

# On the Robustness of Randomized Ensembles to Adversarial Perturbations

**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 ×



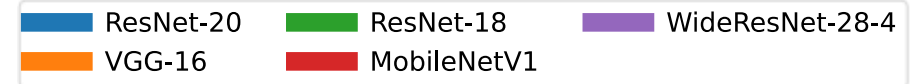
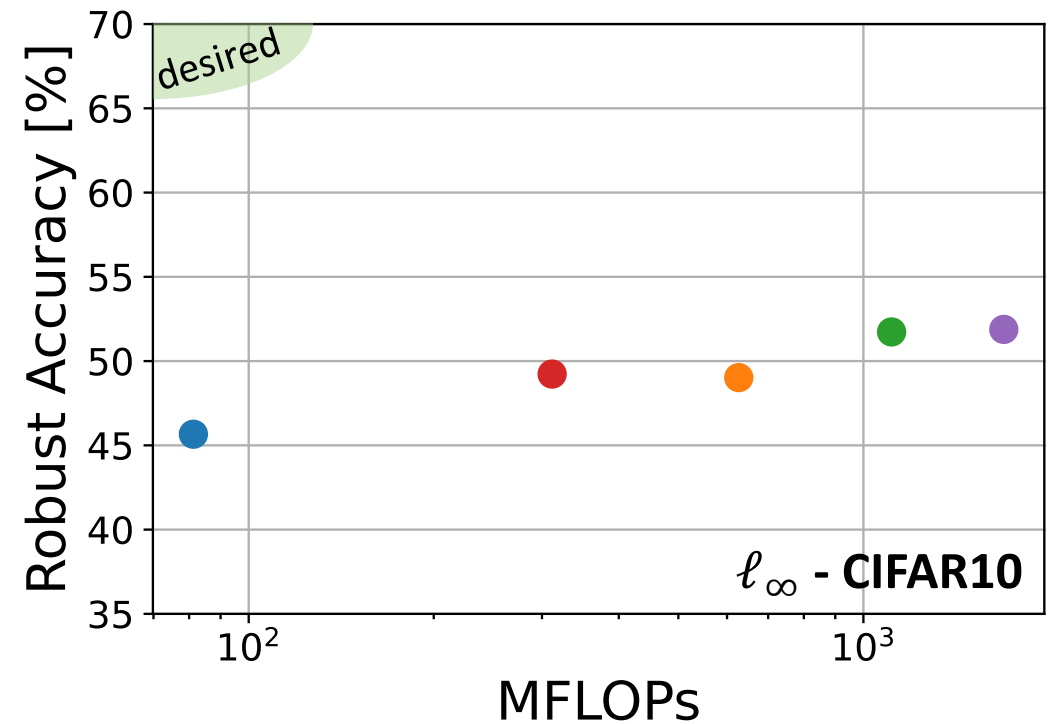
=

adversarial sample



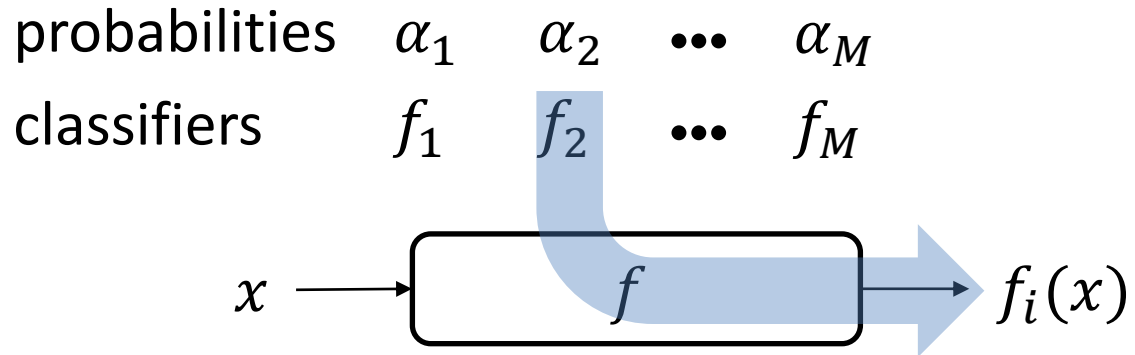
decision: 'gibbon'

