# Better Aggregation in Test-Time Augmentation

**Divya Shanmugam, Davis Blalock, Guha Balakrishnan, John Guttag**

# TTA is the aggregation of predictions across transformations of an image.

Traditionally:



↓

Model

↓

"bakery"  ✗

# TTA is the aggregation of predictions across transformations of an image.

Traditionally:



↓

Model

↓

"bakery"  ❌

With TTA:



↓

Model

↓

{"bakery", "bakery", "sandwich", "sandwich", "sandwich"}

↓

"sandwich"  ✅

# TTA produces more accurate and robust predictions than the original model *without retraining*



{"bakery", "bakery", "sandwich", "sandwich", "sandwich"}

"sandwich" ✅

# TTA produces more accurate and robust predictions than the original model *without retraining*

Two choices:
1. Selecting augmentations
2. Aggregating the resulting predictions



Model

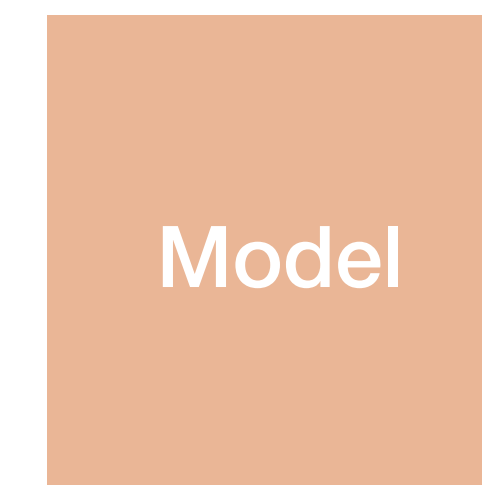{"bakery", "bakery", "sandwich", "sandwich", "sandwich"}

"sandwich" ✅

# TTA produces more accurate and robust predictions than the original model *without retraining*

Two choices:
1. Selecting augmentations
2. Aggregating the resulting predictions

Common augmentations include **flips, crops, and scales**, and predictions are typically aggregated via a **simple average.**



Model

{"bakery", "bakery", "sandwich", "sandwich", "sandwich"}

"sandwich" ✅

# TTA is widely applied.



Publications (total)

# TTA is widely applied.



Parallel Structure Deep Neural Network Using CNN and RNN with an Attention Mechanism for Breast Cancer Histology Image Classification.

Hongdou Yao, Xuejie Zhang, Xiaobing Zhou, Shengyan Liu

Publications (total)

# TTA is widely applied.



Parallel Structure Deep Neural Network Using CNN and RNN with an Attention Mechanism for Breast Cancer Histology Image Classification.

Hongdou Yao, Xuejie Zhang, Xiaobing Zhou, Shengyan Liu

Improving convolutional neural networks performance for image classification using test time augmentation: a case study using MURA dataset

Ibrahem Kandel, Mauro Castelli

2021, Health Information Science and Systems - Article

Publications (total)

# TTA is widely applied.



Parallel Structure Deep Neural Network Using CNN and RNN with an Attention Mechanism for Breast Cancer Histology Image Classification.

Hongdou Yao, Xuejie Zhang, Xiaobing Zhou, Shengyan Liu

Improving convolutional neural networks performance for image classification using test time augmentation: a case study using MURA dataset

Ibrahem Kandel, Mauro Castelli

2021, Health Information Science and Systems - Article

Test-time augmentation for deep learning-based cell segmentation on microscopy images

Nikita Moshkov, Botond Mathe, Attila Kertesz-Farkas, Reka Hollandi, Peter Horvath

Publications (total)

# TTA is widely applied.



Parallel Structure Deep Neural Network Using CNN and RNN with an Attention Mechanism for Breast Cancer Histology Image Classification.
Hongdou Yao, Xuejie Zhang, Xiaobing Zhou, Shengyan Liu

Improving convolutional neural networks performance for image classification using test time augmentation: a case study using MURA dataset
Ibrahem Kandel, Mauro Castelli
2021, Health Information Science and Systems - Article

Test-time augmentation for deep learning-based cell segmentation on microscopy images
Nikita Moshkov, Botond Mathe, Attila Kertesz-Farkas, Reka Hollandi, Peter Horvath

Robustness of convolutional neural networks in recognition of pigmented skin lesions
Roman C Maron, Sarah Haggenmüller, Christof von Kalle, Jochen S Utikal, Friedegund Meier, Frank F Gell...
2021, European Journal of Cancer - Article

Publications (total)

11

# Standard approaches to TTA work consistently improve network performance.



ImageNet + Standard Test-Time Augmentation

# Standard approaches to TTA change many predictions from correct to incorrect.

# Our plan



Characterize the errors introduced by TTA.



