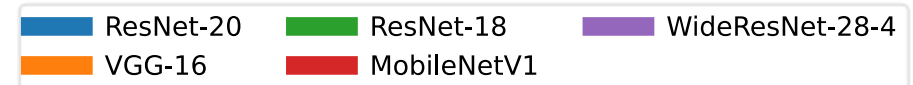
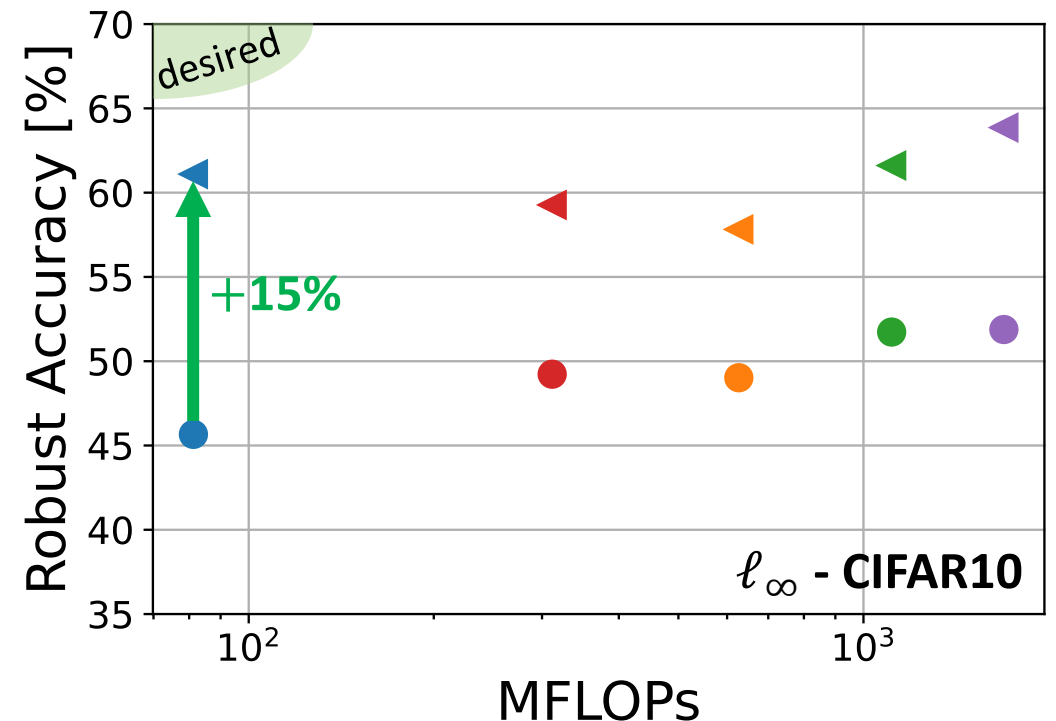
**multiple** classifiers  $f_1, \dots, f_M$



inference: pick **one** at random

**no** increase in # of FLOPs

using two classifiers trained via BAT [Pinot et al, 2020]



