



# Angular Visual Hardness

Beidi Chen

Department of Computer Science, Rice University

Collaborators: Weiyang Liu, Animesh Garg, Zhiding Yu, Anshumali Shrivastava,  
Jan kautz, and Anima Anandkumar



Caltech



**PART**

**0**

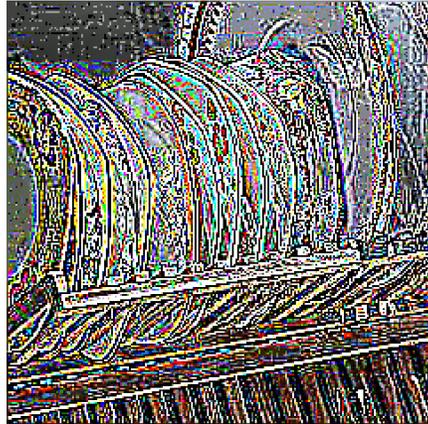
**Motivation**

# Human Visual Hardness

plate rack



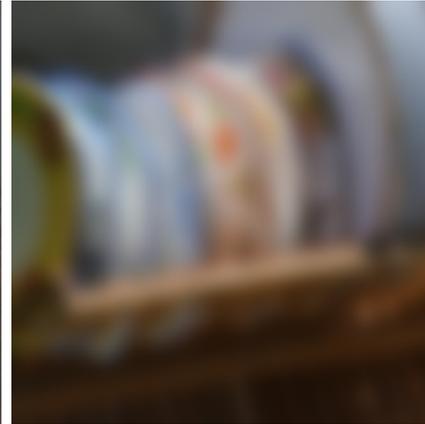
sharpness



contrast



blur

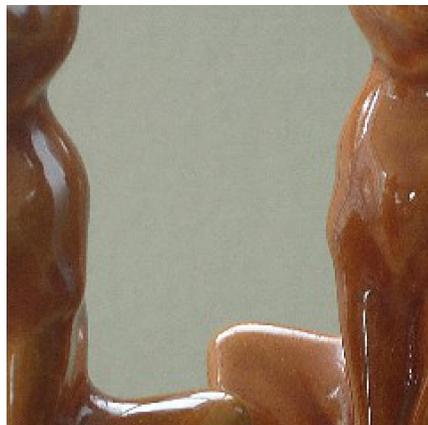


**Image Degradation**

dishwasher



saltshaker



nail



oil filter



**Semantic Ambiguity**

# Gap between human visual system and CNNs

**Hard** for Human and **Easy** for CNNs



Class Name

Nail

0.93

0.2

Human Selection Frequency

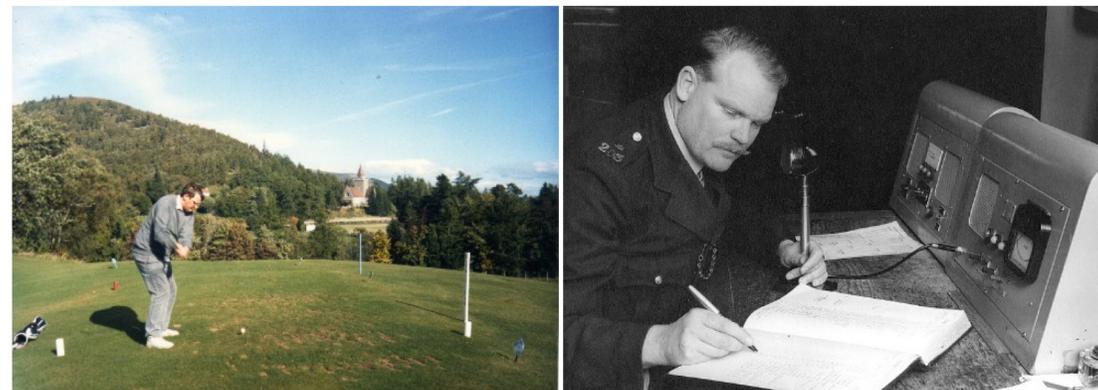
Softmax Score

Oil Filter

0.998

0.2

**Easy** for Human and **Hard** for CNNs



Golf Ball

0.001

1.0

Radio

0.001

1.0

# Agenda

Part 1

Background

Part 2

Discoveries

Part 3

Applications

Part 4

Conclusion

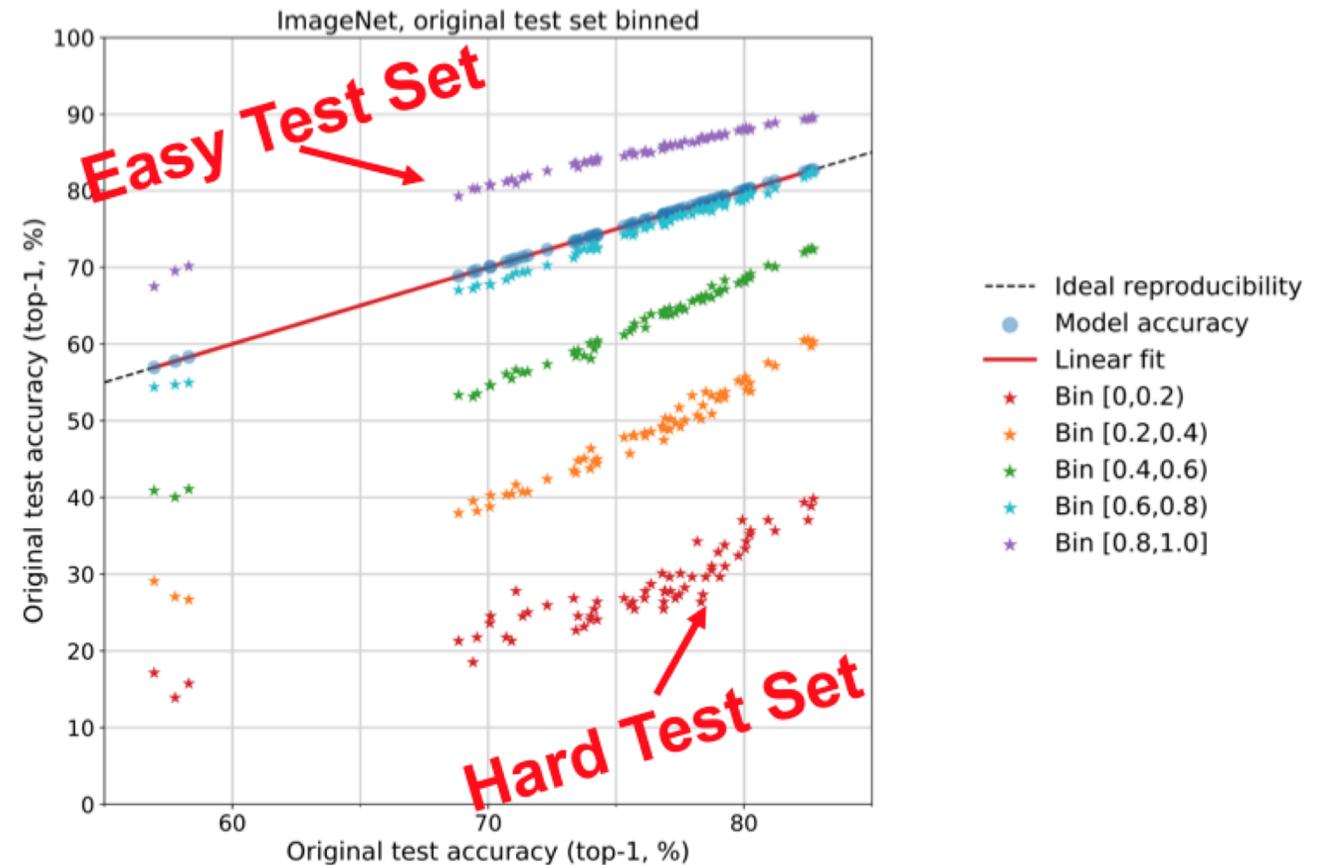
**PART**

**01**

**Background**

# Do ImageNet Classifiers Generalize to ImageNet?

## Human Labeling Interface



# Loss function of CNNs in visual recognition

Softmax cross-entropy loss