Present a new TTA method that addresses these shortcomings.

# Datasets we considered:

ImageNet: 1000 classes, 1.2 million images

# Datasets we considered:
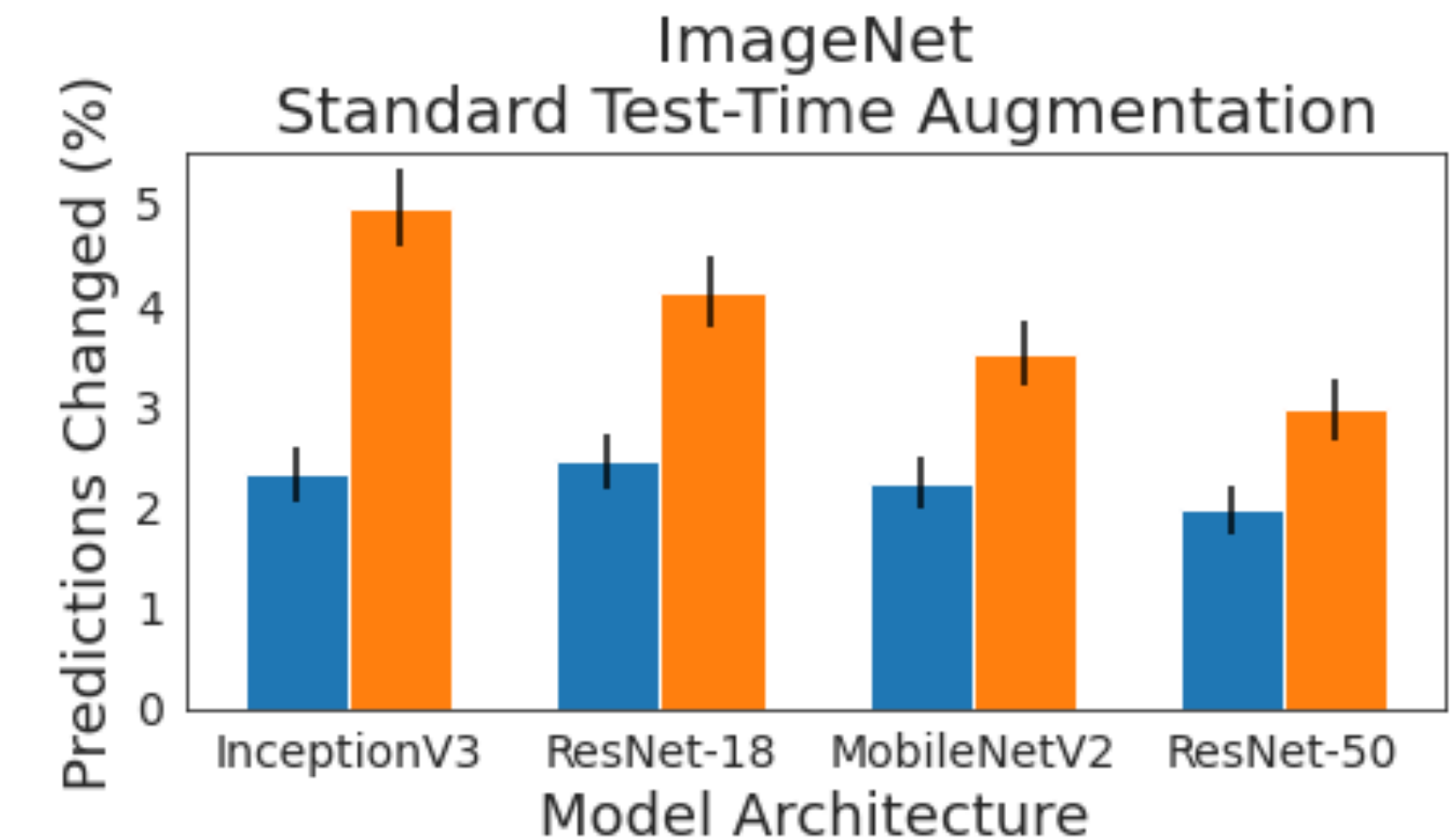
ImageNet: 1000 classes, 1.2 million images
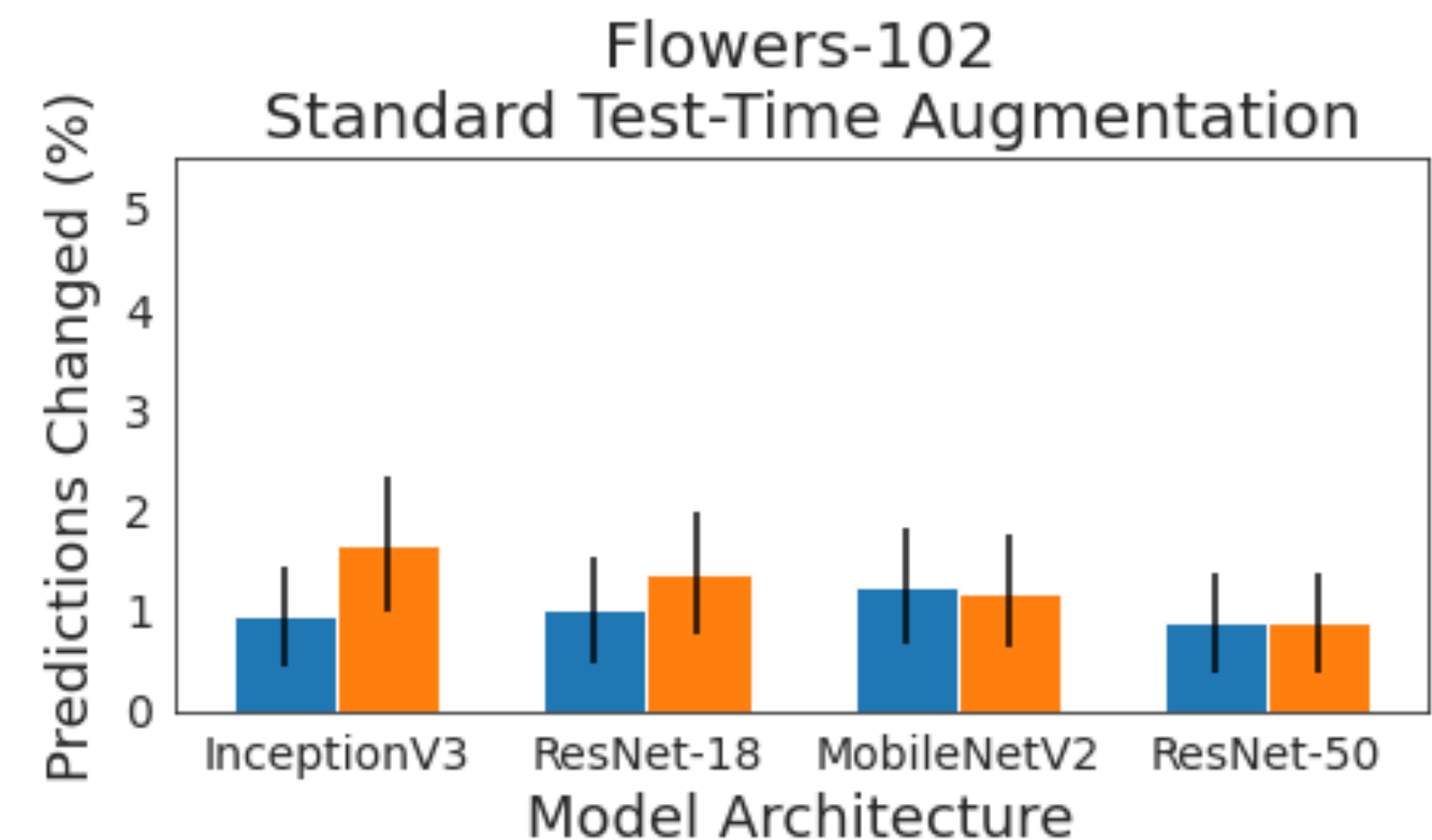