robustness is expensive



# Robustness via Randomized Ensembles

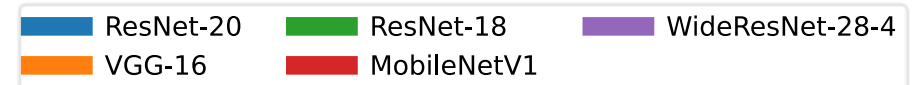
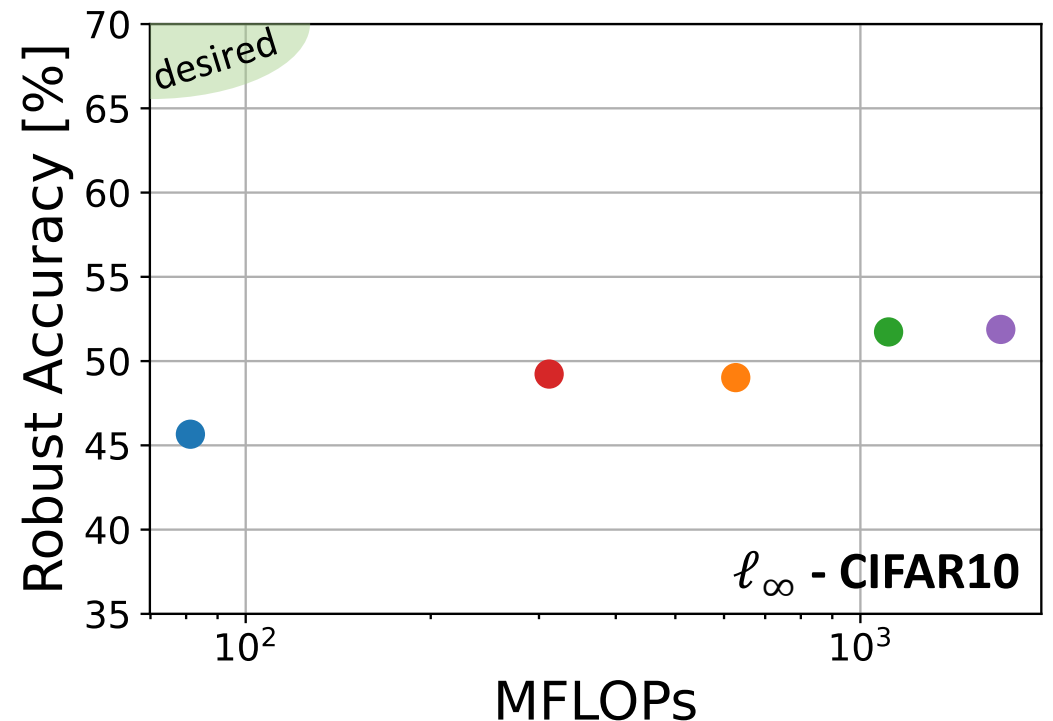
**multiple** classifiers  $f_1, \dots, f_M$



