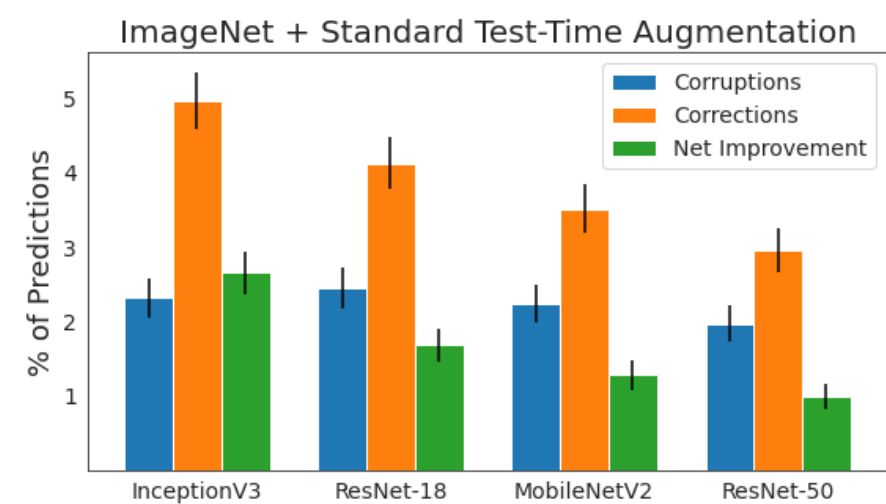# Better Aggregation in Test-Time Augmentation

Divya Shanmugam, Davis Blalock, Guha Balakrishnan, John Guttag

## Motivation

TTA introduces many incorrect predictions.
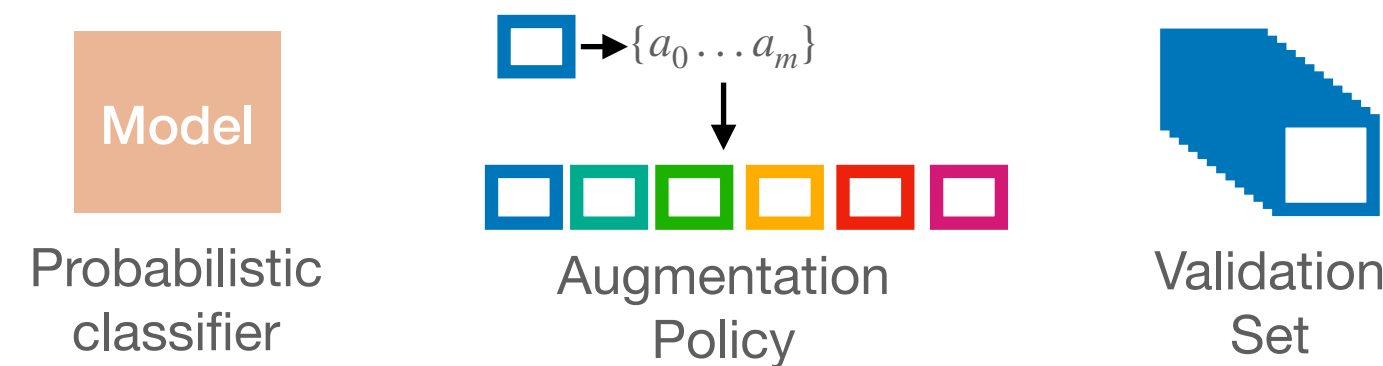
We aim to **characterize the errors introduced by TTA** and **develop a method to address these shortcomings.**

ImageNet + Standard Test-Time Augmentation

(Legend: Corruptions, Corrections, Net Improvement)

Bar chart: % of Predictions for InceptionV3, ResNet-18, MobileNetV2, ResNet-50

## Method

Key idea: learn augmentation specific weights to aggregate predictions.

