

Guest Lecture:

Mixing Neural Network Classifiers to Balance Accuracy and Adversarial Robustness

Presenter: Yatong Bai yatong_bai@berkeley.edu

May 19, 2024

About Myself

- Rising 5th-year Ph.D. candidate at UC Berkeley advised by Professor Somayeh Sojoudi.
- Research focus:
 - Reconciling adversarial robustness and accuracy of classification models.
 - Efficient audio generation through consistency models.
- Teaching:
 - Convex optimization and approximation.

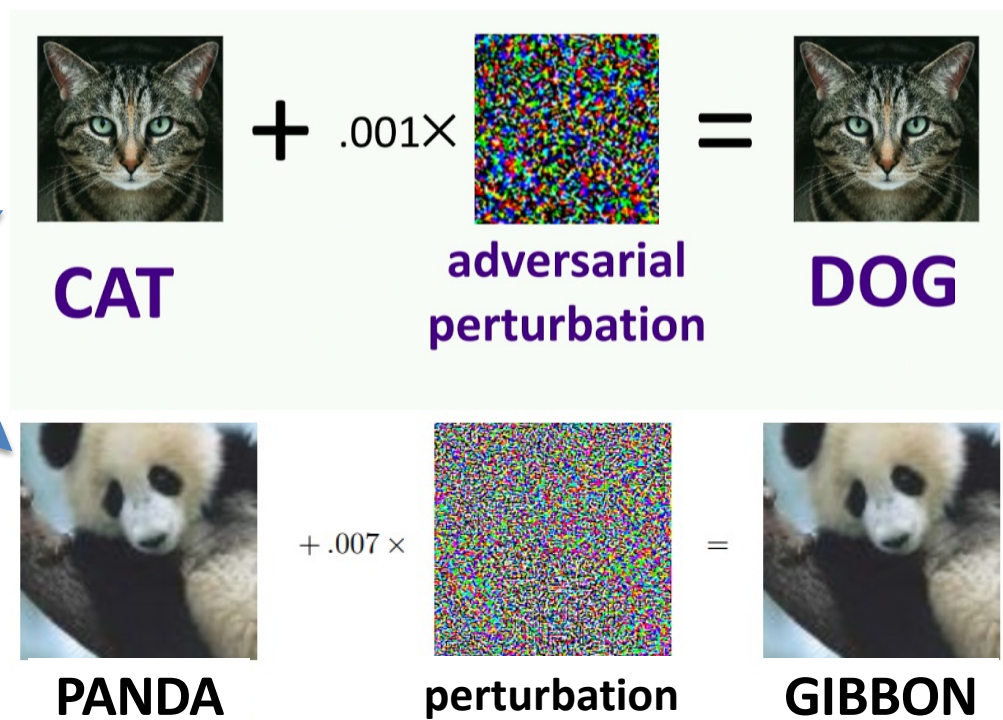
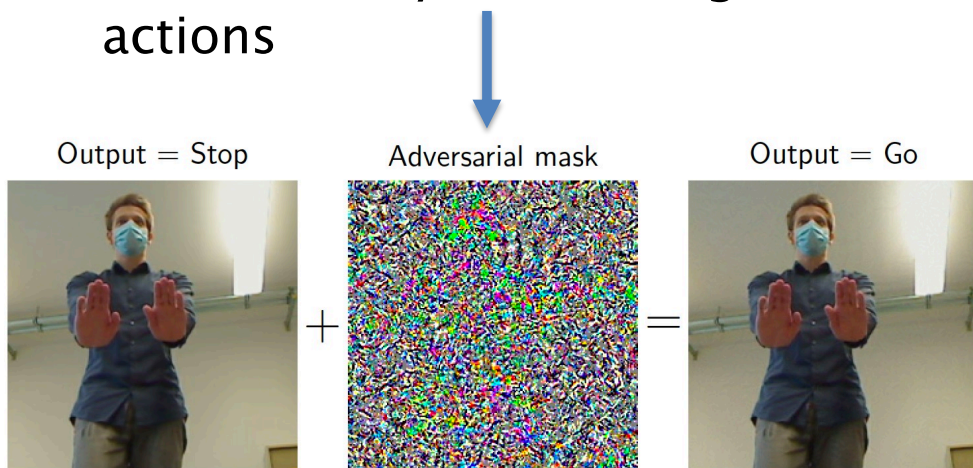


Overview of This Presentation

- Brief intro to adversarial robustness.
- Improving the accuracy–robustness trade–off.
 - Mixing classifiers to balance robustness and accuracy.
 - Adaptive Smoothing: adaptive mixing ratio.
<https://arxiv.org/abs/2301.12554>
 - MixedNUTS: mix in a nonlinear fashion.
<https://arxiv.org/abs/2402.02263>

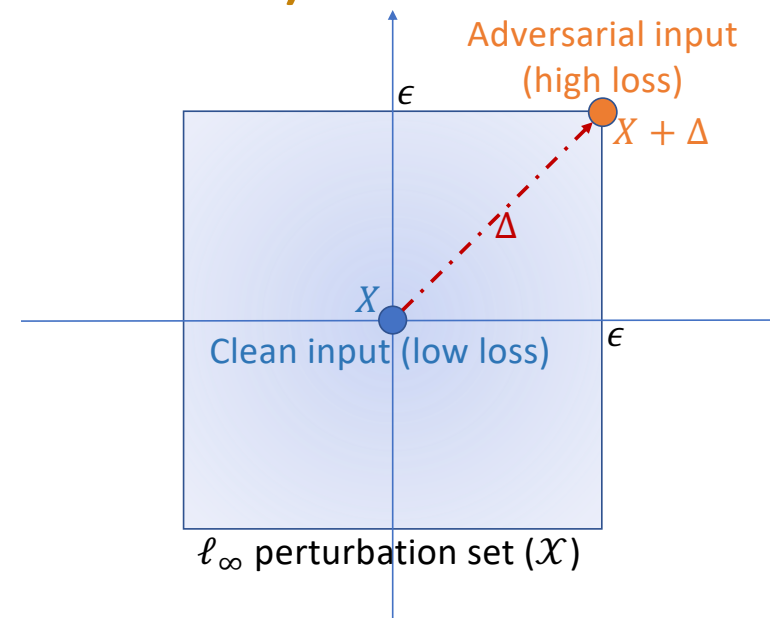
Adversarial Robustness

- Neural networks are vulnerable
 - Small input perturbations elicit unexpected outputs.
- For classifiers: misclassifications.
- For control systems: dangerous actions



Adversarial example generation (An optimization formulation)

- We need a budget for the attack, since the adversarial perturbations should be imperceptible by human.
 - A common uncertainty set is an ℓ_∞ -norm-bounded additive set with radius ϵ :
 - I. e., a cube around each clean input.