$$L = \frac{1}{N} \sum_i L_i = \frac{1}{N} \sum_i -\log \left( \frac{e^{f_{y_i}}}{\sum_j e^{f_j}} \right)$$

$$L_i = -\log \left( \frac{e^{\|\mathbf{W}_{y_i}\| \|\mathbf{x}_i\| \cos(\theta_{y_i})}}{\sum_j e^{\|\mathbf{W}_j\| \|\mathbf{x}_i\| \cos(\theta_j)}} \right)$$

The magnitude information

The angle between feature and classifier

Model Confidence

The diagram illustrates the decomposition of the softmax loss function. The numerator of the fraction inside the log is  $e^{\|\mathbf{W}_{y_i}\| \|\mathbf{x}_i\| \cos(\theta_{y_i})}$ . A blue box highlights  $\|\mathbf{W}_{y_i}\| \|\mathbf{x}_i\|$  and is labeled 'The magnitude information'. A red box highlights  $\cos(\theta_{y_i})$  and is labeled 'The angle between feature and classifier'. A pink box encloses the entire fraction, with a pink arrow pointing to the text 'Model Confidence' below it.

**PART**

**02**

**Discoveries**

# Proposal: Angular Visual Hardness (AVH)

Given a sample  $x$  with label  $y$ :

$$AVH(x) = \frac{\mathcal{A}(x, w_y)}{\sum_{i=1}^C \mathcal{A}(x, w_i)}$$

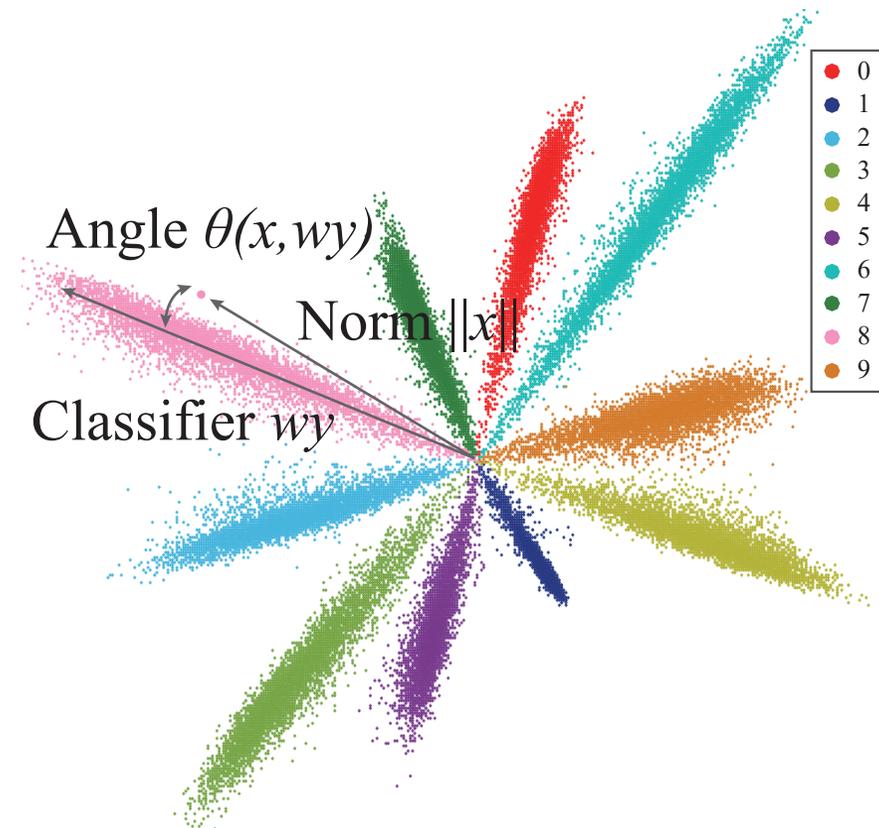
where,

$$\mathcal{A}(u, v) = \arccos\left(\frac{\langle u, v \rangle}{\|u\| \|v\|}\right)$$

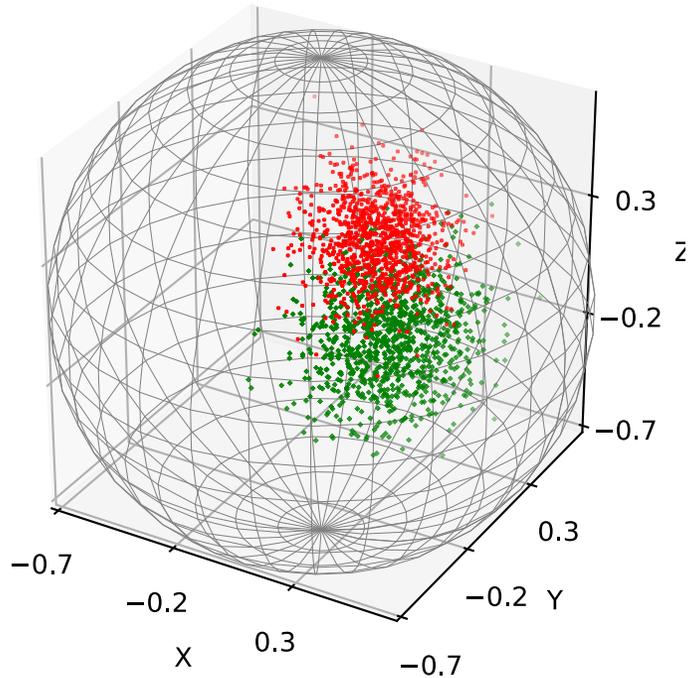
$w_i$  is the classifier for the  $i$ -th class.