Using:

Model — Probabilistic classifier

$\rightarrow \{a_0 \dots a_m\}$ — Augmentation Policy

Validation Set

Learn:

$$\begin{bmatrix} \theta_1 & \dots & \theta_M \end{bmatrix} \begin{bmatrix} a_{11} & \dots & a_{1C} \\ \vdots & \ddots & \\ a_{M1} & & a_{MC} \end{bmatrix}$$

AugTTA

$$\begin{bmatrix} \theta_{11} & \dots & \theta_{1C} \\ \vdots & \ddots & \\ \theta_{M1} & & \theta_{MC} \end{bmatrix} \begin{bmatrix} a_{11} & \dots & a_{1C} \\ \vdots & \ddots & \\ a_{M1} & & a_{MC} \end{bmatrix}$$
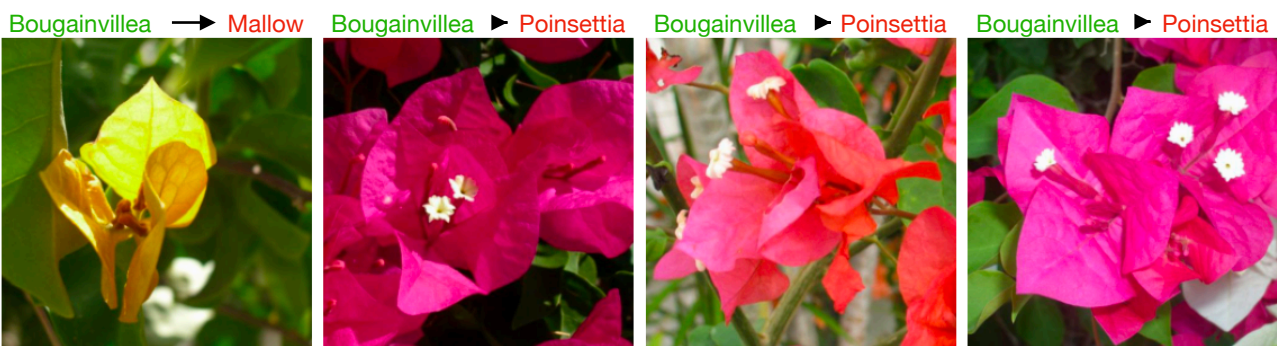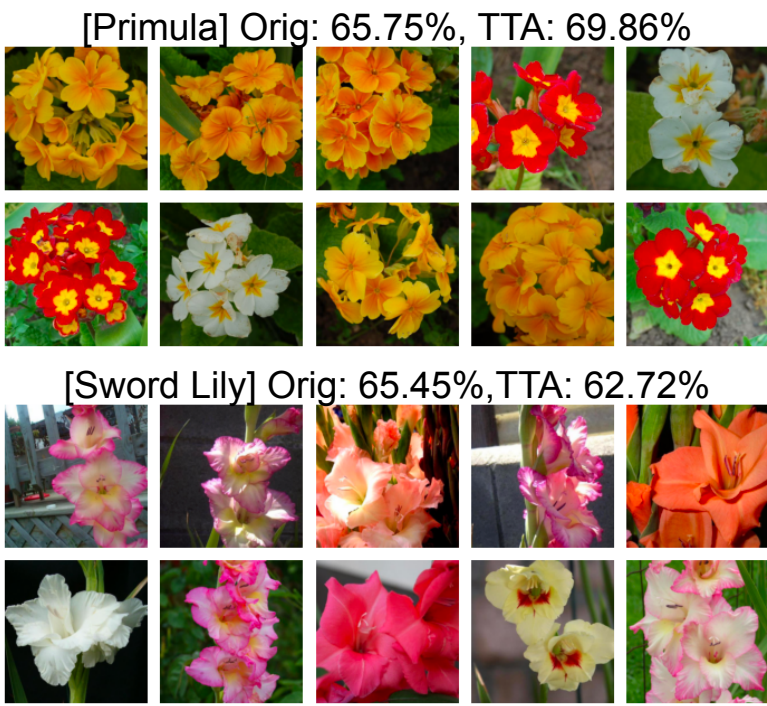
ClassTTA

## Takeaway

We present a new TTA method that uses an **augmentation-specific approach to aggregation** and provides improvements in classification accuracy.

