# Adversarial Vulnerability of Randomized Ensembles
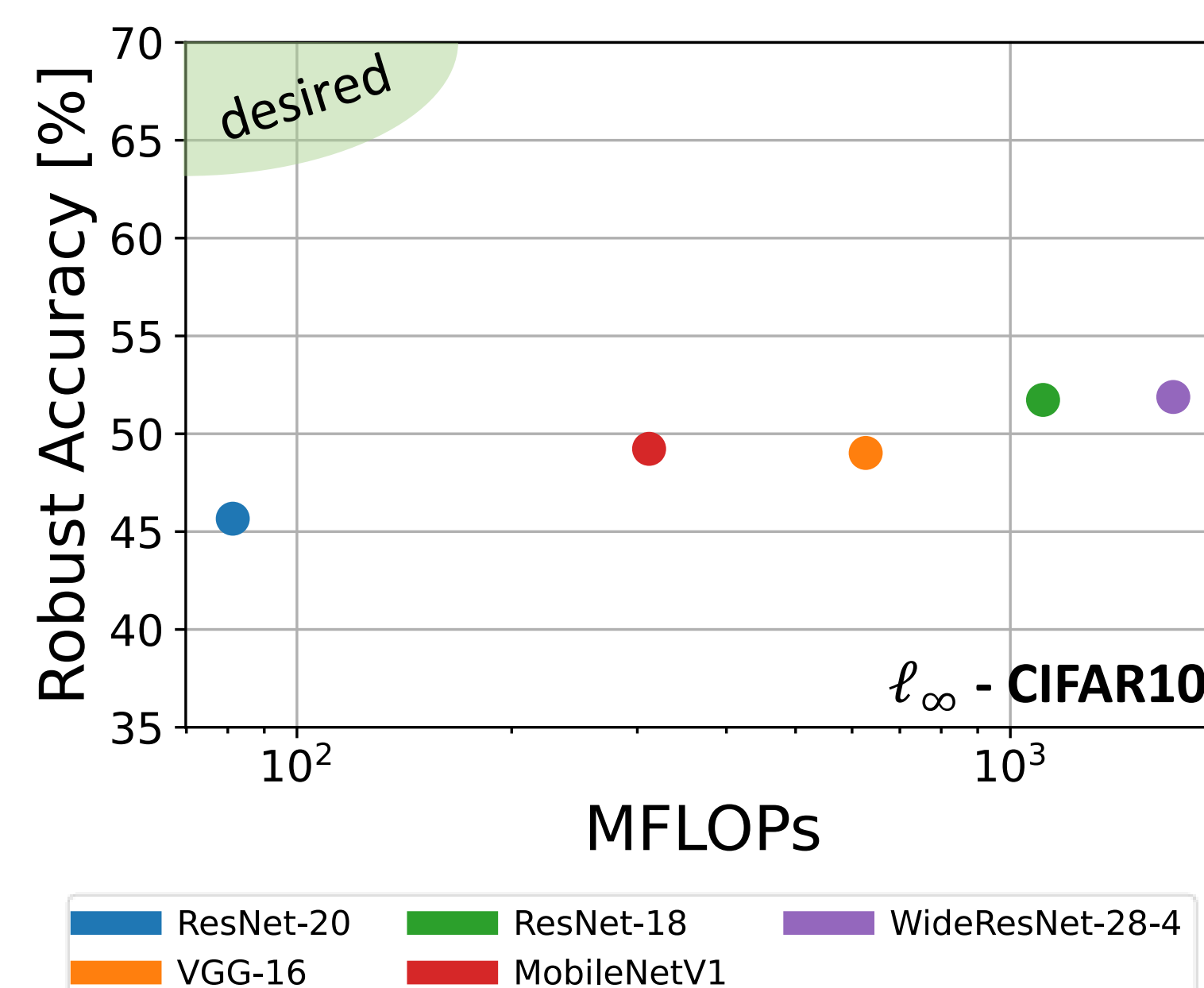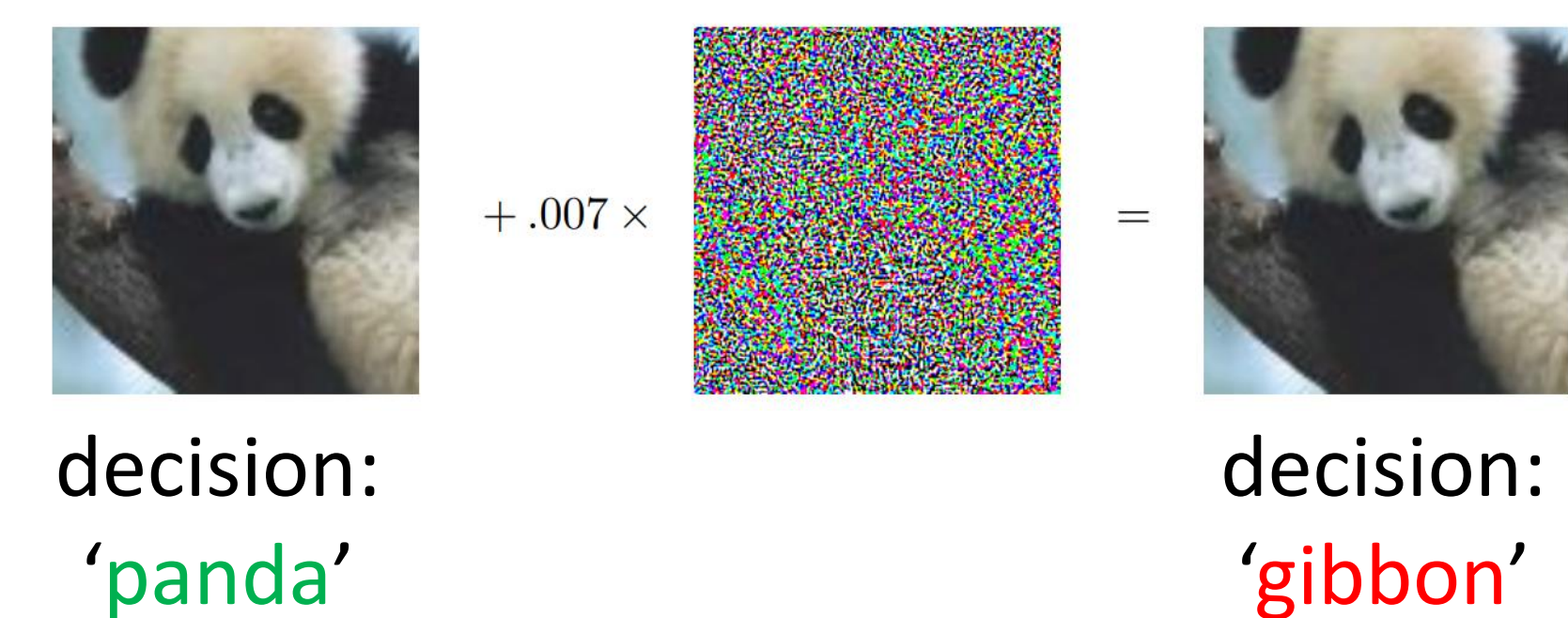
Hassan Dbouk & Naresh Shanbhag - *University of Illinois at Urbana-Champaign*
{hdbouk2,shanbhag}@illinois.edu
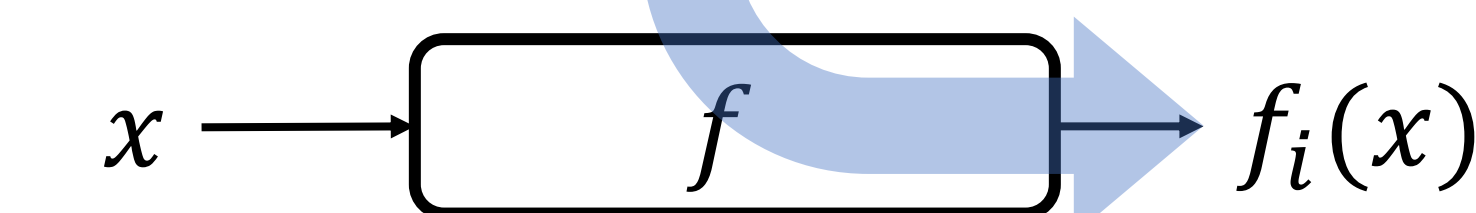
## Motivation