# Robustness via Randomized Ensembles

**multiple** classifiers  $f_1, \dots, f_M$

probabilities  $\alpha_1 \quad \alpha_2 \quad \dots \quad \alpha_M$

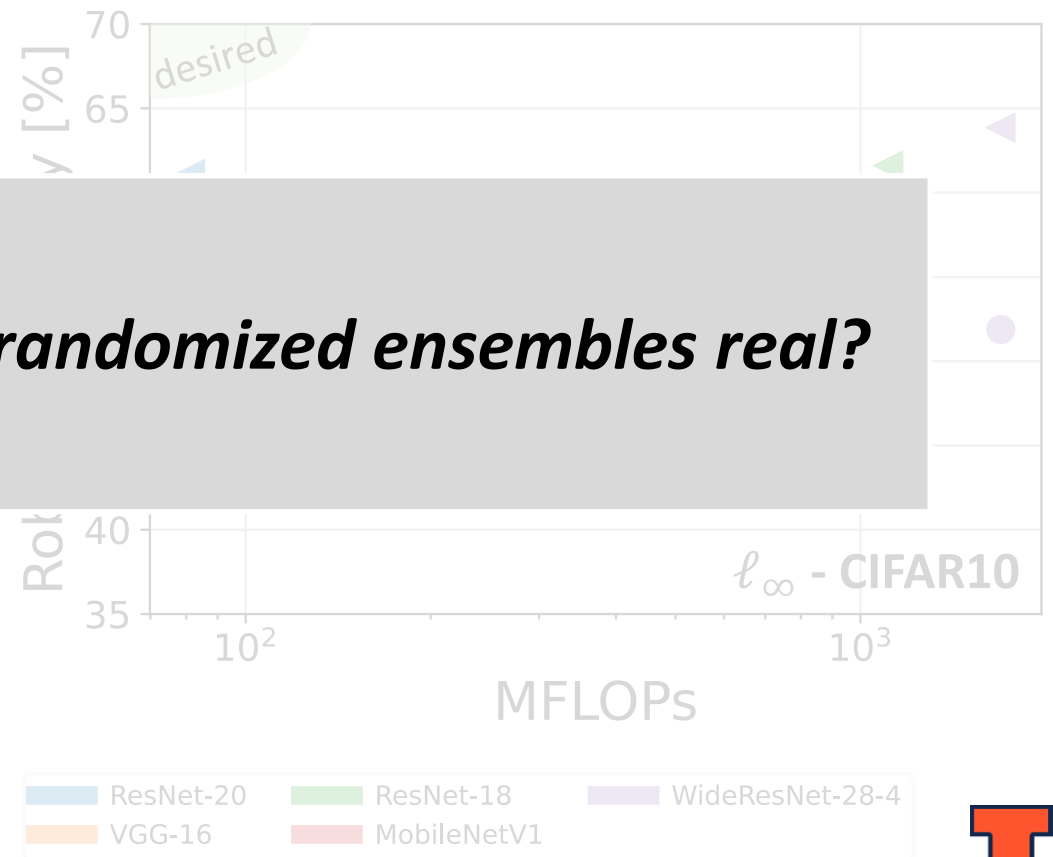
clas

***Are the robustness gains provided by randomized ensembles real?***

inference: pick **one** at random

**no** increase in # of FLOPS

using two classifiers trained via BAT [Pinot et al, 2020]