- The adversarial examples are usually generated via the following optimization problem:

$$\max_{\delta: x+\delta \in \mathcal{X}} \underset{\text{Loss fn}}{\ell} \left(\underbrace{g(x+\delta)}_{\text{NN output for attacked input}}, \underbrace{Y}_{\text{Target output}} \right), \text{ where } g \text{ represents the NN as a function.}$$

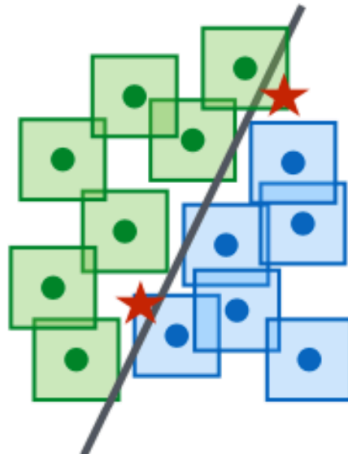
Defending attacks -- Adversarial training (Robust Optimization)

- One defense method: Adversarial training (train with adversarial data) [Madry et al., 2018, Goodfellow et al., 2015] .
 - Train robust models via robust optimization. For an uncertainty set \mathcal{X} , solve the optimization problem

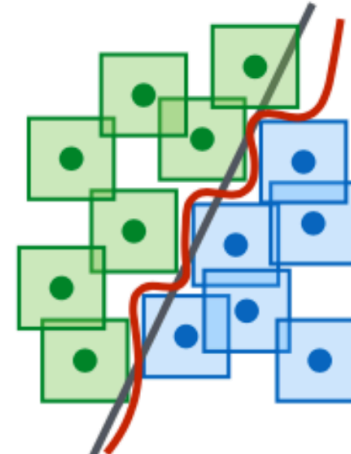
$$\underbrace{\min_{\theta}}_{\text{Optimize NN weights}} \left(\underbrace{\max_{\delta: x+\delta \in \mathcal{X}}}_{\text{Generate attack}} \ell(g_{\theta}(x+\delta), Y) \right) + \underbrace{r_{\theta}}_{\text{Regularization}} \quad (1)$$



Nominal Decision Boundary



Doesn't Separate l_{∞} Norm Balls



Robust Decision Boundary

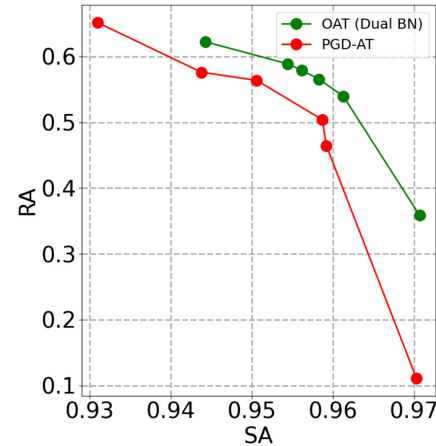
Comparison of decision boundaries of standard training and adversarial training [Madry et al., 2018].

Alternative methods:

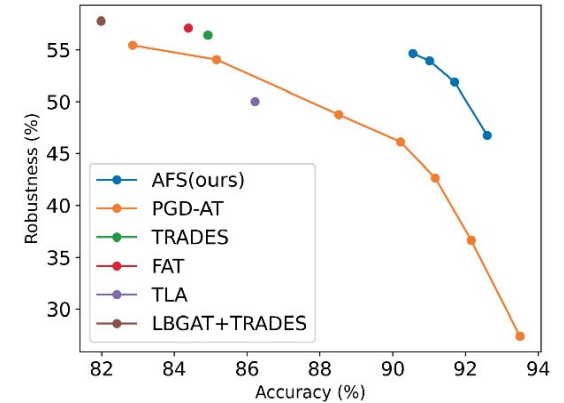
- TRADES, Randomized Smoothing.

Accuracy–Robustness Trade–Off

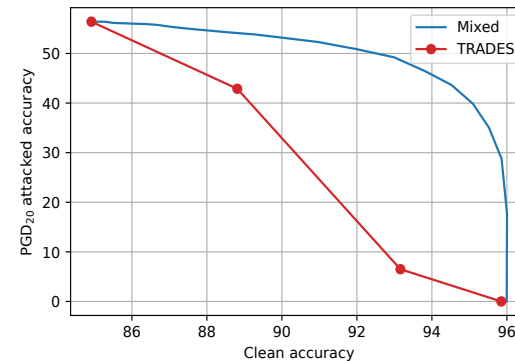
- Robust models often sacrifice clean accuracy.
- Theoretically, robust generalization needs much more training data.
- Existing methods for alleviating the trade-off:
 - Additional real/synthetic training data;
 - Attack purification;
 - Alternative training loss functions.



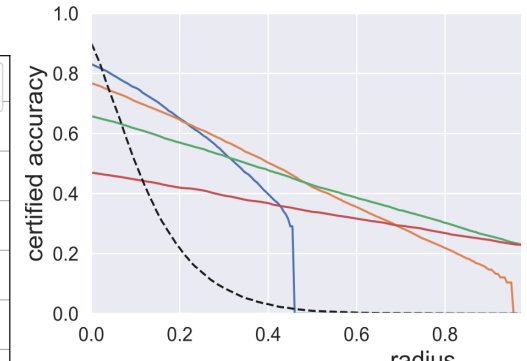
Once-for-All Adversarial Training: In-Situ Tradeoff between Robustness and Accuracy for Free



Towards Both Accurate and Robust Neural Networks Without Extra Data



Improving the Accuracy–Robustness Trade-Off of Classifiers via Adaptive Smoothing