deep nets are <u>vulnerable</u>



decision: 'panda'  +.007 ×  =  decision: 'gibbon'

**robust** and **efficient** inference

robustness is <u>expensive</u>



ResNet-20 | ResNet-18 | WideResNet-28-4
VGG-16 | MobileNetV1

## Robustness via Randomized Ensembles

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

probabilities $\alpha_1 \; \alpha_2 \; \cdots \; \alpha_M$

classifiers $\quad f_1 \; f_2 \; \cdots \; f_M$

$x \longrightarrow f \longrightarrow f_i(x)$
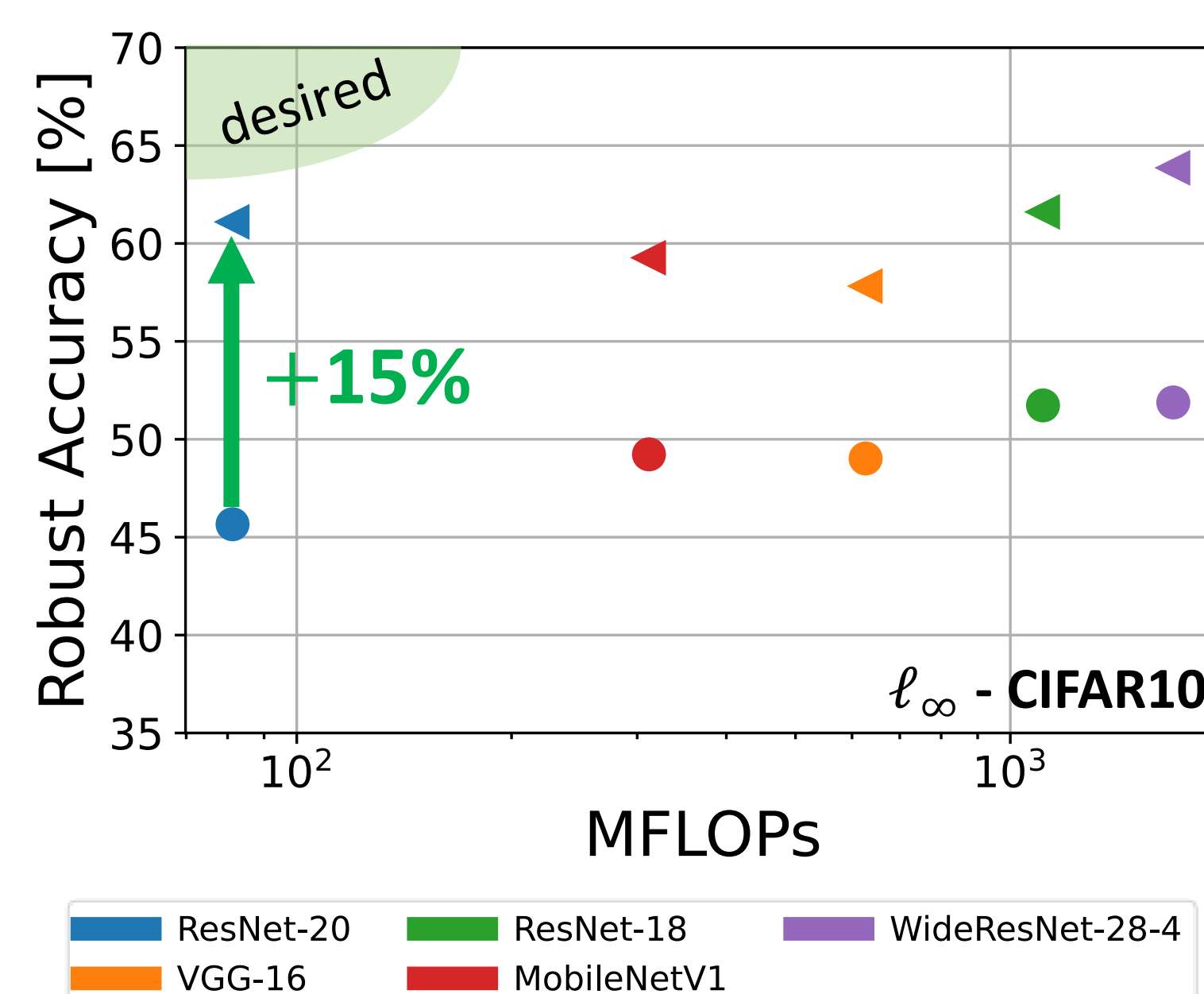
inference: pick **one** at random