inference: pick **one** at random

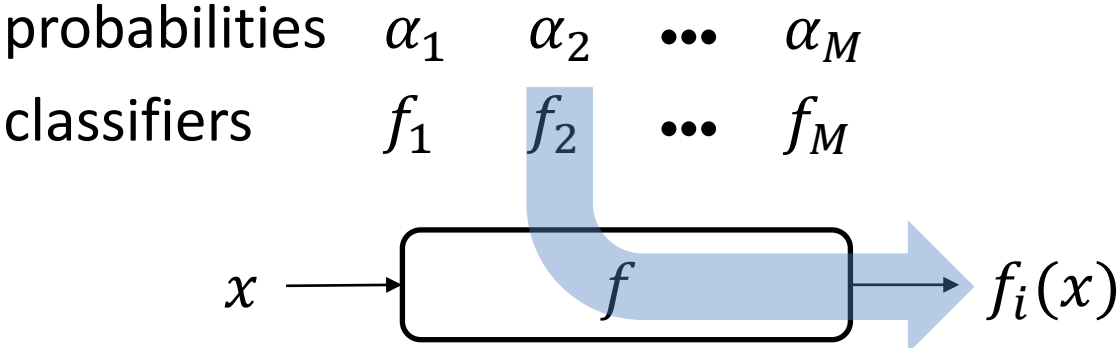
**no** increase in # of FLOPs

robustness is expensive



# Boosted Adversarial Training (BAT) [Pinot et al, ICML'20]

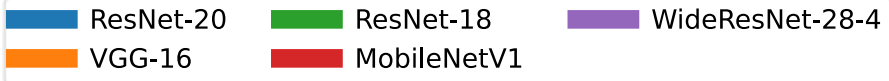
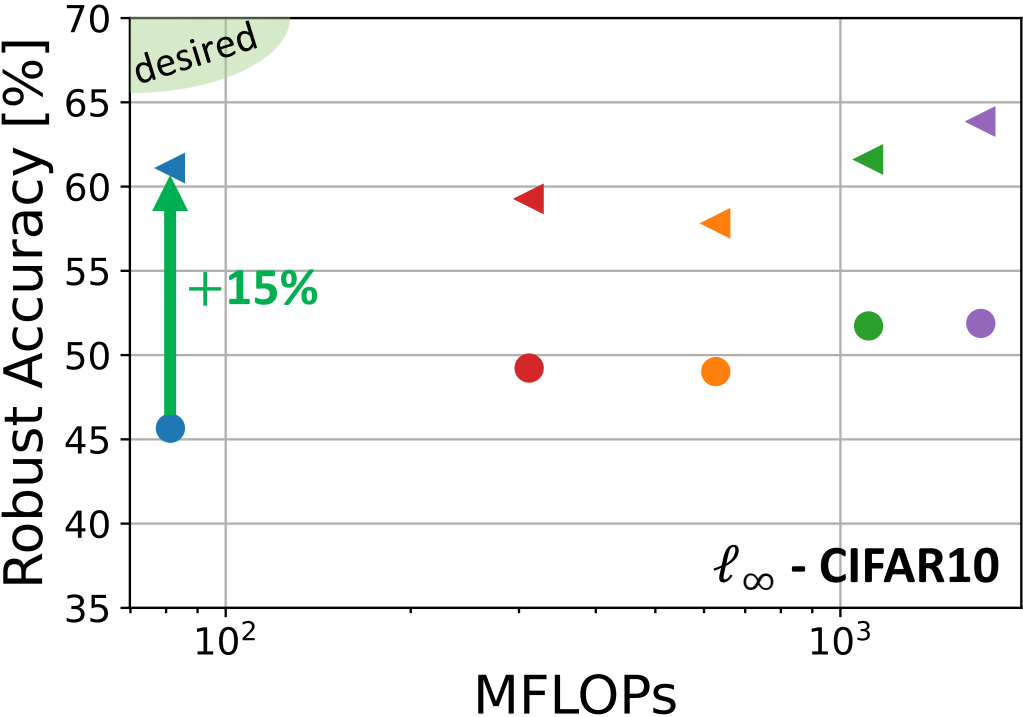
**multiple** classifiers  $f_1, \dots, f_M$