Flowers-102: 102 classes, 1020 images
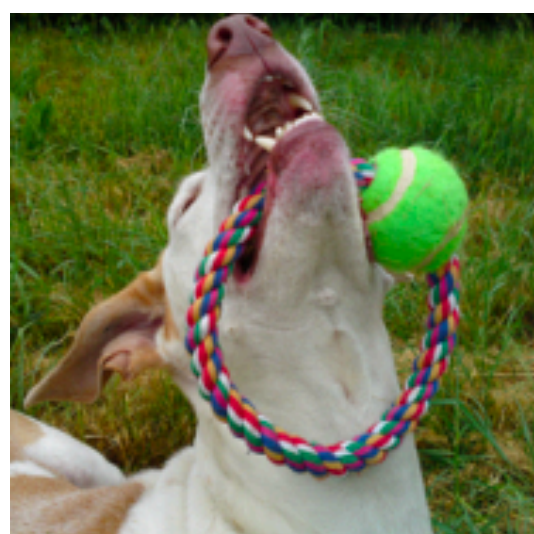
# Datasets we considered:
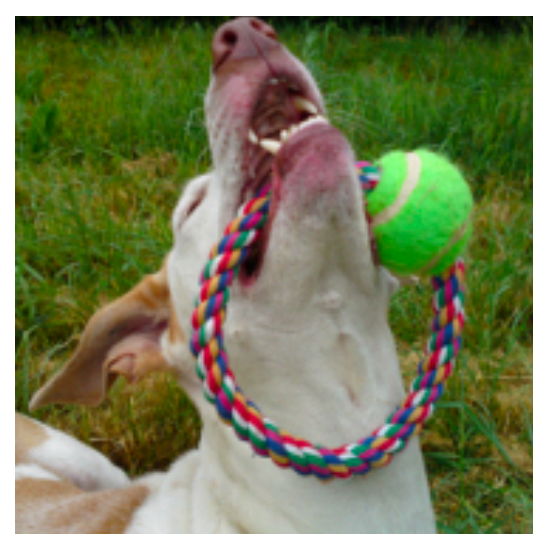
ImageNet: 1000 classes, 1.2 million images



Flowers-102: 102 classes, 1020 images
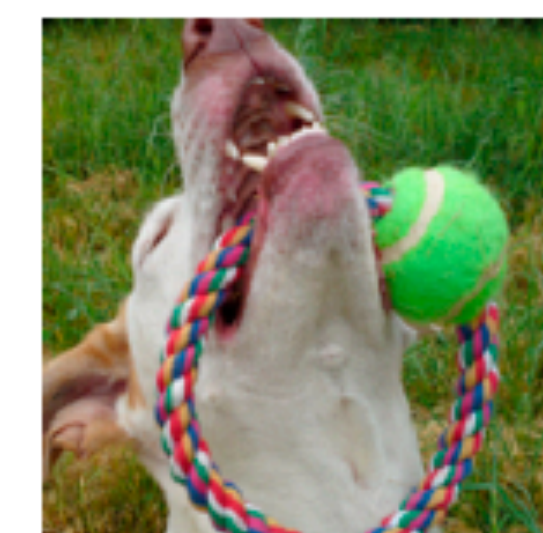
# Understanding why corruptions occur



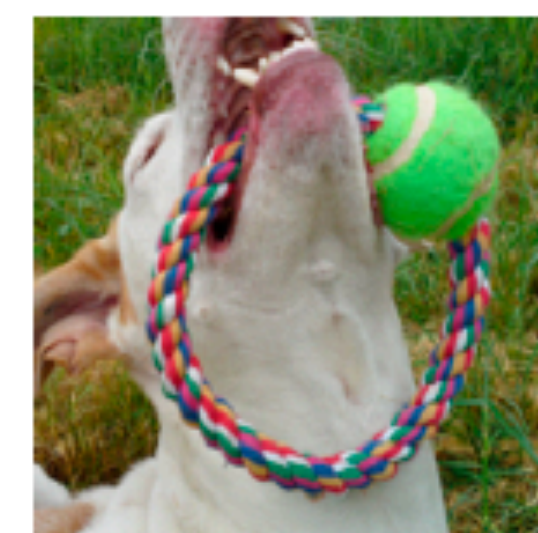True Label:
Ibizan Hound

# Zooming in on images with **multiple classes** favors classes that appear smaller.



Test-Time Augmentations of Original Image
(Flips, Crops, and Scales)

True Label:
Ibizan Hound

TTA Label:
Tennis Ball

# TTA can also benefit classes differently because of **class-dependent variation.**



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

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

Class-specific and dataset-specific attributes can affect the performance of traditional TTA.