## Theoretical Foundation:

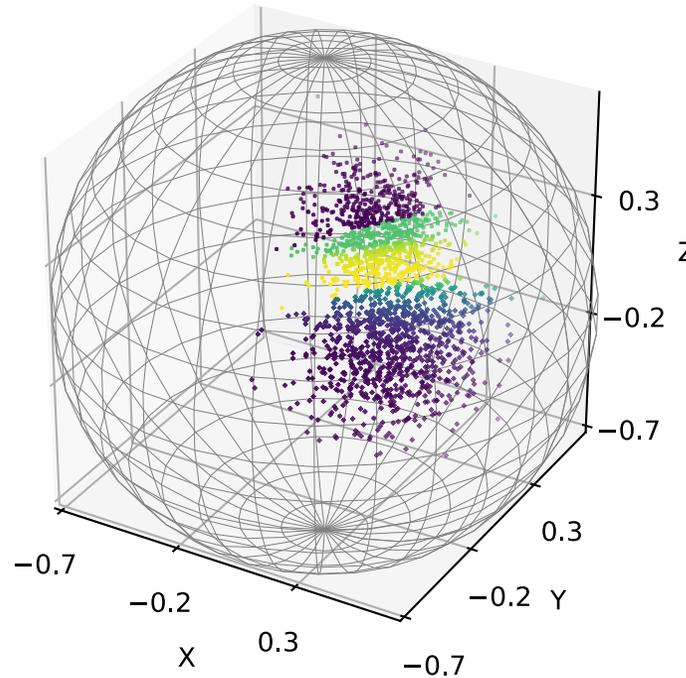
Soudry et al. "The Implicit Bias of Gradient Descent on Separable Data" ICLR 2018



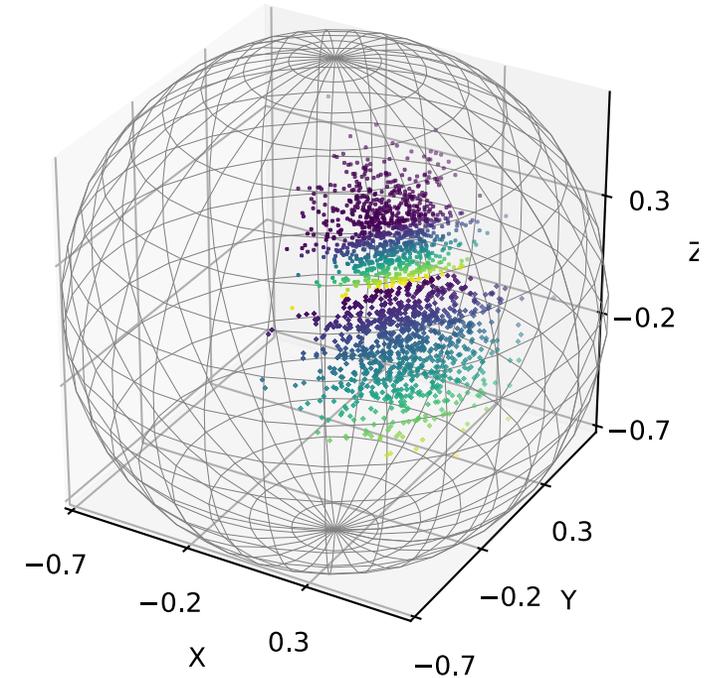
# Simple Example: AVH vs. $||\mathbf{x}||$



Raw data

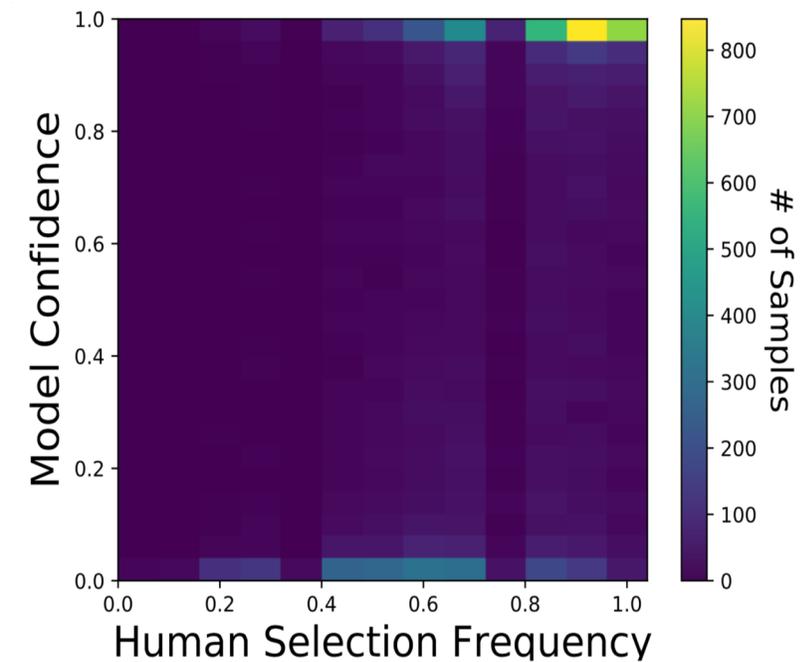
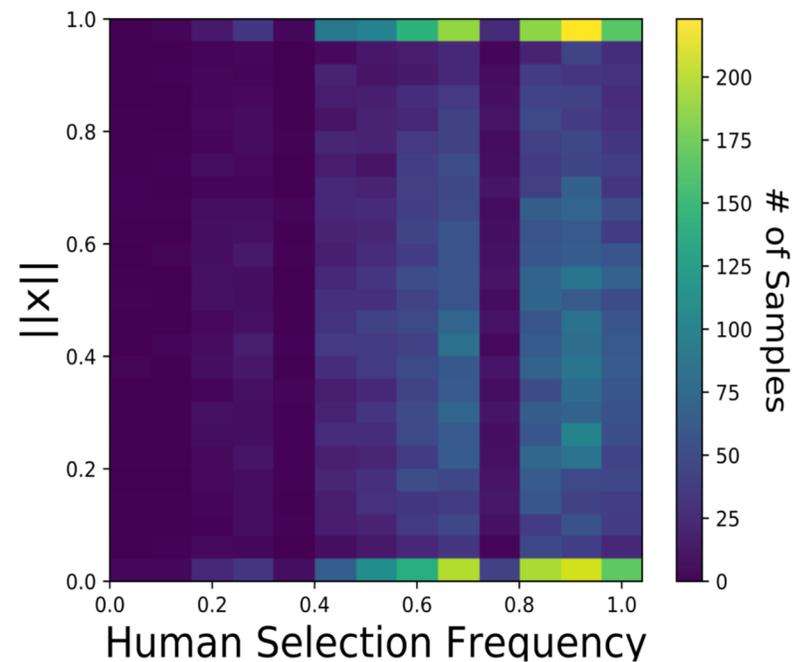
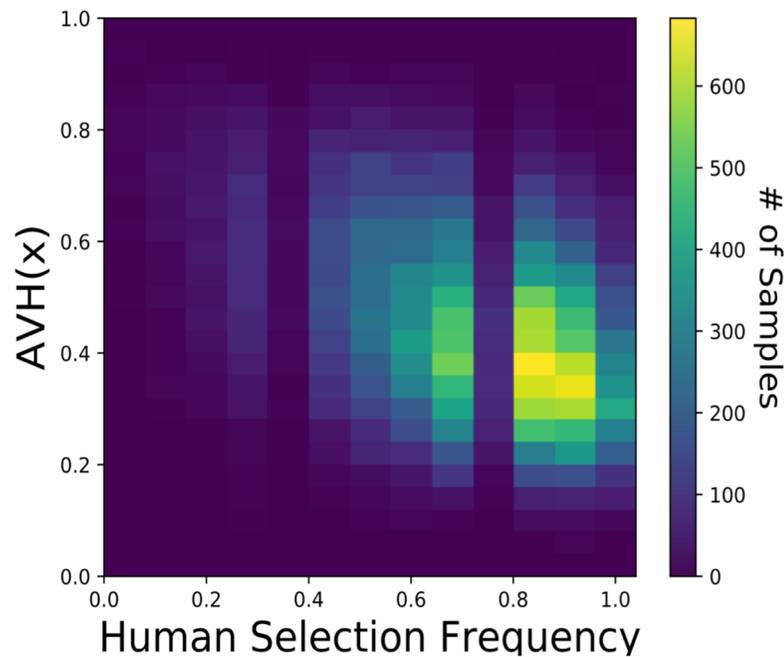