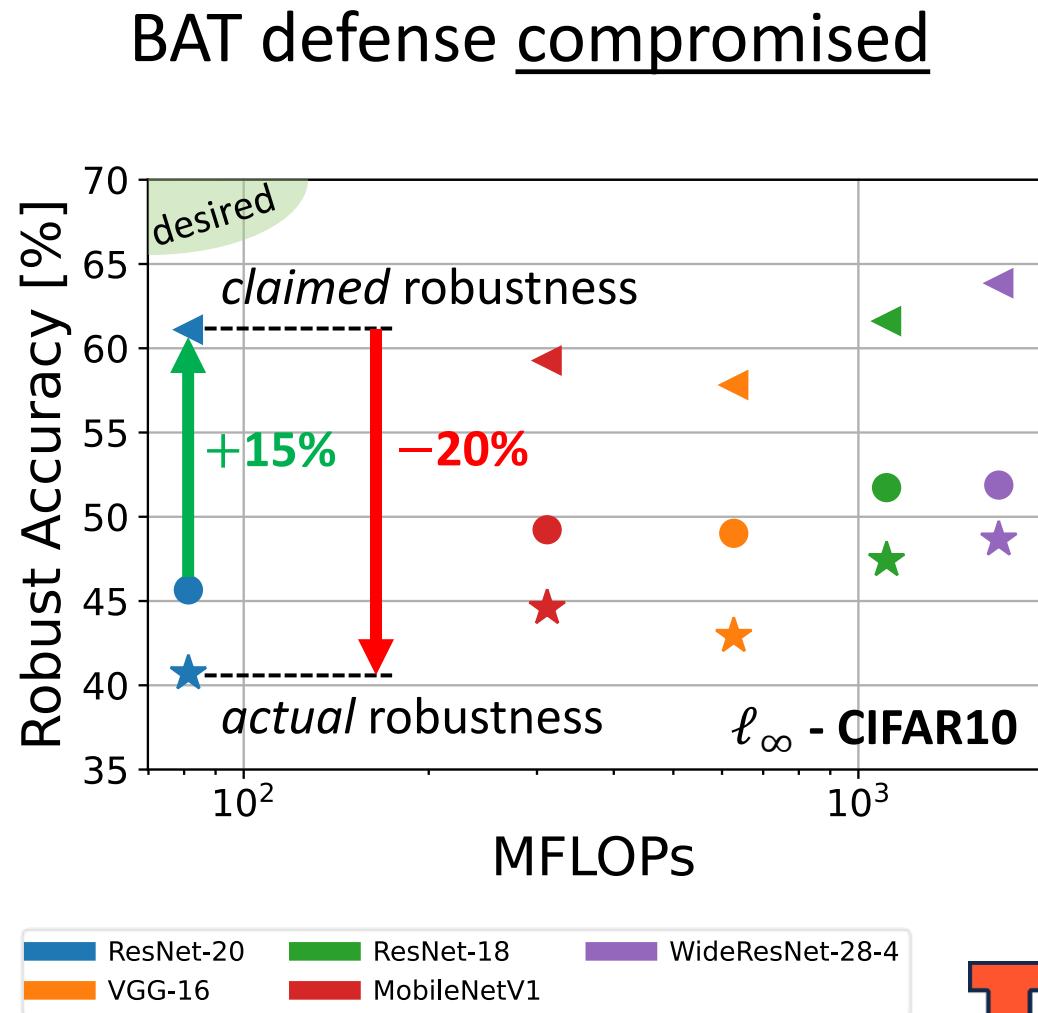