## Analysis

1. Standard TTA changes predictions for classes with smaller distinguishing features, and classes that vary in scale.

Bougainvillea → Mallow | Bougainvillea ▶ Poinsettia | Bougainvillea ▶ Poinsettia | Bougainvillea ▶ Poinsettia

2. TTA harms classification accuracy for classes that exhibit higher variation in the training data.

3. The value of TTA is significantly correlated with the number of examples per class.

[Primula] Orig: 65.75%, TTA: 69.86%

[Sword Lily] Orig: 65.45%, TTA: 62.72%

## Results

1. Our method outperforms others across four datasets, four architectures, and two test-time augmentation policies.

**Standard TTA Policy.**

| Dataset | Model | Original | Max | Mean | GPS | Ours |
|---|---|---|---|---|---|---|
| Flowers102 | MobileNetV2 | $90.28 \pm 0.10$ | $90.17 \pm 0.25$ | $90.47 \pm 0.20$ | $88.28 \pm 0.17$ | $\mathbf{92.62 \pm 0.10}$ |
| Flowers102 | InceptionV3 | $89.28 \pm 0.08$ | $89.59 \pm 0.15$ | $90.07 \pm 0.22$ | $89.93 \pm 0.16$ | $\mathbf{91.16 \pm 0.21}$ |
| Flowers102 | ResNet-18 | $89.78 \pm 0.17$ | $89.47 \pm 0.11$ | $90.21 \pm 0.23$ | $90.01 \pm 0.22$ | $\mathbf{91.02 \pm 0.17}$ |
| Flowers102 | ResNet-50 | $\mathbf{91.72 \pm 0.18}$ | $91.61 \pm 0.08$ | $\mathbf{91.96 \pm 0.27}$ | $92.03 \pm 0.09$ | $92.02 \pm 0.16$ |
| ImageNet | MobileNetV2 | $71.38 \pm 0.06$ | $72.50 \pm 0.13$ | $\mathbf{72.69 \pm 0.06}$ | $72.50 \pm 0.11$ | $72.43 \pm 0.08$ |
| ImageNet | InceptionV3 | $69.66 \pm 0.12$ | $71.8 \pm 0.09$ | $72.45 \pm 0.13$ | $71.57 \pm 0.10$ | $\mathbf{72.79 \pm 0.02}$ |

**Expanded TTA Policy.**

| Dataset | Model | Original | Max | Mean | GPS | Ours |
|---|---|---|---|---|---|---|
| Flowers102 | MobileNetV2 | $90.94 \pm 0.16$ | $86.85 \pm 0.24$ | $91.14 \pm 0.08$ | $91.34 \pm 0.16$ | $\mathbf{92.49 \pm 0.20}$ |
| Flowers102 | InceptionV3 | $89.17 \pm 0.33$ | $87.89 \pm 0.20$ | $89.20 \pm 0.23$ | $89.43 \pm 0.16$ | $\mathbf{91.02 \pm 0.26}$ |
| Flowers102 | ResNet-18 | $89.20 \pm 0.16$ | $83.30 \pm 0.19$ | $89.47 \pm 0.09$ | $\mathbf{89.90 \pm 0.24}$ | $89.78 \pm 0.16$ |
| Flowers102 | ResNet-50 | $92.37 \pm 0.13$ | $89.39 \pm 0.19$ | $92.48 \pm 0.11$ | $92.57 \pm 0.21$ | $\mathbf{93.29 \pm 0.21}$ |
| ImageNet | MobileNetV2 | $71.18 \pm 0.05$ | $67.65 \pm 0.08$ | $71.84 \pm 0.12$ | $72.49 \pm 0.09$ | $\mathbf{72.57 \pm 0.09}$ |
| ImageNet | InceptionV3 | $69.51 \pm 0.08$ | $66.00 \pm 0.13$ | $70.85 \pm 0.11$ | $\mathbf{71.05 \pm 0.08}$ | $71.02 \pm 0.06$ |

2. Learned weights confirm qualitative results, and demonstrate higher variance for classes that exhibit higher variation in the training data.

**Black-Eyed Susan**
Low Variance in Aug. Weights

**Columbine**
High Variance in Aug. Weights

Read the paper for more test-time augmentation insights and instructions to reproduce experiments!