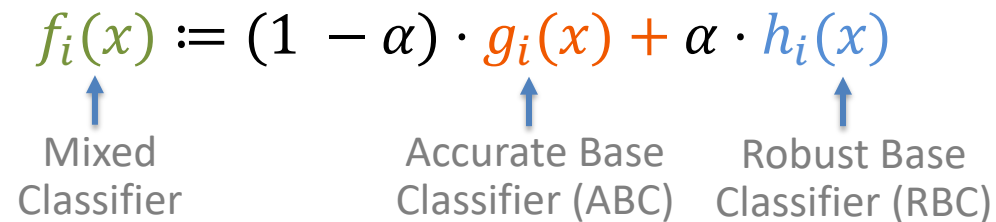


Certified Adversarial Robustness via Randomized Smoothing

Mixing Classifiers for Better Trade-Off

- What if we combine the wisdom of an **accurate model** and a **robust model**?
- Specifically, we “mix” their outputs, resulting in a **mixed classifier**.

$$f_i(x) := (1 - \alpha) \cdot g_i(x) + \alpha \cdot h_i(x)$$

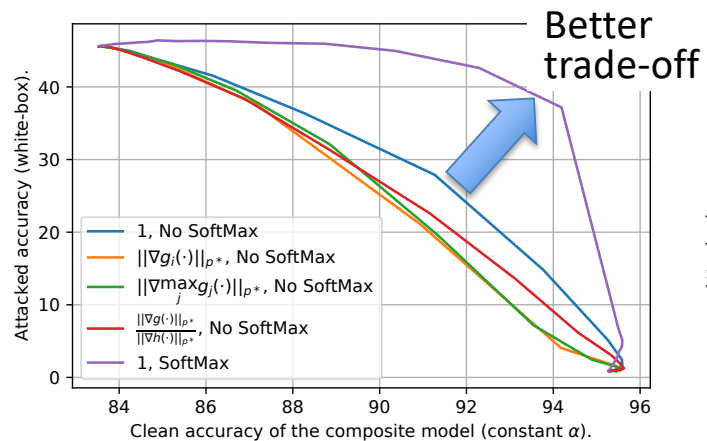


Mixed Classifier Accurate Base Classifier (ABC) Robust Base Classifier (RBC)

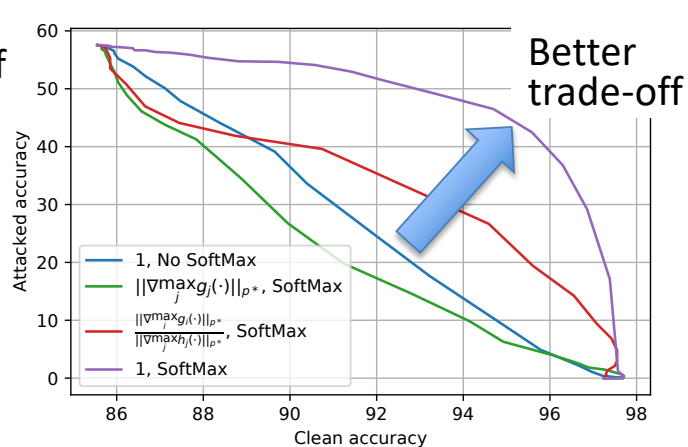
- Should we mix the logits or probabilities?
 - Classifiers often use a “Softmax” operation to convert “logits” $(-\infty, +\infty)$ to prediction probabilities $(0, 1)$.

Empirically comparing the design choices

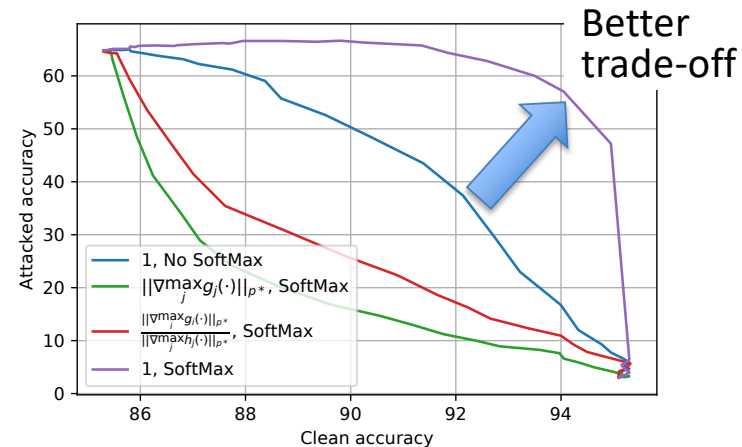
- We compare the cases with various values of α via the clean accuracy versus attacked accuracy plot:



ResNet18+AT, l_∞



ConvNeXT+TRADE, l_∞



ResNet18+AT, l_2

Figure 1: Adaptive PGD₁₀ accuracy versus clean accuracy for the three different choices of $R(x)$ on CIFAR-10.

- Blue: smoothing with logits. Purple: smoothing with probabilities.
- Conclusion: smoothing should be done on probabilities.

Mixing Probabilities is Better

- Conclusion: we should mix the base classifiers' **prediction probabilities**.
- The resulting class-wise mixing formulation is:

The diagram illustrates the derivation of the mixed classifier formula. At the top, the text "Convert back to logits" has a blue arrow pointing down to the \log function in the equation. Below the equation, the text "Mixed Classifier" has a blue arrow pointing up to the $f_i(x)$ term. The text "Accurate Base Classifier (ABC)" has a blue arrow pointing up to the $g(x)_i$ term, and the text "Robust Base Classifier (RBC)" has a blue arrow pointing up to the $h(x)_i$ term. The word "Softmax" is positioned above the equation, with two blue arrows pointing down to the σ functions in the expression.

$$f_i(x) := \log \left((1 - \alpha) \cdot \sigma \circ g(x)_i + \alpha \cdot \sigma \circ h(x)_i \right)$$

Mixed Classifier

Accurate Base Classifier (ABC)

Robust Base Classifier (RBC)

Softmax

Convert back to logits

Intuition for mixing the probabilities

- The robust classifier $h(\cdot)$ is typically smooth or Lipschitz, and we want $g_{\text{CNN}}^\alpha(\cdot)$ to inherit these properties.
- The accurate classifier $g(\cdot)$ is in general non-smooth and non-robust.
- If $g(\cdot) \in [0, 1]$ (probabilities), then the "level of incorrectness" can be bounded. It is then possible for the smoothness of $h(\cdot)$ to overshadow the turbulence of $g(\cdot)$, ultimately making $g_{\text{CNN}}^\alpha(\cdot)$ robust.
-- Will present a Lemma to formalize this.
- If $g(\cdot) \in \mathbb{R}$ (logits), then it can be arbitrarily unsmooth. $h(\cdot)$ may not be possible to correct $g(\cdot)$.