Color map of AVH

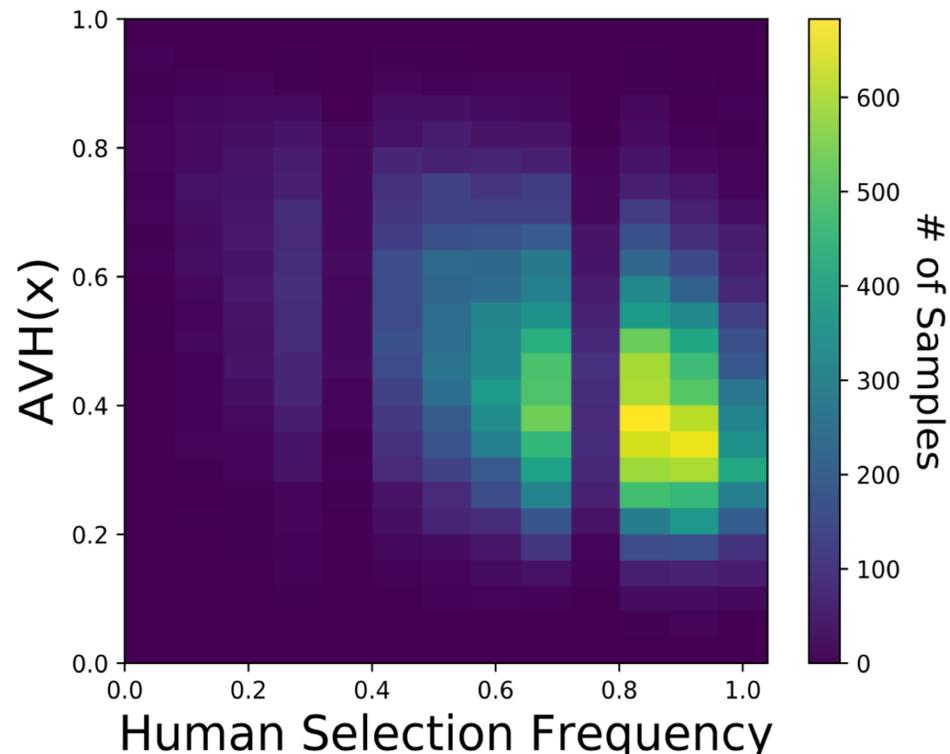


Color map of  $||\mathbf{x}||$

# CNN characteristics vs. human selection frequency



# AVH is well aligned with human frequency

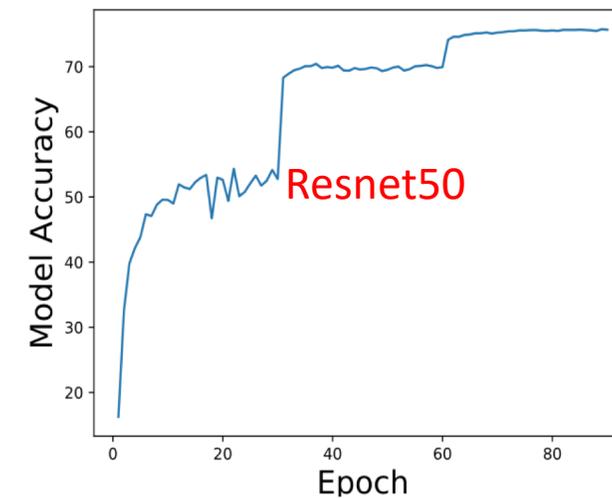
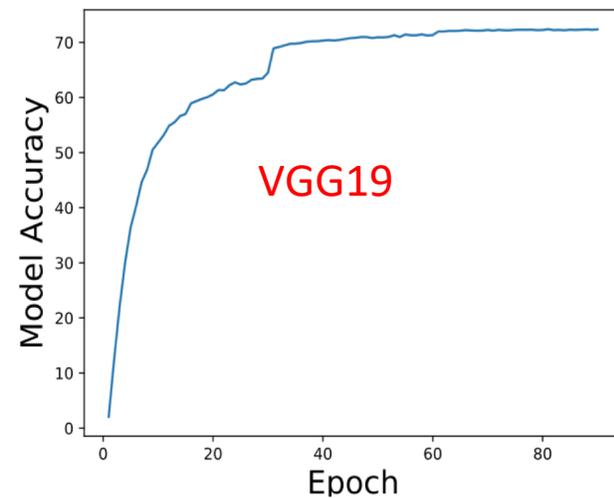
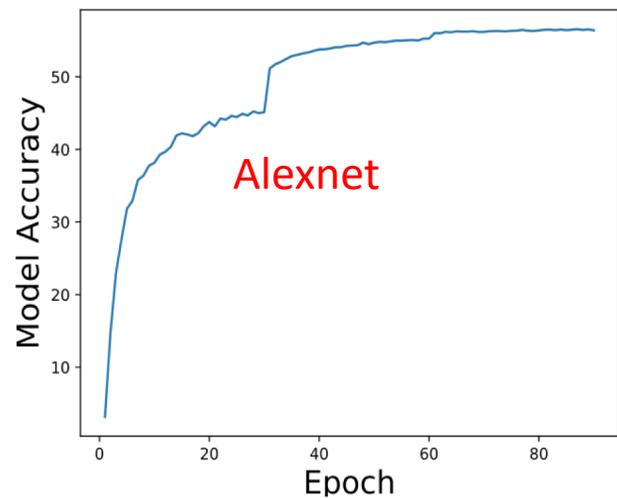
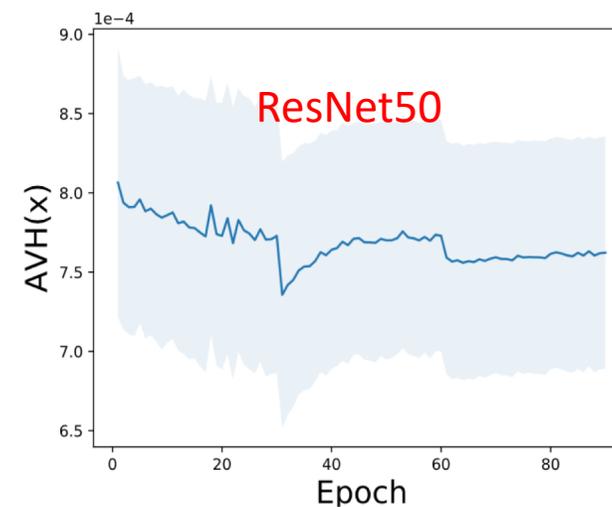
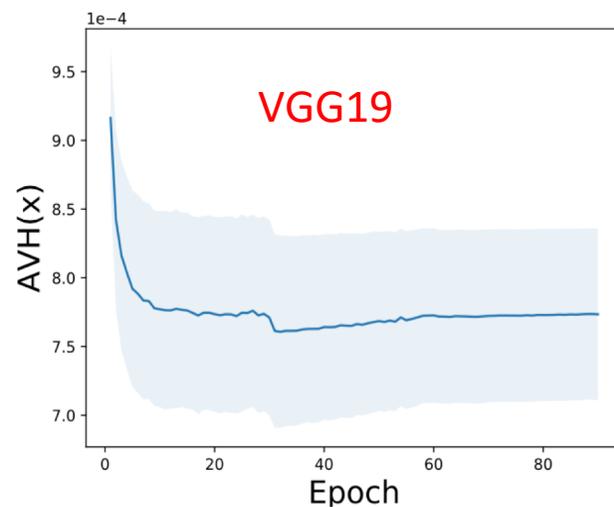
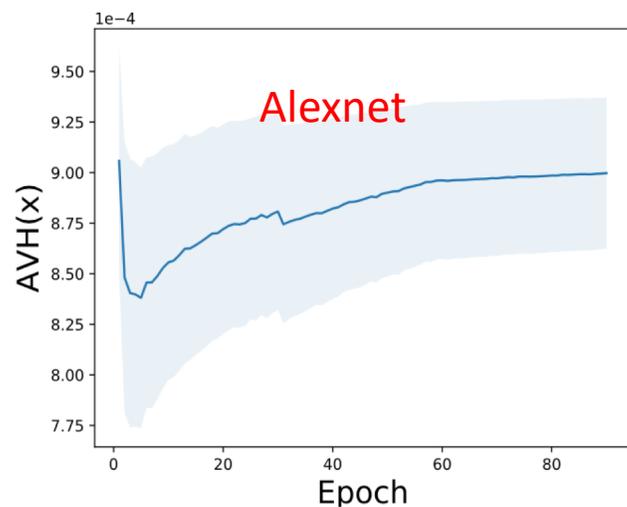


## Spearman rank correlations

	z-score	Total Coef	[0, 0.2]	[0.2, 0.4]	[0.4, 0.6]	[0.6, 0.8]	[0.8, 1.0]
Number of Samples	-	29987	837	2732	6541	11066	8811
AVH	0.377	0.36	0.228	0.125	0.124	0.103	0.094
Model Confidence	0.337	0.325	0.192	0.122	0.102	0.078	0.056
$\ \mathbf{x}\ _2$	-	0.0017	0.0013	0.0007	0.0005	0.0004	0.0003