# Key idea: Learn augmentation-specific weights for aggregating predictions.

# Key idea: Learn augmentation-specific weights for aggregating predictions.

We assume three inputs:

# Key idea: Learn augmentation-specific weights for aggregating predictions.

We assume three inputs:

**1**

Black box classifier

$p(y)$

Model

# Key idea: Learn augmentation-specific weights for aggregating predictions.

We assume three inputs:

**1**

Black box classifier

$p(y)$

**2**

Augmentation policy

$\{a_0, a_1, a_2, a_3 \ldots a_m\}$

# Key idea: Learn augmentation-specific weights for aggregating predictions.

We assume three inputs:

**1** Black box classifier

**2** Augmentation policy

**3** Labeled set of images

$$\{a_0, a_1, a_2, a_3 \ldots a_m\}$$

Model

*p(y)*

# Key idea: Learn augmentation-specific weights for aggregating predictions.

Two models:

1) Learn a weight parameter for each augmentation
2) Learn a weight parameter for each augmentation-class pair

# Key idea: Learn augmentation-specific weights for aggregating predictions.

Two models:

1) Learn a weight parameter for each augmentation
2) Learn a weight parameter for each augmentation-class pair

AugTTA

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

Weights for each augmentation

# Key idea: Learn augmentation-specific weights for aggregating predictions.

Two models:

1) Learn a weight parameter for each augmentation
2) Learn a weight parameter for each augmentation-class pair

AugTTA

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

Weights for each augmentation

ClassTTA

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

Weights for each class-augmentation pair

# Our method in three steps:

**1**

Split the labeled data into data for training the models, and data for deciding which model to use.

# Our method in three steps:

**1**

Split the labeled data into data for training the models, and data for deciding which model to use.

**2**

Learn the parameters for AugTTA and ClassTTA using projected gradient descent to ensure learned weights are non-negative.

# Our method in three steps:

Split the labeled data into data for training the models, and data for deciding which model to use.

Learn the parameters for AugTTA and ClassTTA using projected gradient descent to ensure learned weights are non-negative.

Choose AugTTA or ClassTTA based on performance on the held-out data.

# Our method produces higher Top-1 classification accuracy than existing work.

**Standard TTA Policy.**

| Dataset | Model | Original | Max | Mean | GPS | Ours |
|---------|-------|----------|-----|------|-----|------|
| Flowers102 | MobileNetV2 | $90.28 \pm 0.10$ | | | | |

# Our method produces higher Top-1 classification accuracy than existing work.

**Standard TTA Policy.**

| Dataset | Model | Original | Max | Mean | GPS | Ours |
|---------|-------|----------|-----|------|-----|------|
| Flowers102 | MobileNetV2 | $90.28 \pm 0.10$ | $90.17 \pm 0.25$ | | | |

# Our method produces higher Top-1 classification accuracy than existing work.

**Standard TTA Policy.**

| Dataset | Model | Original | Max | Mean | GPS | Ours |
|---|---|---|---|---|---|---|
| Flowers102 | MobileNetV2 | $90.28 \pm 0.10$ | $90.17 \pm 0.25$ | $90.47 \pm 0.20$ | | |

# Our method produces higher Top-1 classification accuracy than existing work.

**Standard TTA Policy.**

| Dataset | Model | Original | Max | Mean | GPS | Ours |
|---|---|---|---|---|---|---|
| Flowers102 | MobileNetV2 | $90.28 \pm 0.10$ | $90.17 \pm 0.25$ | $90.47 \pm 0.20$ | $88.28 \pm 0.17$ | |

# Our method produces higher Top-1 classification accuracy than existing work.

**Standard TTA Policy.**

| Dataset | Model | Original | Max | Mean | GPS | Ours |
|---------|-------|----------|-----|------|-----|------|
| Flowers102 | MobileNetV2 | $90.28 \pm 0.10$ | $90.17 \pm 0.25$ | $90.47 \pm 0.20$ | $88.28 \pm 0.17$ | $\mathbf{92.62 \pm 0.10}$ |