Certiably robust with a margin (Theoretically guaranteed robustness)

To facilitate the proof for certified robust radii, we first introduce the notion "robust with a margin".

Definition

Consider an arbitrary input $x \in \mathbb{R}^d$ and let $y = \arg \max_i h_i(x)$, $\mu \in [0, 1]$, and $r \geq 0$. Then, $h(\cdot)$ is said to be certiably robust at x with margin μ and radius r if $h_y(x + \delta) \geq h_i(x + \delta) + \mu$ for all $i \neq y$ and all $\delta \in \mathbb{R}^d$ such that $\|\delta\|_p \leq r$.

Lemma

Let $x \in \mathbb{R}^d$ and $r \geq 0$.

If it holds that $\alpha \in [\frac{1}{2}, 1]$ and $h(\cdot)$ is certiably robust at x with margin $\frac{1-\alpha}{\alpha}$ and radius r , then the smoothed classifier $g_{\text{CNN}}^\alpha(\cdot)$ is robust in the sense that $\arg \max_i g_{\text{CNN},i}^\alpha(x + \delta) = \arg \max_i h_i(x)$ for all $\delta \in \mathbb{R}^d$ such that $\|\delta\|_p \leq r$.

- Intuition: if $h(\cdot)$ is robust and confident, then it can override whatever $g(\cdot)$ predicts.

Certiably robust with a margin -- Proof

Lemma

(Restated.) If it holds that $\alpha \in [\frac{1}{2}, 1]$ and $h(\cdot)$ is certiably robust at x with margin $\frac{1-\alpha}{\alpha}$ and radius r , then $\arg \max_i g_{\text{CNN},i}^\alpha(x + \delta) = \arg \max_i h_i(x)$ for all $\delta \in \mathbb{R}^d$ such that $\|\delta\|_p \leq r$.

Proof

Since $\alpha \in [\frac{1}{2}, 1]$, it holds that $\frac{1-\alpha}{\alpha} \in [0, 1]$.

Suppose that $h(\cdot)$ is certiably robust at x with margin $\frac{1-\alpha}{\alpha}$ and radius r .

Let $y = \arg \max_i h_i(x)$. Consider an arbitrary $i \in [c] \setminus \{y\}$ and $\delta \in \mathbb{R}^d$ such that $\|\delta\|_p \leq r$. It holds that

$$\begin{aligned} \exp(g_{\text{CNN},y}^\alpha(x + \delta)) - \exp(g_{\text{CNN},i}^\alpha(x + \delta)) &= (1 - \alpha)(g_y(x + \delta) - g_i(x + \delta)) + \alpha(h_y(x + \delta) - h_i(x + \delta)) \\ (\text{Because } g_i(x + \delta) \in [0, 1]) \quad &\geq (1 - \alpha)(0 - 1) + \alpha(h_y(x + \delta) - h_i(x + \delta)) \\ &\geq (\alpha - 1) + \alpha \left(\frac{1-\alpha}{\alpha}\right) = 0. \end{aligned}$$

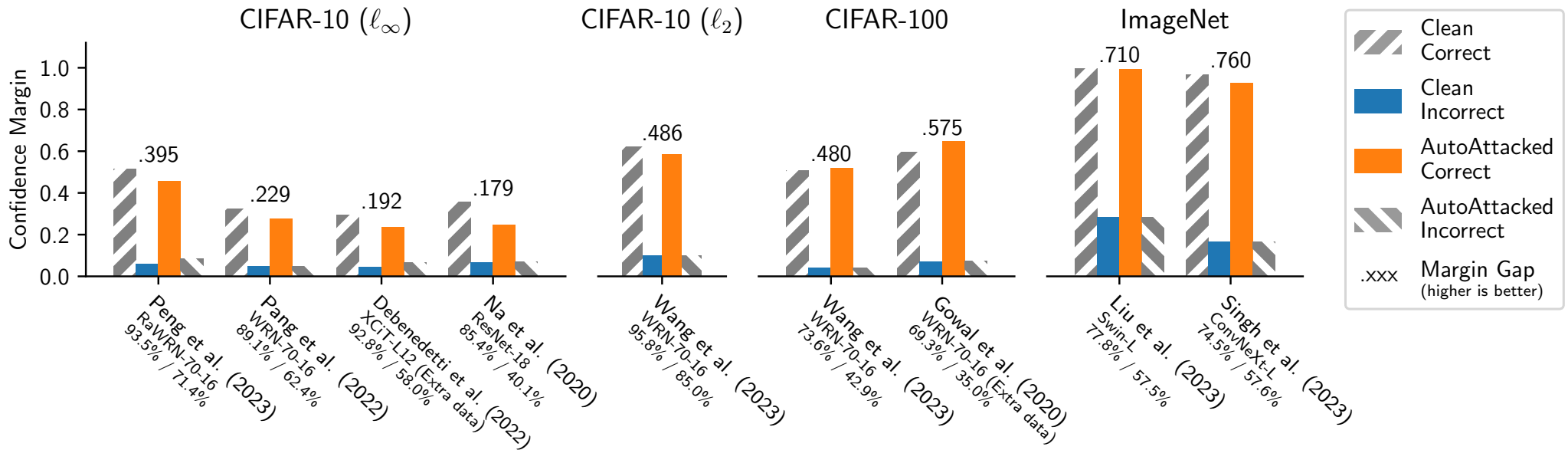
Thus, it holds that $g_{\text{CNN},y}^\alpha(x + \delta) \geq g_{\text{CNN},i}^\alpha(x + \delta)$ for all $i \neq y$, and thus $\arg \max_i g_{\text{CNN},i}^\alpha(x + \delta) = y = \arg \max_i h_i(x)$. □

Mechanism for Improved Accuracy Trade-Off

- Empirically robust models are more confident when correct than when incorrect, even on attacked data.