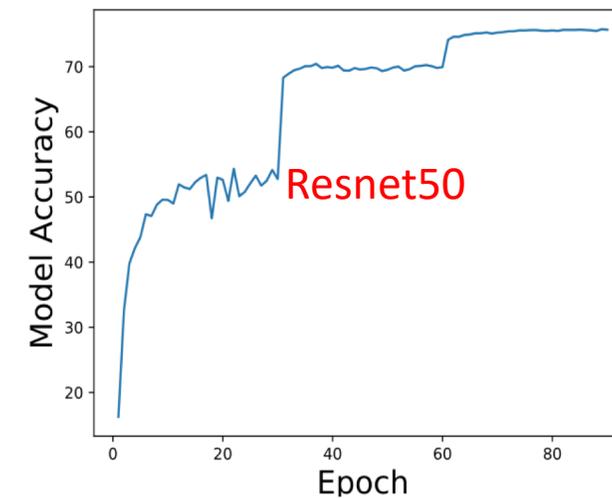
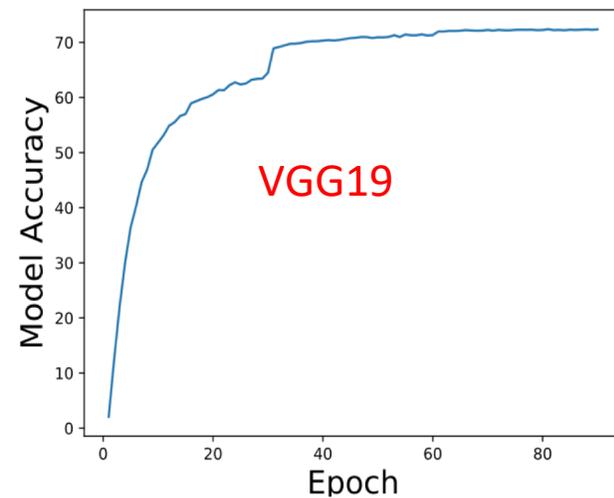
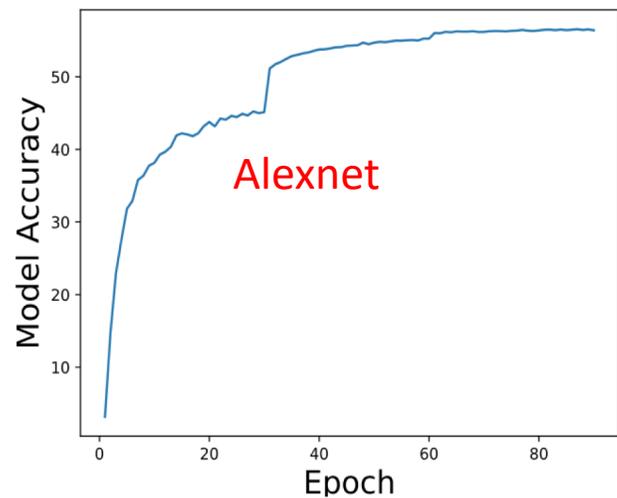
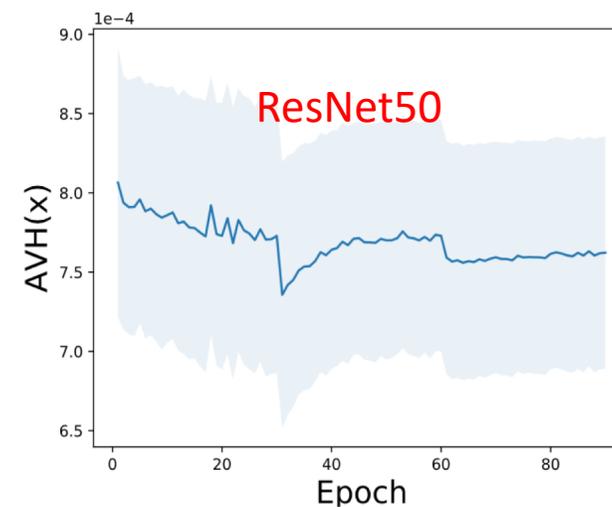
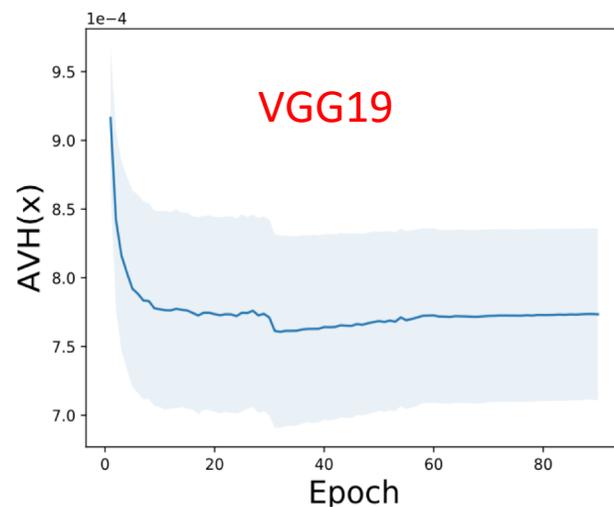
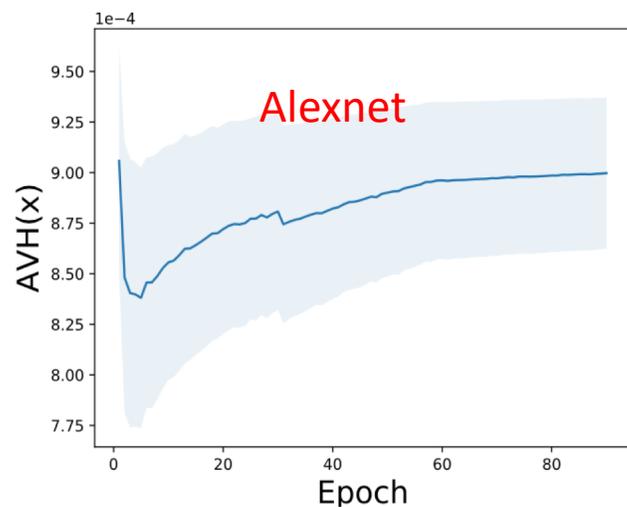
# Discovery 1

AVH hits a plateau very early even when the accuracy or loss is still improving