# Our method produces higher Top-1 classification accuracy than existing work.

**Standard TTA Policy.**

| Dataset | Model | Original | Max | Mean | GPS | Ours |
|---|---|---|---|---|---|---|
| Flowers102 | MobileNetV2 | $90.28 \pm 0.10$ | $90.17 \pm 0.25$ | $90.47 \pm 0.20$ | $88.28 \pm 0.17$ | $\textbf{92.62} \pm \textbf{0.10}$ |
| Flowers102 | InceptionV3 | $89.28 \pm 0.08$ | $89.59 \pm 0.15$ | $90.07 \pm 0.22$ | $89.93 \pm 0.16$ | $\textbf{91.16} \pm \textbf{0.21}$ |
| Flowers102 | ResNet-18 | $89.78 \pm 0.17$ | $89.47 \pm 0.11$ | $90.21 \pm 0.23$ | $90.01 \pm 0.22$ | $\textbf{91.02} \pm \textbf{0.17}$ |
| Flowers102 | ResNet-50 | $\textbf{91.72} \pm \textbf{0.18}$ | $91.61 \pm 0.08$ | $\textbf{91.96} \pm \textbf{0.27}$ | $\textbf{92.03} \pm \textbf{0.09}$ | $92.02 \pm 0.16$ |
| ImageNet | MobileNetV2 | $71.38 \pm 0.06$ | $72.50 \pm 0.13$ | $\textbf{72.69} \pm \textbf{0.06}$ | $72.50 \pm 0.11$ | $72.43 \pm 0.08$ |
| ImageNet | InceptionV3 | $69.66 \pm 0.12$ | $71.8 \pm 0.09$ | $72.45 \pm 0.13$ | $71.57 \pm 0.10$ | $\textbf{72.79} \pm \textbf{0.02}$ |
| ImageNet | ResNet-18 | $69.37 \pm 0.1$ | $70.26 \pm 0.13$ | $\textbf{71.02} \pm \textbf{0.13}$ | $70.8 \pm 0.1$ | $\textbf{71.06} \pm \textbf{0.10}$ |
| ImageNet | ResNet-50 | $75.78 \pm 0.08$ | $76.62 \pm 0.08$ | $\textbf{76.91} \pm \textbf{0.09}$ | $\textbf{76.73} \pm \textbf{0.11}$ | $\textbf{76.75} \pm \textbf{0.14}$ |
| CIFAR100 | CNN-7 | $74.15 \pm 0.18$ | $75.00 \pm 0.31$ | $75.48 \pm 0.11$ | $75.45 \pm 0.21$ | $\textbf{75.92} \pm \textbf{0.20}$ |
| STL10 | CNN-5 | $77.92 \pm 0.19$ | $77.76 \pm 0.22$ | $\textbf{78.58} \pm \textbf{0.25}$ | $\textbf{78.32} \pm \textbf{0.17}$ | $\textbf{78.52} \pm \textbf{0.31}$ |

# TTA + smaller networks can exceed original performance of larger networks.

**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 | $\mathbf{91.72 \pm 0.18}$ | $91.61 \pm 0.08$ | $\mathbf{91.96 \pm 0.27}$ | $\mathbf{92.03 \pm 0.09}$ | $92.02 \pm 0.16$ |
| 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 | $69.66 \pm 0.12$ | $71.8 \pm 0.09$ | $72.45 \pm 0.13$ | $71.57 \pm 0.10$ | $\mathbf{72.79 \pm 0.02}$ |
| ImageNet | ResNet-18 | $69.37 \pm 0.1$ | $70.26 \pm 0.13$ | $\mathbf{71.02 \pm 0.13}$ | $70.8 \pm 0.1$ | $\mathbf{71.06 \pm 0.10}$ |
| ImageNet | ResNet-50 | $75.78 \pm 0.08$ | $76.62 \pm 0.08$ | $\mathbf{76.91 \pm 0.09}$ | $\mathbf{76.73 \pm 0.11}$ | $\mathbf{76.75 \pm 0.14}$ |
| CIFAR100 | CNN-7 | $74.15 \pm 0.18$ | $75.00 \pm 0.31$ | $75.48 \pm 0.11$ | $75.45 \pm 0.21$ | $\mathbf{75.92 \pm 0.20}$ |
| STL10 | CNN-5 | $77.92 \pm 0.19$ | $77.76 \pm 0.22$ | $\mathbf{78.58 \pm 0.25}$ | $\mathbf{78.32 \pm 0.17}$ | $\mathbf{78.52 \pm 0.31}$ |