Definition 1. Consider a model $h : \mathbb{R}^d \rightarrow \mathbb{R}^c$, an arbitrary input $x \in \mathbb{R}^d$, and its associated predicted label $\hat{y} \in [c]$. The *confidence margin* is defined as $m_h(x) := \sigma \circ h_{\hat{y}}(x) - \max_{i \neq \hat{y}} \sigma \circ h_i(x)$.

- Some examples (SOTA models on various datasets):



Mechanism for Improved Accuracy Trade-Off

- When α is slightly greater than 0.5:
 - On clean data, $g(\cdot)$ is better than $h(\cdot)$.
Since $h(\cdot)$ is unconfident when making mistakes, it can be corrected by $g(\cdot)$;
 - On attacked data, $h(\cdot)$ is better than $g(\cdot)$.
Since $h(\cdot)$ is confident in correct predictions, it can overcome $g(\cdot)$.

↑
Accurate Base
Classifier (ABC)

↑
Robust Base
Classifier (RBC)

- When α is slightly greater than 0.5:
 - On clean data, $g(\cdot)$ is better than $h(\cdot)$.
Since $h(\cdot)$ is unconfident when making mistakes, it can be corrected by $g(\cdot)$;
 - On attacked data, $h(\cdot)$ is better than $g(\cdot)$.
Since $h(\cdot)$ is confident in correct predictions, it can overcome $g(\cdot)$.

Adaptive Smoothing: Flexible Mixing Ratio

- Recall the mixed classifier formulation:

$$f_i(x) := \log \left((1 - \alpha) \cdot \sigma \circ g(x)_i + \alpha \cdot \sigma \circ h(x)_i \right)$$

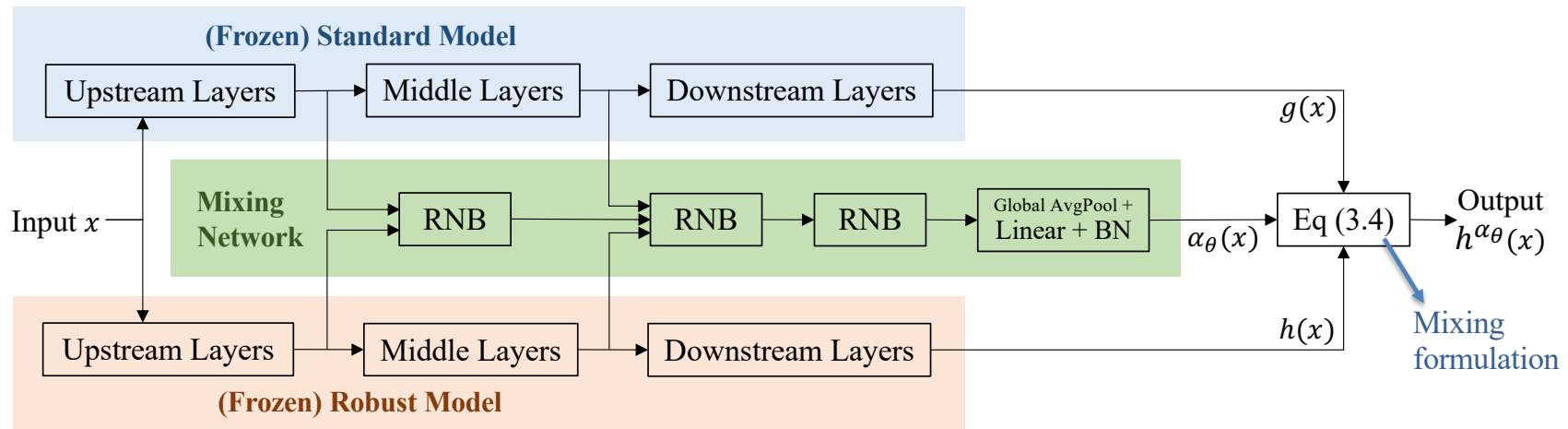
Diagram illustrating the mixed classifier formulation:

- The term $f_i(x)$ is labeled as the **Mixed Classifier**.
- The operation \log is labeled as **Convert back to logits**.
- The operation σ is labeled as **Softmax**.
- The term $g(x)_i$ is labeled as the **Accurate Base Classifier (ABC)**.
- The term $h(x)_i$ is labeled as the **Robust Base Classifier (RBC)**.

- It makes sense to **make the mixing ratio α a function of x** .

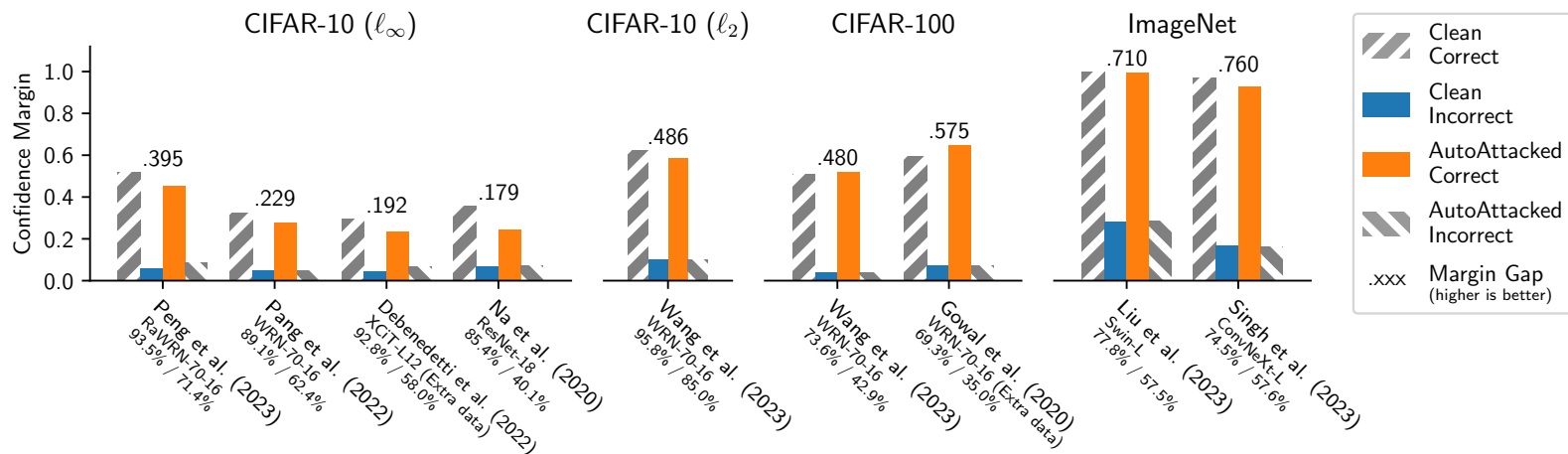
Adaptive Smoothing: Flexible Mixing Ratio

- It makes sense to make the mixing ratio α a function of x .
 - Make $\alpha(x)$ **small** and prefer the **ABC** $g(x)$ when x is **natural** (no attack).
 - Make $\alpha(x)$ **large** and prefer the **RBC** $h(x)$ when x is **adversarial**.
- Parameterizing $\alpha(x)$: an additional neural network module.