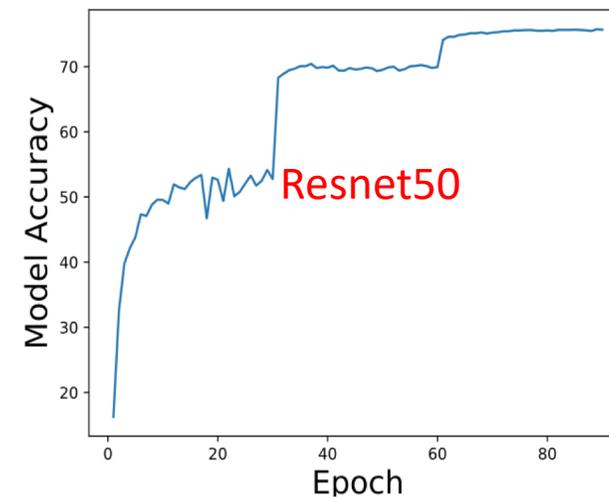
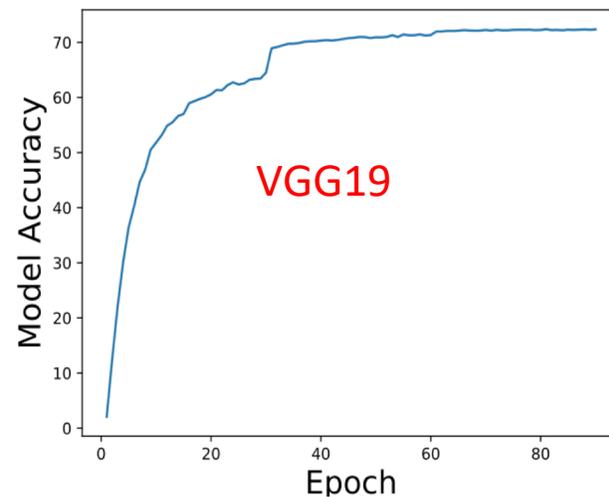
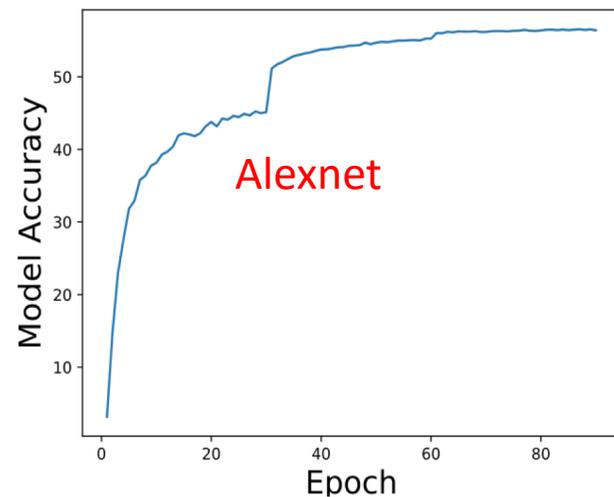
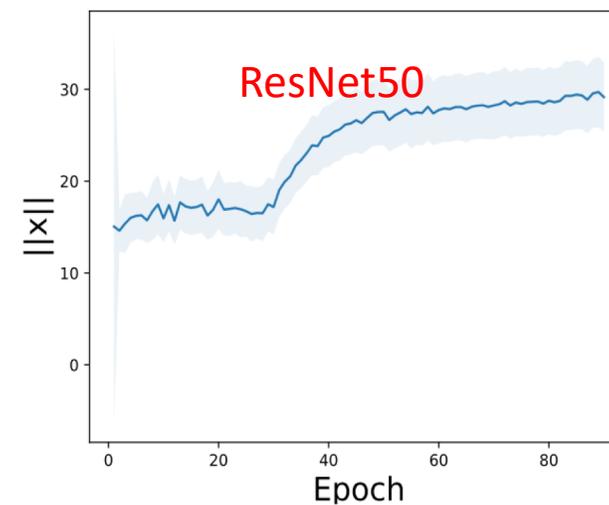
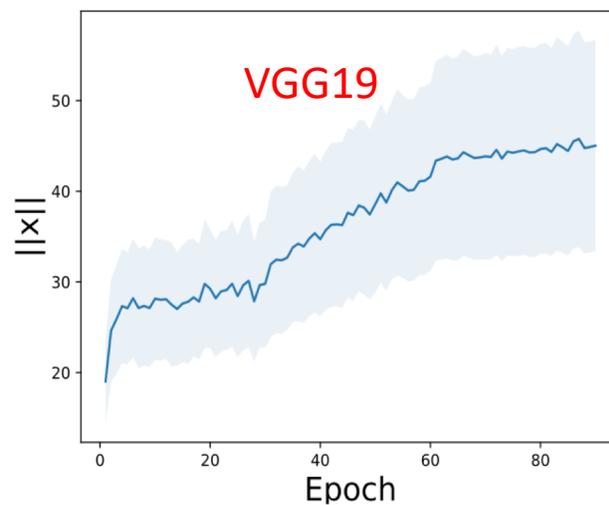
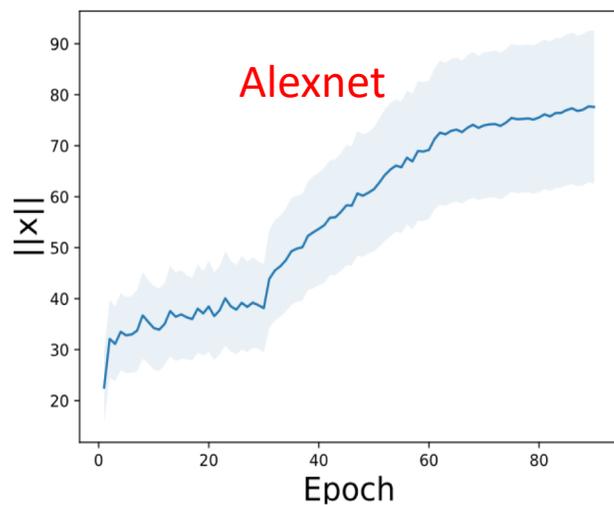
# Discovery 2