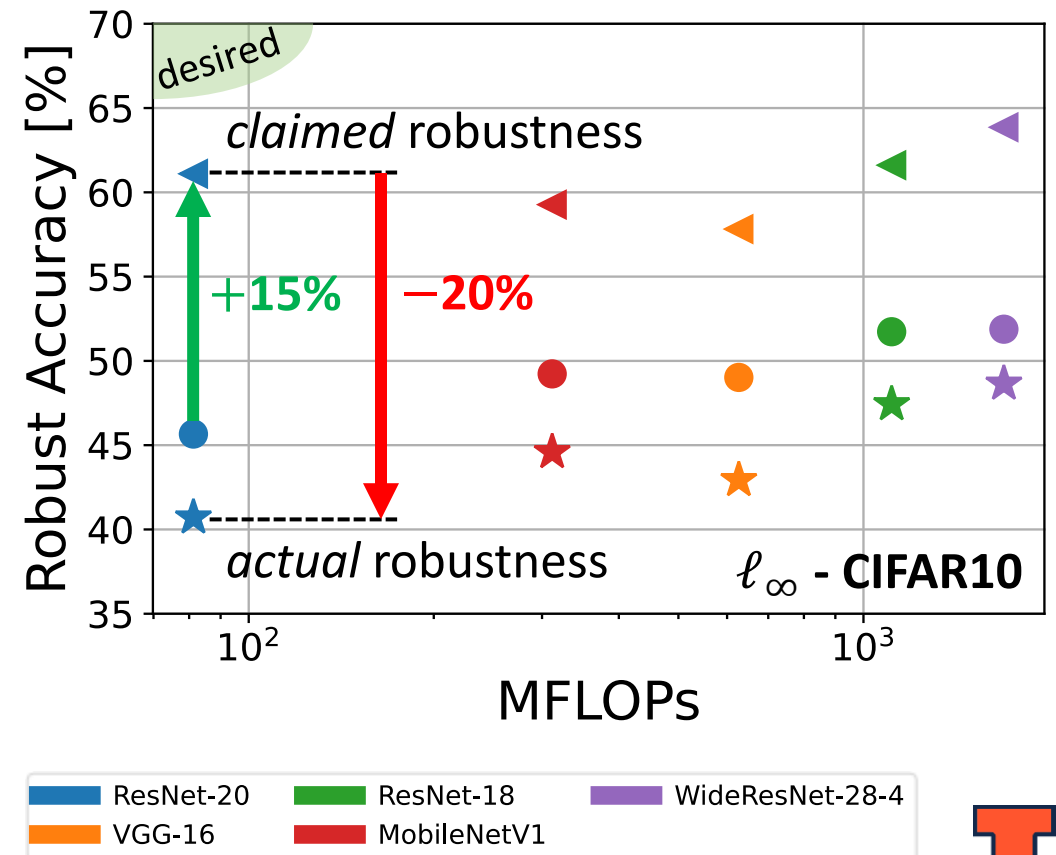


# This Work: Revealing the Vulnerability

## main contributions

- show that adaptive PGD is ill-suited for evaluating robustness of RECs
  - no guarantees even for linear classifiers
- propose a provably consistent and efficient adversarial attack algorithm – **ARC: Attacking Randomized ensembles of Classifiers**
- demonstrate that existing randomized ensembles defenses are in fact more vulnerable than standard AT

## BAT defense compromised



# The ARC Algorithm – Binary Linear Classifiers

---

**Algorithm 1** The ARC Algorithm for BLCs

---

1: **Input:** REC  $(\mathcal{F}, \alpha)$ , labeled data-point  $(\mathbf{x}, y)$ , norm  $p$ , and radius  $\epsilon$ .

2: **Output:** Adversarial perturbation  $\delta$  such that  $\|\delta\|_p \leq \epsilon$ .

3: Initialize  $\delta \leftarrow \mathbf{0}$ ,  $v \leftarrow L(\mathbf{x}, y, \alpha)$ ,  $q \leftarrow \frac{p}{p-1}$

4: Define  $\mathcal{I}$  such that  $\alpha_i \geq \alpha_j \forall i, j \in \mathcal{I}$  and  $i \leq j$ .

5: **for**  $i \in \mathcal{I}$  **do**

6:   */\* optimal unit  $\ell_p$  norm adversarial direction for  $f_i$  \*/*

7:    $\mathbf{g} \leftarrow -y \frac{|\mathbf{w}_i|^{q-1} \odot \text{sgn}(\mathbf{w}_i)}{\|\mathbf{w}_i\|_q^{q-1}}$

8:   */\* shortest  $\ell_p$  distance between  $\mathbf{x}$  and  $f_i$  \*/*

9:    $\zeta \leftarrow \frac{|f_i(\mathbf{x})|}{\|\mathbf{w}_i\|_q}$

10:   **if**  $\zeta \geq \epsilon \vee i = 1$  **then**

11:      $\beta \leftarrow \epsilon$

12:   **else**

13:      $\beta \leftarrow \frac{\epsilon}{\epsilon - \zeta} \left| \frac{y \mathbf{w}_i^T \delta}{\|\mathbf{w}_i\|_q} + \zeta \right| + \rho$

14:   **end if**

15:    $\hat{\delta} \leftarrow \epsilon \frac{\delta + \beta \mathbf{g}}{\|\delta + \beta \mathbf{g}\|_p}$    *▷ candidate  $\hat{\delta}$  such that  $\|\hat{\delta}\|_p = \epsilon$*