MixedNUTS: Nonlinear Mixed Classifier

- **Recall:** Mixed classifiers rely on the RBC $h(\cdot)$'s benign confidence properties.
 - More confident in correct examples than incorrect ones.



- Confidence can be adjusted without changing predictions.
 - (e.g., temperature scaling).
- Can we augment the benign properties to improve the mixed classifier?

MixedNUTS: Nonlinear Mixed Classifier

- How to augment the benign properties?
- Apply a non-linear transformation $M(\cdot)$ to RBC $h(\cdot)$'s logits before Softmax and mixing.
 - Notation: $h^M(x) = M(h(x))$.
 - Temperature scaling is a special case where $M(\cdot)$ is linear.
- Apply temperature scaling to ABC $g(\cdot)$'s logits before Softmax and mixing.
 - Ablation study shows that zero temperature (one-hot probabilities) works the best.

MixedNUTS: Nonlinear Mixed Classifier

- Goal: optimize $M(\cdot)$'s clean accuracy for a given robust accuracy r_f .

$$\begin{aligned} & \max_{M \in \mathcal{M}, \alpha \in [1/2, 1]} \mathbb{P}_{(X, Y) \sim \mathcal{D}} \left[\arg \max_i f_i^M(X) = Y \right] \quad (2) \\ & \text{s. t. } \mathbb{P}_{(X, Y) \sim \mathcal{D}} \left[\arg \max_i f_i^M(X + \delta_{f_i}^*(X)) = Y \right] \geq r_{f^M}, \end{aligned}$$

Maximize mixed classifier clean accuracy while maintaining robust accuracy

- Consider the approximate problem

$$\begin{aligned} & \min_{M \in \mathcal{M}, \alpha \in [1/2, 1]} \mathbb{P}_{X \sim \mathcal{X}_{ic}} \left[m_{h^M}(X) \geq \frac{1-\alpha}{\alpha} \right] \\ & \text{s. t. } \mathbb{P}_{Z \sim \mathcal{X}_{ca}} \left[\underline{m}_{h^M}^*(Z) \geq \frac{1-\alpha}{\alpha} \right] \geq \beta, \end{aligned} \quad (3)$$

Minimize $h^M(\cdot)$'s confidence margin at mispredicted clean data while maintaining $h^M(\cdot)$'s margin at correctly predicted worst-case adversarial data

where \mathcal{X}_{ic} is the distribution formed by clean examples incorrectly classified by $h^M(\cdot)$, \mathcal{X}_{ca} is the distribution formed by attacked examples correctly classified by $h^M(\cdot)$, X, Z are the random variables drawn from these distributions, and $\beta \in [0, 1]$ controls the desired level of robust accuracy with respect to the robust accuracy of $h(\cdot)$.

- The approximate problem decouples the optimization from $g(\cdot)$.

Quality of Approximation

- Original goal:

$$\begin{aligned} & \max_{M \in \mathcal{M}, \alpha \in [1/2, 1]} \mathbb{P}_{(X, Y) \sim \mathcal{D}} \left[\arg \max_i f_i^M(X) = Y \right] \quad (2) \\ \text{s. t. } & \mathbb{P}_{(X, Y) \sim \mathcal{D}} \left[\arg \max_i f_i^M(X + \delta_{f^M}^*(X)) = Y \right] \geq r_{f^M}, \end{aligned}$$

- Approximate problem:

$$\begin{aligned} & \min_{M \in \mathcal{M}, \alpha \in [1/2, 1]} \mathbb{P}_{X \sim \mathcal{X}_{ic}} \left[m_{h^M}(X) \geq \frac{1-\alpha}{\alpha} \right] \quad (3) \\ \text{s. t. } & \mathbb{P}_{Z \sim \mathcal{X}_{ca}} \left[\underline{m}_{h^M}^*(Z) \geq \frac{1-\alpha}{\alpha} \right] \geq \beta, \end{aligned}$$

- The objectives are equivalent, (3)'s constraint is more conservative

Assumption 4.1. On unattacked clean data, if $h^M(\cdot)$ makes a correct prediction, then $g(\cdot)$ is also correct.

Assumption 4.2. The transformation $M(\cdot)$ does not change the predicted class due to, e.g., monotonicity. Namely, it holds that $\arg \max_i M(h(x))_i = \arg \max_i h_i(x)$ for all x .

Theorem 4.3. Suppose that Assumption 4.2 holds. Let r_h denote the robust accuracy of $h(\cdot)$. If $\beta \geq r_{f^M}/r_h$, then a solution to (3) is feasible for (2).

Theorem 4.4. Suppose that Assumption 4.1 holds. Furthermore, consider an input random variable X and suppose that the margin of $h^M(X)$ is independent of whether $g(X)$ is correct. Then, minimizing the objective of (3) is equivalent to maximizing the objective of (2).

Nonlinear Transformation Parameterization

- **Step 1: Layer Norm (LN)**
 - Nonlinear transformations' effect depends on the logits range.
 - LN unifies the range.
 - For each image x , we standardize the logits $h(x)$ to have zero mean and variance one.
- **Step 2: Clamp**
 - We use a ReLU-like function to clamp the logits smaller than a positive threshold toward zero.
 - Introduce the threshold parameter c .
 - Since correct predictions have greater margins, clamping enlarges the margin difference between correct and incorrect examples.
 - We select GELU based on ablation studies.