AVH is an indicator of model's generalization ability



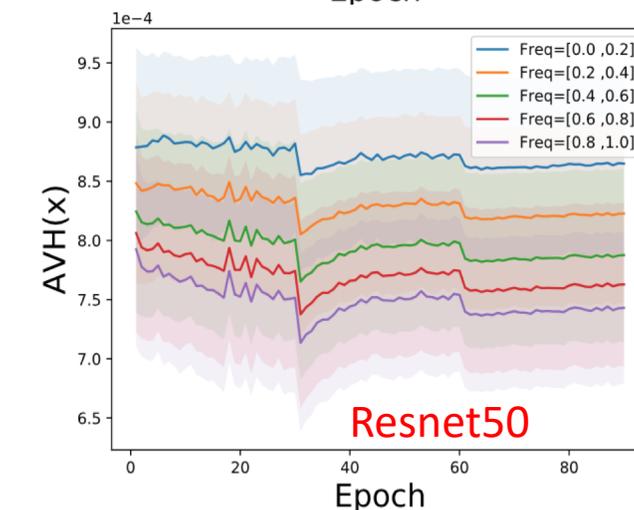
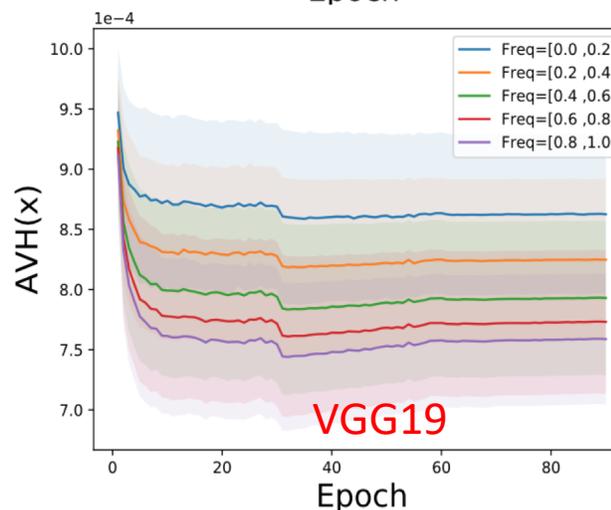
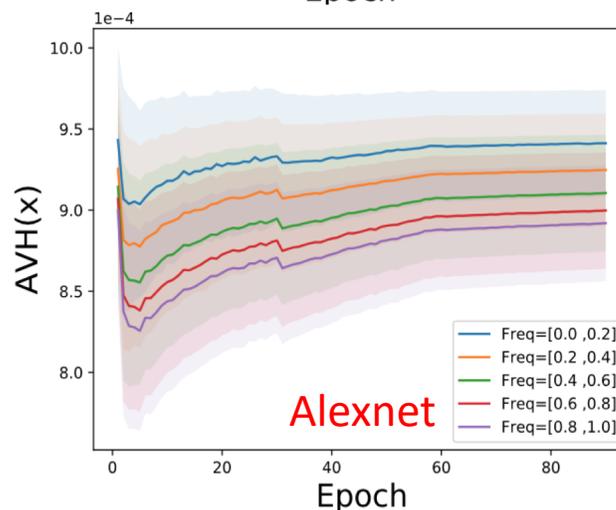
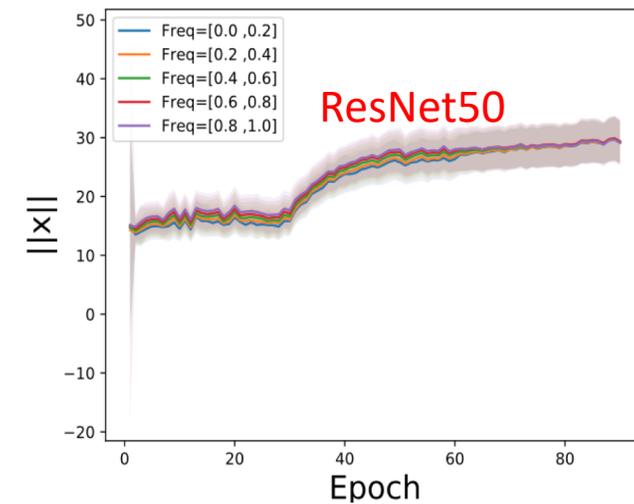
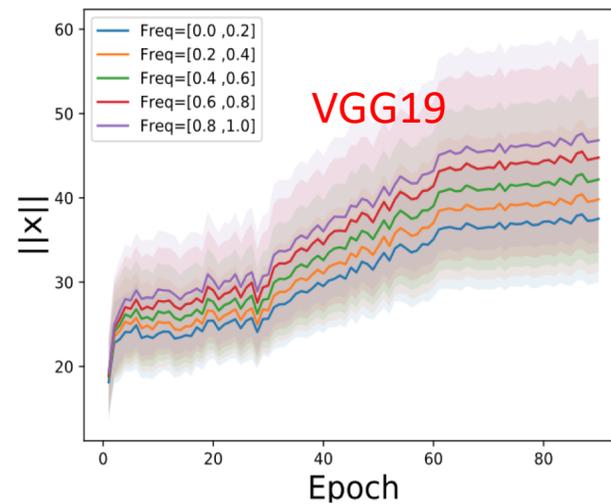
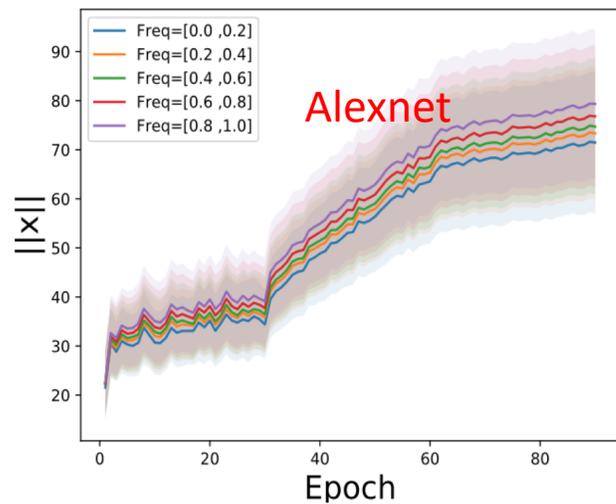
# Discovery 3