16:    $\hat{v} \leftarrow L(\mathbf{x} + \hat{\delta}, y, \alpha)$

17:   */\* if robustness does not increase, update  $\delta$  \*/*

18:   **if**  $\hat{v} \leq v$  **then**

19:      $\delta \leftarrow \hat{\delta}$ ,  $v \leftarrow \hat{v}$

20:   **end if**

21: **end for**

---

- greedily iterate over all classifiers once





# The ARC Algorithm – Binary Linear Classifiers

---

## Algorithm 1 The ARC Algorithm for BLCs

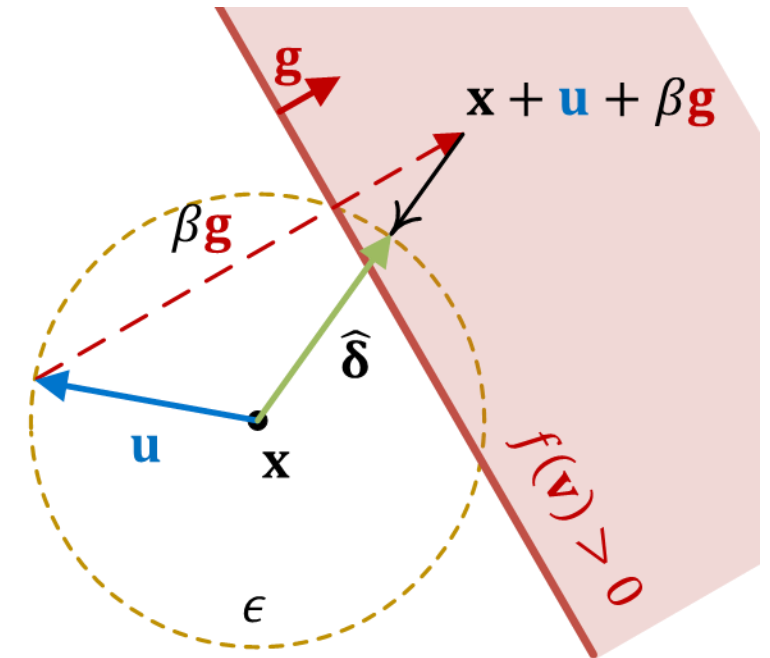
---

- 1: **Input:** REC  $(\mathcal{F}, \alpha)$ , labeled data-point  $(\mathbf{x}, y)$ , norm  $p$ , and radius  $\epsilon$ .
- 2: **Output:** Adversarial perturbation  $\delta$  such that  $\|\delta\|_p \leq \epsilon$ .
- 3: Initialize  $\delta \leftarrow \mathbf{0}$ ,  $v \leftarrow L(\mathbf{x}, y, \alpha)$ ,  $q \leftarrow \frac{p}{p-1}$
- 4: Define  $\mathcal{I}$  such that  $\alpha_i \geq \alpha_j \forall i, j \in \mathcal{I}$  and  $i \leq j$ .
- 5: **for**  $i \in \mathcal{I}$  **do**
- 6:   */\* optimal unit  $\ell_p$  norm adversarial direction for  $f_i$  \*/*
- 7:    $\mathbf{g} \leftarrow -y \frac{|\mathbf{w}_i|^{q-1} \odot \text{sgn}(\mathbf{w}_i)}{\|\mathbf{w}_i\|_q^{q-1}}$
- 8:   */\* shortest  $\ell_p$  distance between  $\mathbf{x}$  and  $f_i$  \*/*
- 9:    $\zeta \leftarrow \frac{|f_i(\mathbf{x})|}{\|\mathbf{w}_i\|_q}$
- 10:   **if**  $\zeta \geq \epsilon \vee i = 1$  **then**
- 11:      $\beta \leftarrow \epsilon$
- 12:   **else**
- 13:      $\beta \leftarrow \frac{\epsilon}{\epsilon - \zeta} \left| \frac{y \mathbf{w}_i^T \delta}{\|\mathbf{w}_i\|_q} + \zeta \right| + \rho$
- 14:   **end if**
- 15:    $\hat{\delta} \leftarrow \epsilon \frac{\delta + \beta \mathbf{g}}{\|\delta + \beta \mathbf{g}\|_p}$    *▷ candidate  $\hat{\delta}$  such that  $\|\hat{\delta}\|_p = \epsilon$*
- 16:    $\hat{v} \leftarrow L(\mathbf{x} + \hat{\delta}, y, \alpha)$
- 17:   */\* if robustness does not increase, update  $\delta$  \*/*
- 18:   **if**  $\hat{v} \leq v$  **then**
- 19:      $\delta \leftarrow \hat{\delta}$ ,  $v \leftarrow \hat{v}$
- 20:   **end if**
- 21: **end for**