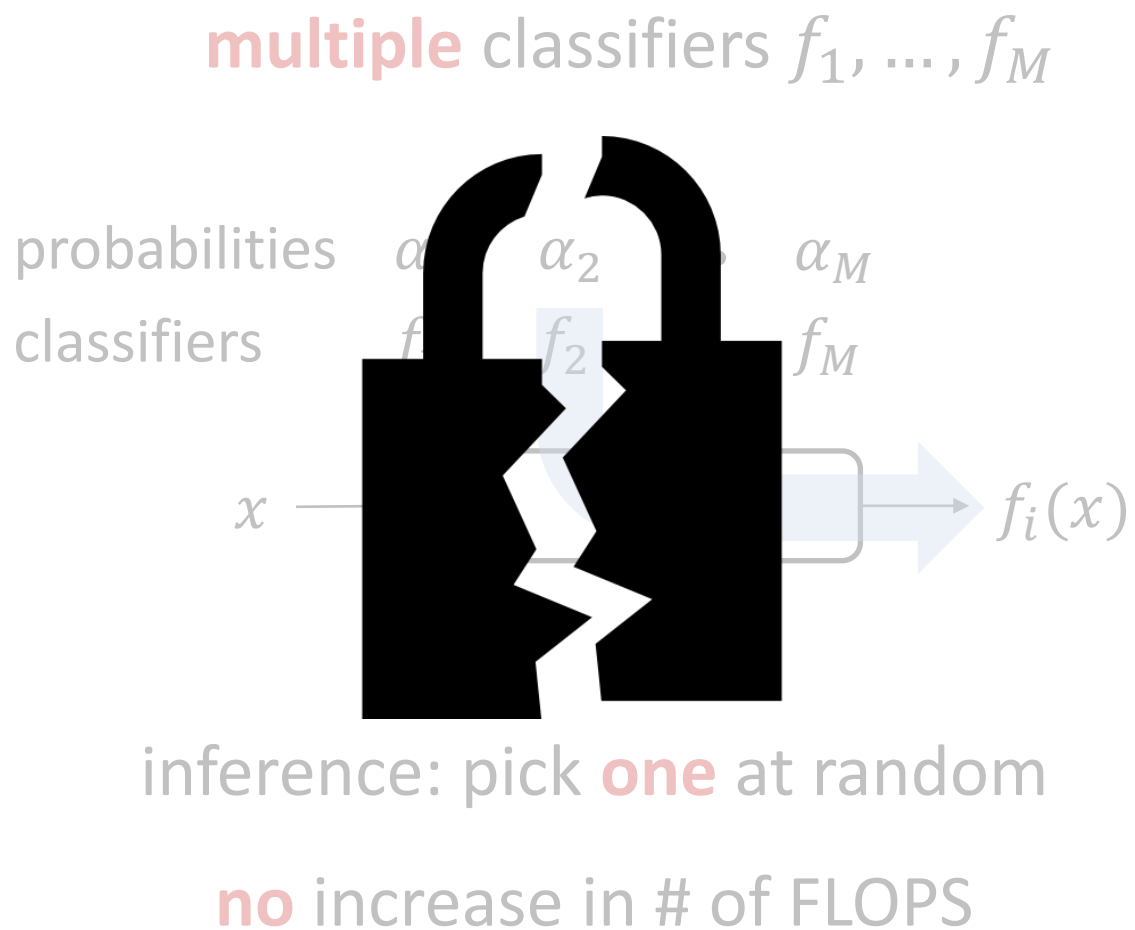
inference: pick **one** at random