So far, $h^M(x) = \text{GELU}(\text{LN}(h(x)) + c)$

Nonlinear Transformation Parameterization

- **Step 3: Exponentiation**

- Amplify large logits (common in correct predictions) to further enlarge the margin difference.
- Use absolute value to preserve logit sign.
- Introduce the exponent parameter p .

- **Step 4: Temperature Scaling**

- Softmax “saturates” with large logits.
- Temperature scaling allows for adjusting the level of saturation.
- Introduce the scale parameter s .

Final formulation:

$$h^{\text{Clamp},c}(x) = \text{Clamp}(\text{LN}(h(x)) + c)$$
$$h^M_c(x) = s \cdot |h^{\text{Clamp},c}(x)|^p \cdot \text{sgn}(h^{\text{Clamp},c}(x))$$

Optimizing s, p, c, α

- The resulting problem is then

$$\begin{aligned} \min_{s, p, c, \alpha \in \mathbb{R}} \quad & \mathbb{P}_{X \sim \mathcal{X}_{ic}} \left[m_{h^{\text{map}, s, p, c}}(X) \geq \frac{1-\alpha}{\alpha} \right] \\ \text{s. t.} \quad & \mathbb{P}_{Z \sim \mathcal{X}_{ca}} \left[m_{h^{\text{map}, s, p, c}}^*(Z) \geq \frac{1-\alpha}{\alpha} \right] \geq \beta \\ & s \geq 0, \quad p \geq 0, \quad 1/2 \leq \alpha \leq 1. \end{aligned}$$

- $\beta = 0.985$ works well in practice.

- Only three degrees of freedom.
 - Because the robust accuracy constraint is always active.
- Algorithm: grid search over s, p, c and calculate α via the constraint.
- Approximation for efficiency:
 - Use $h(\cdot)$ as a surrogate for $h^M(\cdot)$ in margin calculations, so that grid search doesn't need to include attack.

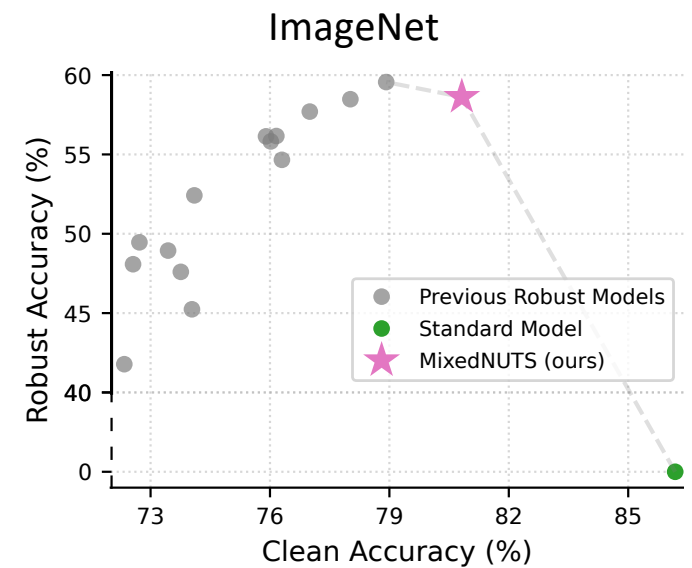
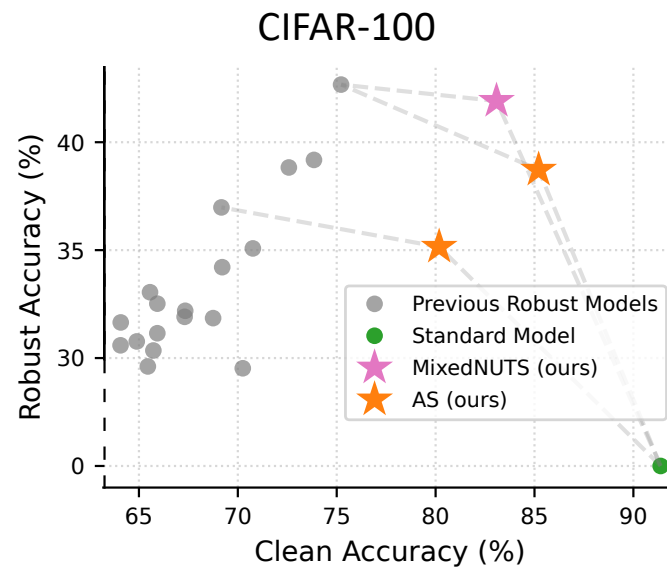
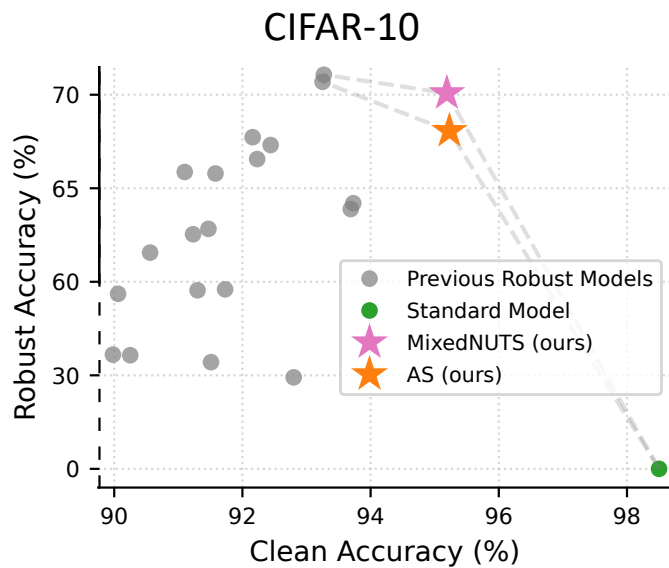
Optimizing s, p, c, α

Algorithm 1 Algorithm for optimizing s, p, c , and α .