**no** increase in # of FLOPS

using <u>two classifiers</u> trained via BAT [Pinot et al, 2020]



+15%

ResNet-20 | ResNet-18 | WideResNet-28-4
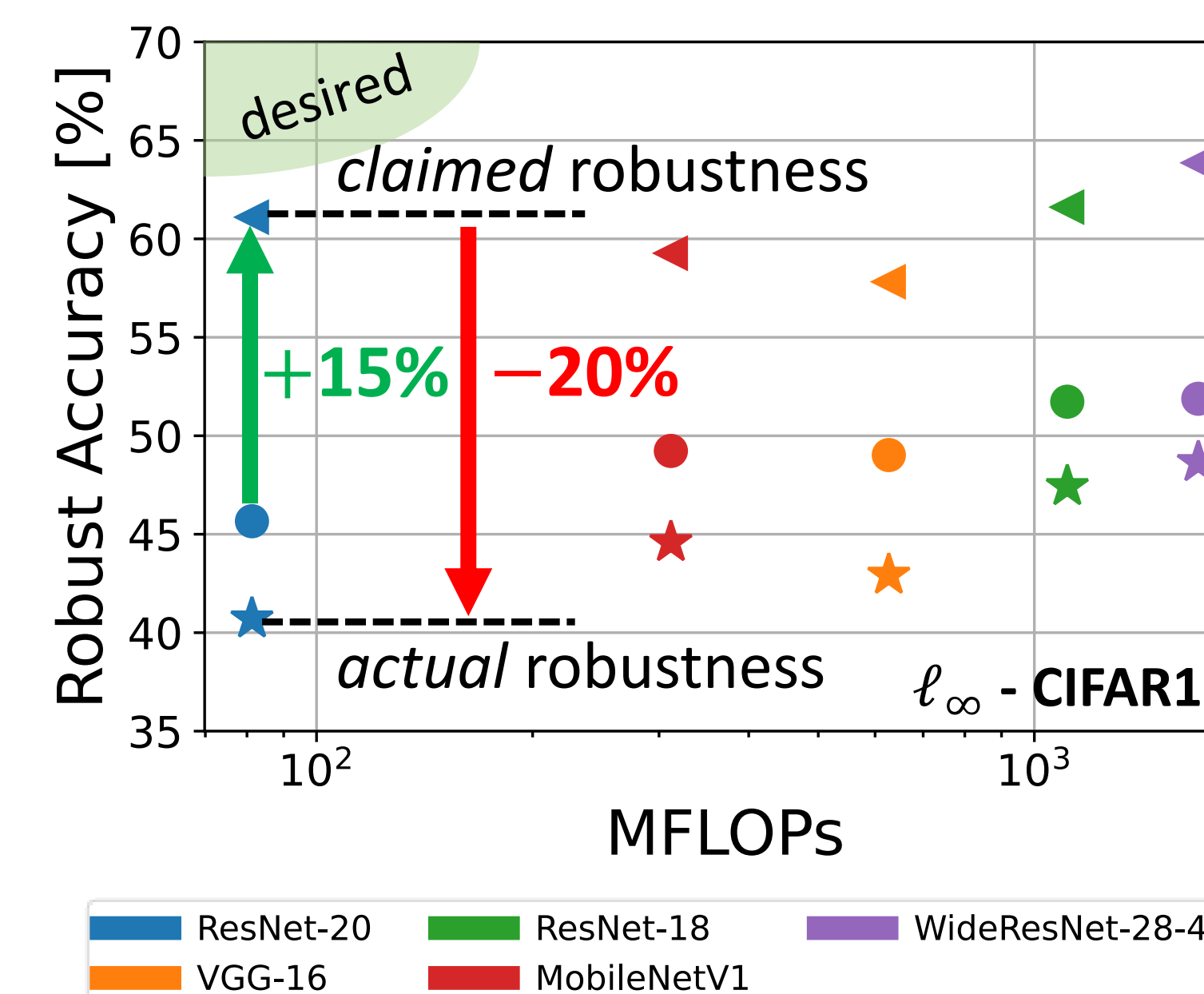VGG-16 | MobileNetV1