---

- greedily iterate over all classifiers once
- novel adaptive step size computation:

smallest  $\beta > 0$  such that  
 $\hat{\delta} = \gamma(\mathbf{u} + \beta \mathbf{g})$   
 can fool  $f$



# The ARC Algorithm – Binary Linear Classifiers

---

## Algorithm 1 The ARC Algorithm for BLCs

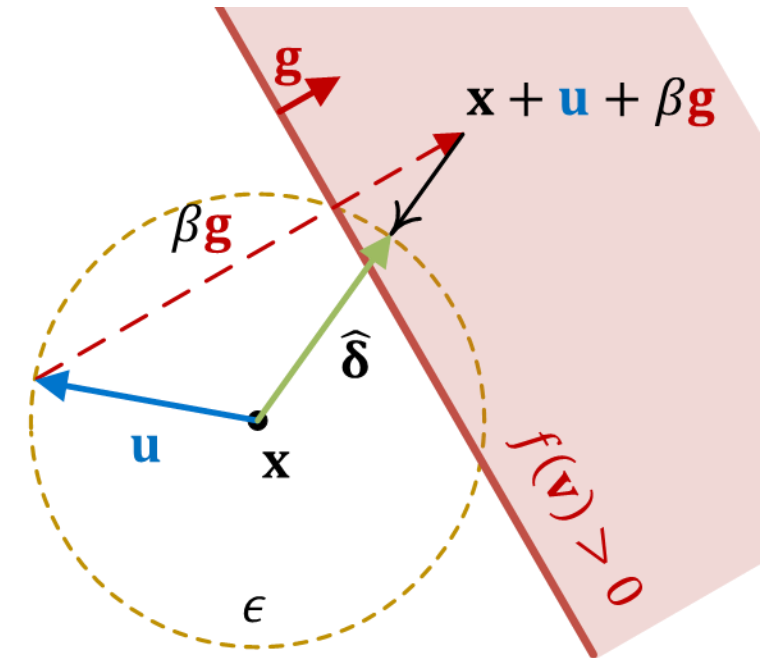
---

- 1: **Input:** REC  $(\mathcal{F}, \alpha)$ , labeled data-point  $(\mathbf{x}, y)$ , norm  $p$ , and radius  $\epsilon$ .
- 2: **Output:** Adversarial perturbation  $\delta$  such that  $\|\delta\|_p \leq \epsilon$ .
- 3: Initialize  $\delta \leftarrow \mathbf{0}$ ,  $v \leftarrow L(\mathbf{x}, y, \alpha)$ ,  $q \leftarrow \frac{p}{p-1}$
- 4: Define  $\mathcal{I}$  such that  $\alpha_i \geq \alpha_j \forall i, j \in \mathcal{I}$  and  $i \leq j$ .
- 5: **for**  $i \in \mathcal{I}$  **do**
- 6:   */\* optimal unit  $\ell_p$  norm adversarial direction for  $f_i$  \*/*
- 7:    $\mathbf{g} \leftarrow -y \frac{|\mathbf{w}_i|^{q-1} \odot \text{sgn}(\mathbf{w}_i)}{\|\mathbf{w}_i\|_q^{q-1}}$
- 8:   */\* shortest  $\ell_p$  distance between  $\mathbf{x}$  and  $f_i$  \*/*
- 9:    $\zeta \leftarrow \frac{|f_i(\mathbf{x})|}{\|\mathbf{w}_i\|_q}$
- 10:   **if**  $\zeta \geq \epsilon \vee i = 1$  **then**
- 11:      $\beta \leftarrow \epsilon$
- 12:   **else**
- 13:      $\beta \leftarrow \frac{\epsilon}{\epsilon - \zeta} \left| \frac{y \mathbf{w}_i^T \delta}{\|\mathbf{w}_i\|_q} + \zeta \right| + \rho$
- 14:   **end if**
- 15:    $\hat{\delta} \leftarrow \epsilon \frac{\delta + \beta \mathbf{g}}{\|\delta + \beta \mathbf{g}\|_p}$    *▷ candidate  $\hat{\delta}$  such that  $\|\hat{\delta}\|_p = \epsilon$*
- 16:    $\hat{v} \leftarrow L(\mathbf{x} + \hat{\delta}, y, \alpha)$
- 17:   */\* if robustness does not increase, update  $\delta$  \*/*
- 18:   **if**  $\hat{v} \leq v$  **then**
- 19:      $\delta \leftarrow \hat{\delta}$ ,  $v \leftarrow \hat{v}$
- 20:   **end if**
- 21: **end for**