**no** increase in # of FLOPs

using two classifiers trained via BAT



# Our Prior Work [ICML'22]: Revealing the vulnerability



# Fundamental Questions on RECs

- when does randomization help in improving robustness?
- what are the limits of RECs?
- how to find the optimal sampling probability?
- how do we train robust RECs in practice?

# Summary of Theoretical Contributions

- derived a theoretical framework for analyzing the robustness of RECs
- derived fundamental results on:
  - necessary and sufficient conditions for RECs to be useful

**Thm 1.**  $\forall f_1$  &  $f_2$  with adv. risks  $\eta_1$  &  $\eta_2$ . If:  $\mathbb{P}\{z \in \mathcal{R}_1\} > |\eta_1 - \eta_2| \rightarrow \eta(\alpha^*) = \frac{1}{2}(\eta_1 + \eta_2 - \mathbb{P}\{z \in \mathcal{R}_1\})$

- theoretical robustness limits of RECs

**Thm 2.**  $\forall \{f_i\}_{i=1}^M$  with adv. risks  $\{\eta_i\}_{i=1}^M$ . We can tightly bound the REC adv. risk:  $\min_{k \in [M]} \left\{ \frac{\eta_k}{k} \right\} \leq \eta(\alpha) \leq \eta_M$

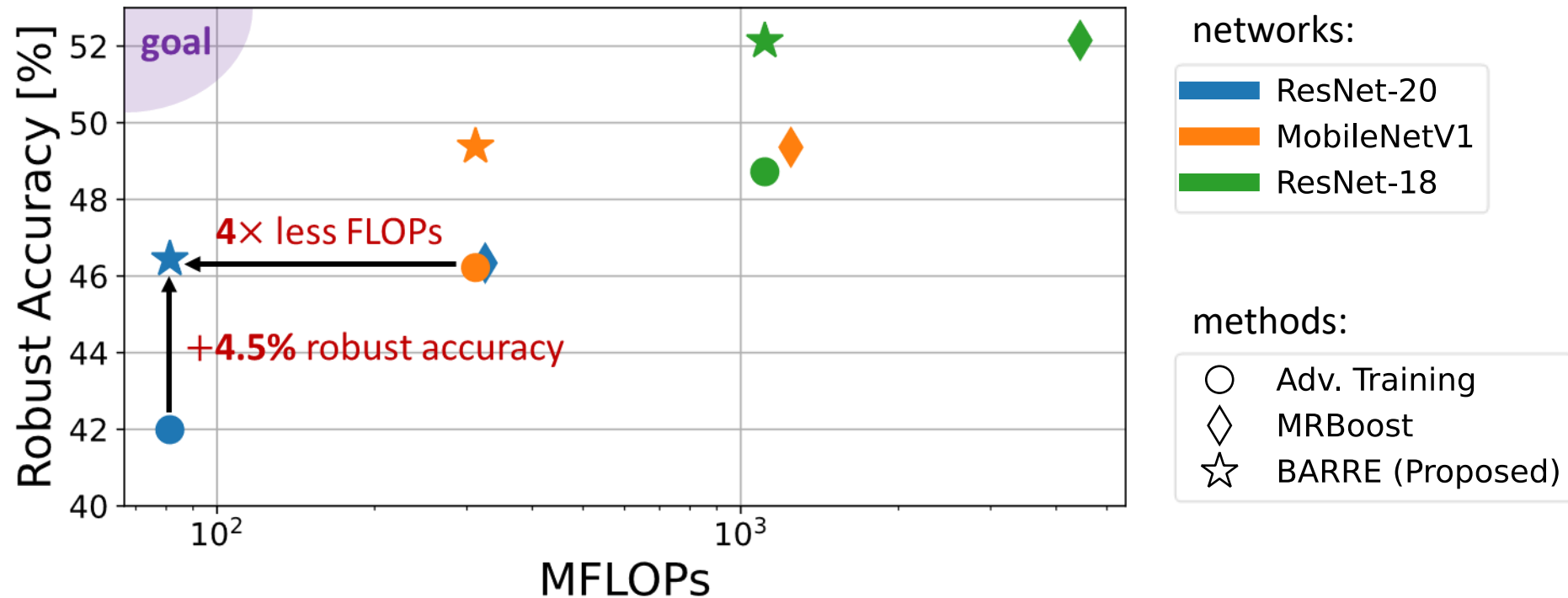
- efficient and optimal methods for finding the optimal sampling probability

**Thm 3.** The OSP algorithm output  $\alpha_T$  satisfies  $0 \leq \hat{\eta}(\alpha_T) - \hat{\eta}(\alpha^*) \xrightarrow{T \rightarrow \infty} 0$ , for all initial conditions





# Proposed Training Algorithm: BARRE



- robustness against  $\ell_\infty$  norm-bounded adversaries – CIFAR-10
- BARRE: drastically improve **robust accuracy** while maintaining **complexity**

# Thank You!

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

## **Acknowledgement:**

This work was supported by the Center for the Co-Design of Cognitive Systems (CoCoSys) funded by the Semiconductor Research Corporation (SRC) and the Defense Advanced Research Projects Agency (DARPA), and SRC's Artificial Intelligence Hardware (AIHW) program.