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

## Revealing the Vulnerability

**main** contributions

• show that adaptive PGD (APGD) is <u>ill-suited</u> for evaluating robustness

• propose a provably consistent and efficient adversarial <u>attack</u> algorithm – **ARC: A**ttacking **R**andomized ensembles of **C**lassifiers

• demonstrate that existing randomized ensembles defenses are in fact more vulnerable than standard AT

BAT defense **compromised**



*claimed* robustness
+15% −20%
*actual* robustness

ResNet-20 | ResNet-18 | WideResNet-28-4
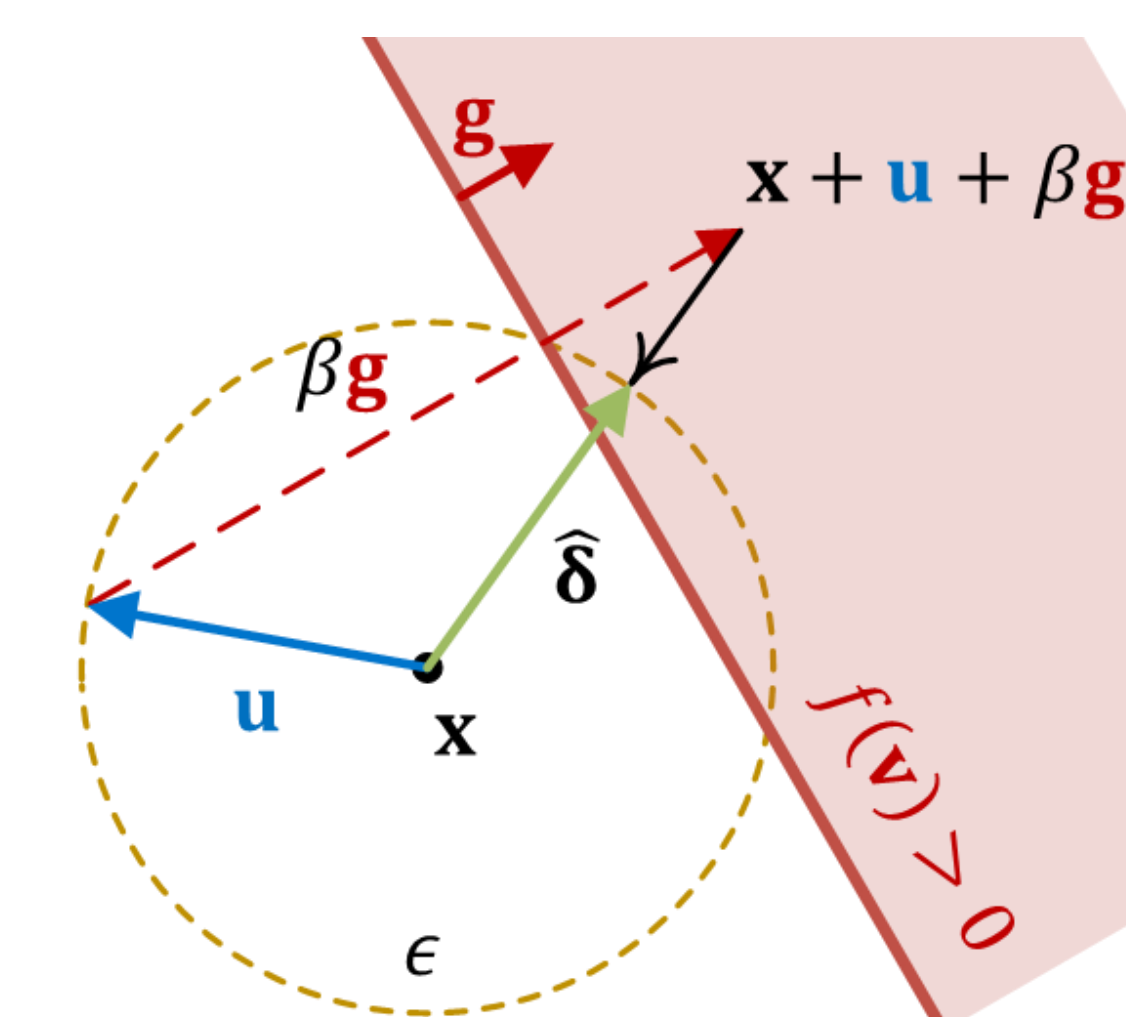VGG-16 | MobileNetV1