---

- greedily iterate over all classifiers once
- novel adaptive step size computation:

smallest  $\beta > 0$  such that  
 $\hat{\delta} = \gamma(\mathbf{u} + \beta \mathbf{g})$   
 can fool  $f$



- extend to multiclass differentiable classifiers



# The ARC Algorithm – Binary Linear Classifiers

---

## Algorithm 1 The ARC Algorithm for BLCs

---

```

1: Input: REC  $(\mathcal{F}, \alpha)$ , labeled data-point  $(\mathbf{x}, y)$ , norm  $p$ ,
   and radius  $\epsilon$ .
2: Output: Adversarial perturbation  $\delta$  such that  $\|\delta\|_p \leq \epsilon$ .
3: Initialize  $\delta \leftarrow \mathbf{0}$ ,  $v \leftarrow L(\mathbf{x}, y, \alpha)$ ,  $q \leftarrow \frac{p}{p-1}$ 
4: Define  $\mathcal{I}$  such that  $\alpha_i \geq \alpha_j \forall i, j \in \mathcal{I}$  and  $i \leq j$ .
5: for  $i \in \mathcal{I}$  do
6:
7:
8:
9:    $\zeta \leftarrow \frac{\|\mathbf{w}_i\|_q}{\|\mathbf{w}_i\|_q}$ 
10:  if  $\zeta \geq \epsilon \vee i = 1$  then
11:     $\beta \leftarrow \epsilon$ 
12:  else
13:     $\beta \leftarrow \frac{\epsilon}{\epsilon - \zeta} \left| \frac{y \mathbf{w}_i^T \delta}{\|\mathbf{w}_i\|_q} + \zeta \right| + \rho$ 
14:  end if
15:   $\hat{\delta} \leftarrow \epsilon \frac{\delta + \beta \mathbf{g}}{\|\delta + \beta \mathbf{g}\|_p}$   $\triangleright$  candidate  $\hat{\delta}$  such that  $\|\hat{\delta}\|_p = \epsilon$ 
16:   $\hat{v} \leftarrow L(\mathbf{x} + \hat{\delta}, y, \alpha)$ 
17:  /* if robustness does not increase, update  $\delta$ 
18:  if  $\hat{v} \leq v$  then
19:     $\delta \leftarrow \hat{\delta}$ ,  $v \leftarrow \hat{v}$ 
20:  end if
21: end for

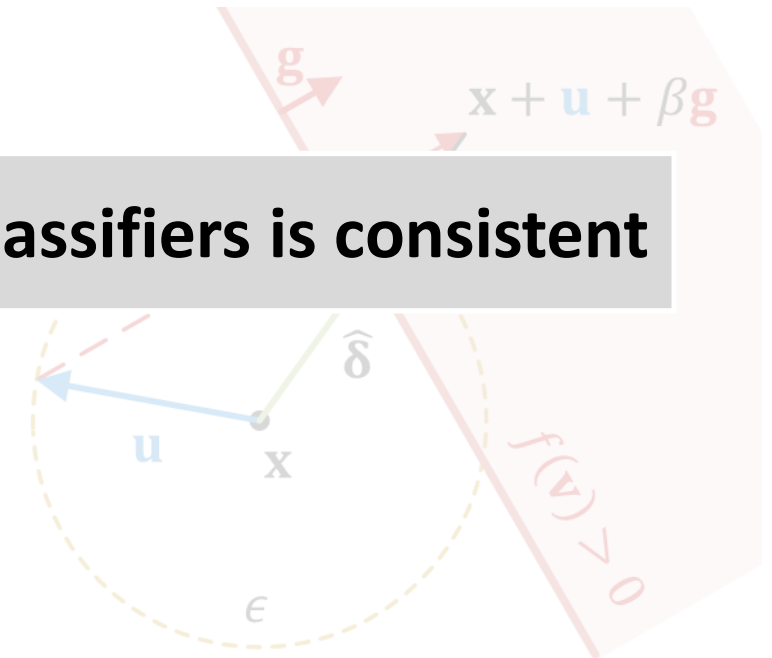
```