# TTA + smaller networks can exceed original performance of larger networks.

**Standard TTA Policy.**

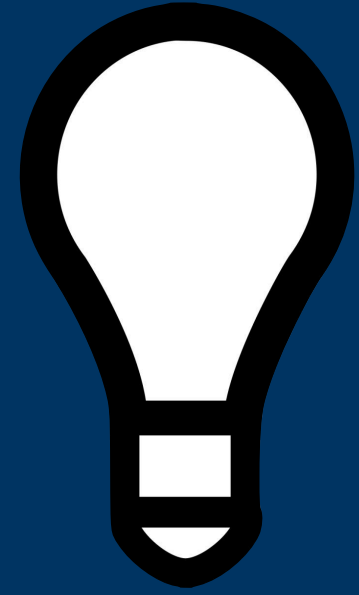| Dataset | Model | Original | Max | Mean | GPS | Ours |
|---------|-------|----------|-----|------|-----|------|
| Flowers102 | MobileNetV2 | $90.28 \pm 0.10$ | $90.17 \pm 0.25$ | $90.47 \pm 0.20$ | $88.28 \pm 0.17$ | $\mathbf{92.62 \pm 0.10}$ |
| Flowers102 | InceptionV3 | $89.28 \pm 0.08$ | $89.59 \pm 0.15$ | $90.07 \pm 0.22$ | $89.93 \pm 0.16$ | $\mathbf{91.16 \pm 0.21}$ |
| 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 | $\mathbf{91.72 \pm 0.18}$ | $91.61 \pm 0.08$ | $\mathbf{91.96 \pm 0.27}$ | $\mathbf{92.03 \pm 0.09}$ | $92.02 \pm 0.16$ |
| 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 | $69.66 \pm 0.12$ | $71.8 \pm 0.09$ | $72.45 \pm 0.13$ | $71.57 \pm 0.10$ | $\mathbf{72.79 \pm 0.02}$ |
| ImageNet | ResNet-18 | $69.37 \pm 0.1$ | $70.26 \pm 0.13$ | $\mathbf{71.02 \pm 0.13}$ | $70.8 \pm 0.1$ | $\mathbf{71.06 \pm 0.10}$ |
| ImageNet | ResNet-50 | $75.78 \pm 0.08$ | $76.62 \pm 0.08$ | $\mathbf{76.91 \pm 0.09}$ | $\mathbf{76.73 \pm 0.11}$ | $76.75 \pm 0.14$ |
| CIFAR100 | CNN-7 | $74.15 \pm 0.18$ | $75.00 \pm 0.31$ | $75.48 \pm 0.11$ | $75.45 \pm 0.21$ | $\mathbf{75.92 \pm 0.20}$ |
| STL10 | CNN-5 | $77.92 \pm 0.19$ | $77.76 \pm 0.22$ | $\mathbf{78.58 \pm 0.25}$ | $\mathbf{78.32 \pm 0.17}$ | $\mathbf{78.52 \pm 0.31}$ |

# The weights learned by ClassTTA reflect variation in the training data.

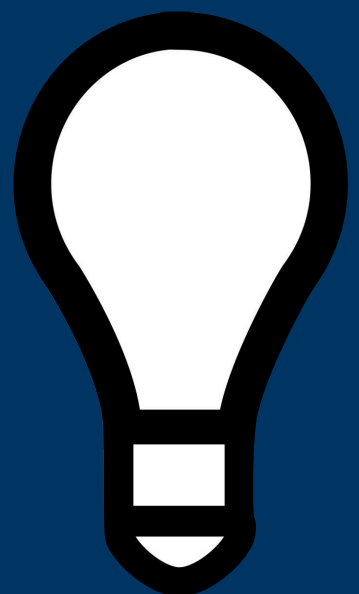Low Variance in Augmentation Weights

High Variance in Augmentation Weights

Our method improves classification accuracy and is **nearly free** in terms of model size, training time, and implementation burden.

The learned weights shed light on 1) dataset-specific and class-specific robustness to specific augmentations  and 2) which classes exhibit higher variation in the training data.

# In summary:

* Class-specific and dataset-specific attributes have systematic effects on the performance of common approaches to TTA.

* We share insights on when TTA is likely to be successful and which classes are negatively affected by the use of TTA.

* We develop a method that increases the classification accuracy of a pre-trained network.

## Visit our poster to learn more!

(or email me at divyas@mit.edu)