## 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:    $\mathbf{g} \leftarrow -y \frac{|\mathbf{w}_i|^{q-1} \odot \mathrm{sgn}(\mathbf{w}_i)}{\|\mathbf{w}_i\|_q^{q-1}}$
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** $\zeta \geq \epsilon \lor i = 1$ **then**
11:      $\beta \leftarrow \epsilon$
12:    **else**
13:      $\beta \leftarrow \frac{\epsilon}{\epsilon - \zeta} \left| \frac{y \mathbf{w}_i^\mathsf{T} \delta}{\|\mathbf{w}_i\|_q} + \zeta \right| + \rho$
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 <u>once</u>
• novel <u>adaptive step size</u> computation:



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

• extend to multiclass differentiable classifiers (e.g., neural nets)

**Theorem:** the ARC algorithm for BLCs is **consistent**

## Experimental Results – ARC vs. APGD

**Ensembles trained via BAT [Pinot et al., 2020]**

**varying networks - CIFAR-10**

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

**varying 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 | −14.51 |
| | $\ell_\infty$ | 8/255 | 53.55 | 65.99 | 52.01 | −13.98 |
| CIFAR-10 | $\ell_2$ | 128/255 | 62.43 | 69.21 | 55.44 | −13.77 |
| | $\ell_\infty$ | 8/255 | 45.66 | 61.10 | 40.71 | −20.39 |
| CIFAR-100 | $\ell_2$ | 128/255 | 34.60 | 41.91 | 28.92 | −12.99 |
| | $\ell_\infty$ | 8/255 | 22.29 | 33.37 | 17.45 | −15.92 |
| IMAGENET | $\ell_2$ | 128/255 | 47.61 | 49.62 | 42.09 | −7.53 |
| | $\ell_\infty$ | 4/255 | 24.33 | 35.92 | 19.54 | −16.38 |

• BAT defense **compromised**
• ARC **outperforms** APGD across various datasets, norms, and network topologies

## Summary & Next Steps

• demonstrated <u>theoretically</u> and <u>empirically</u> that **ARC** is better suited for evaluating the robustness of randomized ensembles

• existing randomized ensembles defenses are more **vulnerable** to $\ell_p$-bounded perturbations than adversarially trained models.

• our work advocates the need for improved randomized defense methods including <u>certifiable defenses</u>