---

- greedily iterate over all classifiers once
- novel adaptive step size computation:

**Theorem: the ARC algorithm for binary linear classifiers is consistent**

$\hat{\delta} = \gamma(\mathbf{u} + \beta \mathbf{g})$   
can fool  $f$



- extend to multiclass differentiable classifiers



# Results Summary

Table 1. Comparison between ARC and adaptive PGD when attacking randomized ensembles trained via BAT (Pinot et al., 2020) across various network architectures and norms on the CIFAR-10 dataset. We use the standard radii  $\epsilon_2 = 128/255$  and  $\epsilon_\infty = 8/255$  for  $\ell_2$  and  $\ell_\infty$ -bounded perturbations, respectively.

NETWORK	NORM	ROBUST ACCURACY [%]			
		AT ( $M = 1$ )	REC ( $M = 2$ )		
		PGD	APGD	ARC	DIFF
RESNET-20	$\ell_2$	62.43	69.21	55.44	<b>-13.77</b>
	$\ell_\infty$	45.66	61.10	40.71	<b>-20.39</b>
MOBILENETV1	$\ell_2$	66.39	67.92	59.43	<b>-8.49</b>
	$\ell_\infty$	49.23	59.27	44.59	<b>-14.68</b>
VGG-16	$\ell_2$	66.08	66.96	59.20	<b>-7.76</b>
	$\ell_\infty$	49.02	57.82	42.93	<b>-14.89</b>
RESNET-18	$\ell_2$	69.16	70.16	65.88	<b>-4.28</b>
	$\ell_\infty$	51.73	61.61	47.43	<b>-14.18</b>
WIDERESNET-28-4	$\ell_2$	69.91	71.48	62.95	<b>-8.53</b>
	$\ell_\infty$	51.88	63.86	48.65	<b>-15.21</b>

Table 2. Comparison between ARC and adaptive PGD when attacking randomized ensembles trained via BAT (Pinot et al., 2020) across various datasets and norms. We use ResNet-18 for ImageNet and ResNet-20 for SVHN, CIFAR-10, and CIFAR-100 datasets.

DATASET	NORM	RADIUS ( $\epsilon$ )	ROBUST ACCURACY [%]			
			AT ( $M = 1$ )	REC ( $M = 2$ )		
			PGD	APGD	ARC	DIFF
SVHN	$\ell_2$	128/255	68.35	74.66	60.15	<b>-14.51</b>
	$\ell_\infty$	8/255	53.55	65.99	52.01	<b>-13.98</b>
CIFAR-10	$\ell_2$	128/255	62.43	69.21	55.44	<b>-13.77</b>
	$\ell_\infty$	8/255	45.66	61.10	40.71	<b>-20.39</b>
CIFAR-100	$\ell_2$	128/255	34.60	41.91	28.92	<b>-12.99</b>
	$\ell_\infty$	8/255	22.29	33.37	17.45	<b>-15.92</b>
IMAGENET	$\ell_2$	128/255	47.61	49.62	42.09	<b>-7.53</b>
	$\ell_\infty$	4/255	24.33	35.92	19.54	<b>-16.38</b>

BAT defense compromised

- ARC outperforms APGD across various datasets, norms, and network topologies



# Next Steps

- develop a complete theoretical framework for better understanding randomized ensembles of classifiers
- how can we design truly robust randomized ensembles in practice?



# Thank You!

code available at <https://github.com/hsndbk4/ARC>

**Acknowledgement:**

This work was supported by the Center for Brain-Inspired Computing (C-BRIC) and Artificial Intelligence Hardware (AIHW) funded by the Semiconductor Research Corporation (SRC) and the Defense Advanced Research Projects Agency (DARPA).