- 1: Given an image set, save the predicted logits associated with mispredicted clean images $\{h^{\text{LN}}(x) : x \in \tilde{\mathcal{X}}_{ic}\}$.
- 2: Run MMAA on $h^{\text{LN}}(\cdot)$ and save the logits of correctly classified perturbed inputs $\{h^{\text{LN}}(x) : x \in \tilde{\mathcal{A}}_{ca}\}$.
- 3: Initialize candidate values $s_1, \dots, s_l, p_1, \dots, p_m, c_1, \dots, c_n$.
- 4: **for** s_i for $i = 1, \dots, l$ **do**
- 5: **for** p_j for $j = 1, \dots, m$ **do**
- 6: **for** c_k for $k = 1, \dots, n$ **do**
- 7: Obtain mapped logits $\{h^{M_{\tilde{c}_k^i}}(x) : x \in \tilde{\mathcal{A}}_{ca}\}$.
- 8: Calculate the margins from the mapped logits $\{m_{h^{M_{\tilde{c}_k^i}}}(x) : x \in \tilde{\mathcal{A}}_{ca}\}$.
- 9: Store the bottom $1 - \beta$ -quantile of the margins as $q_{1-\beta}^{ijk}$ (corresponds to $\frac{1-\alpha}{\alpha}$ in (6)).
- 10: Record the current objective $o^{ijk} \leftarrow \mathbb{P}_{X \in \tilde{\mathcal{X}}_{ic}} [m_{h^{M_{\tilde{c}_k^i}}}(X) \geq q_{1-\beta}^{ijk}]$.
- 11: **end for**
- 12: **end for**
- 13: **end for**
- 14: Find optimal indices $(i^*, j^*, k^*) = \arg \min_{i,j,k} o^{ijk}$.
- 15: Recover optimal mixing weight $\alpha^* := 1/(1+q_{1-\beta}^{i^*j^*k^*})$.
- 16: **return** $s^* := s_{i^*}, p^* := p_{j^*}, c^* := c_{k^*}, \alpha^*$.

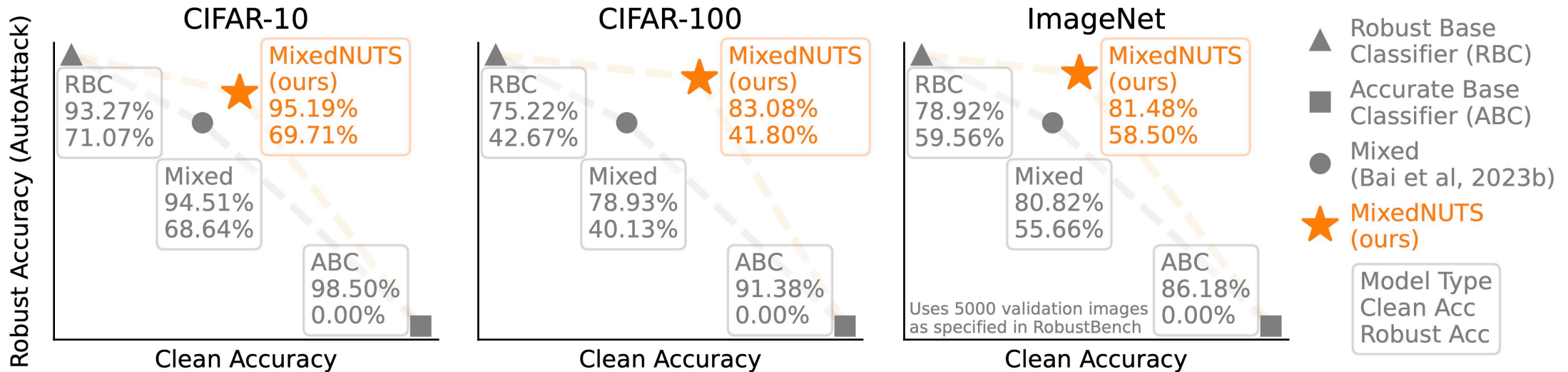
Main Experiment Result

- Mixed classifiers achieve state-of-the-art accuracy-robustness trade-off.



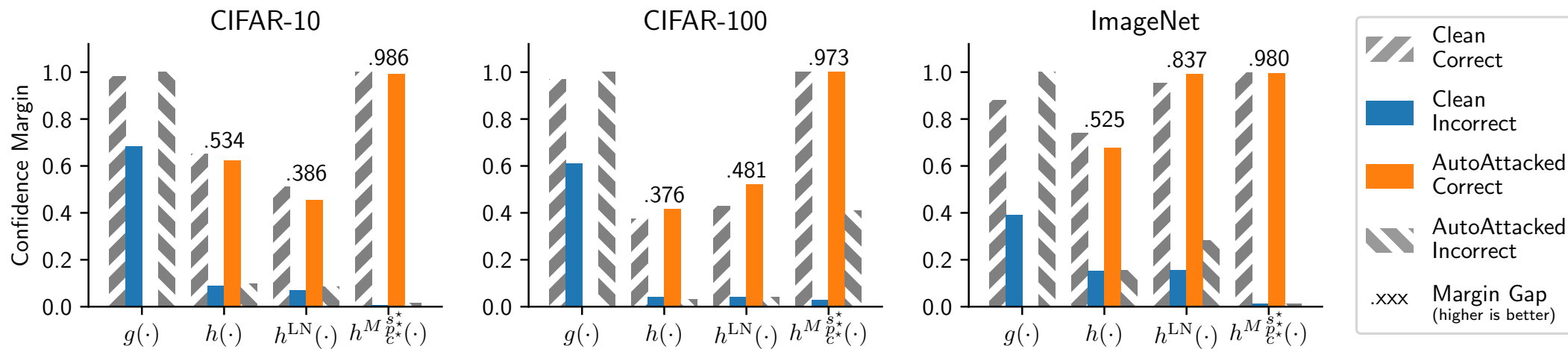
Main Experiment Result

- MixedNUTS' nonlinear logit transformations improve the accuracy-robustness trade-off.



Augmented Benign Margin Property

- MixedNUTS' nonlinear logit transformation augments the RBC's benign confidence margin properties.



Future – Beyond Adversarial Robustness

- Beyond adversarial robustness:
 - Generalized case: Model A specializes in Distribution A ; Model B specializes in Distribution B ; Distributions A, B share the same classes.
- Beyond classification:
 - Language models: output the probabilities of candidate next word tokens.
 - Existing models use mixtures of experts (MoE) to save computation (not all weights are activated).

Thank you!

Adaptive Smoothing: <https://arxiv.org/abs/2301.12554>

MixedNUTS: <https://arxiv.org/abs/2402.02263>

Presenter: Yatong Bai yatong_bai@berkeley.edu

May 19, 2024