The norm of feature embeddings keeps increasing during training



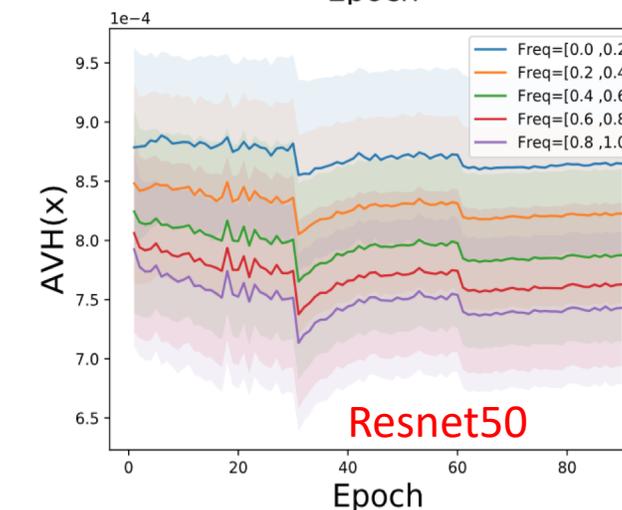
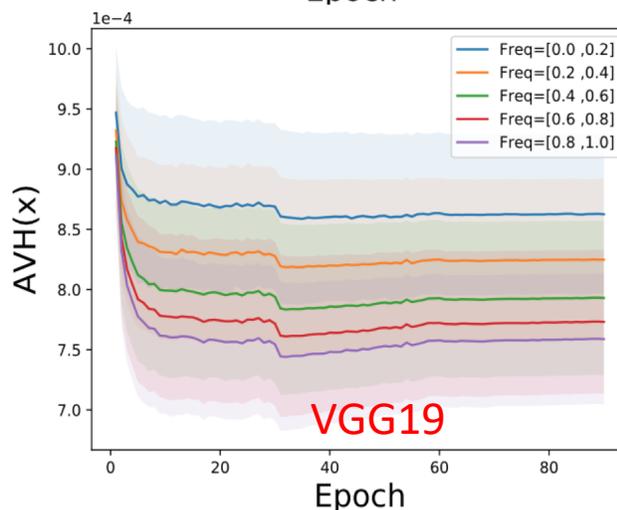
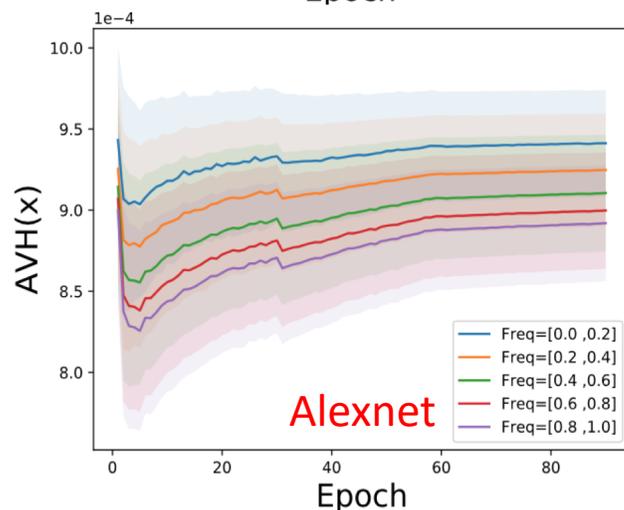
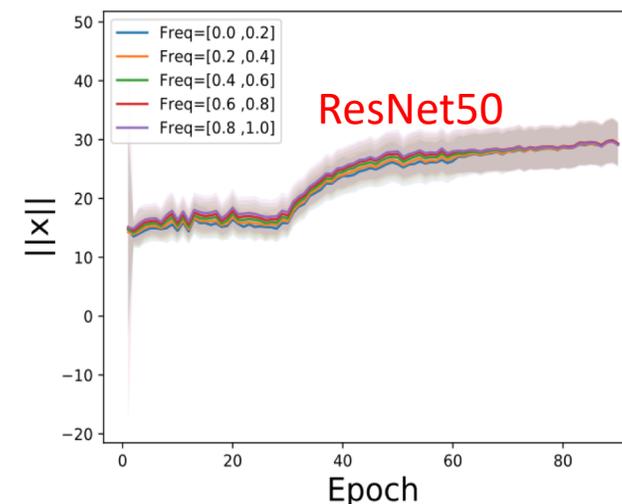
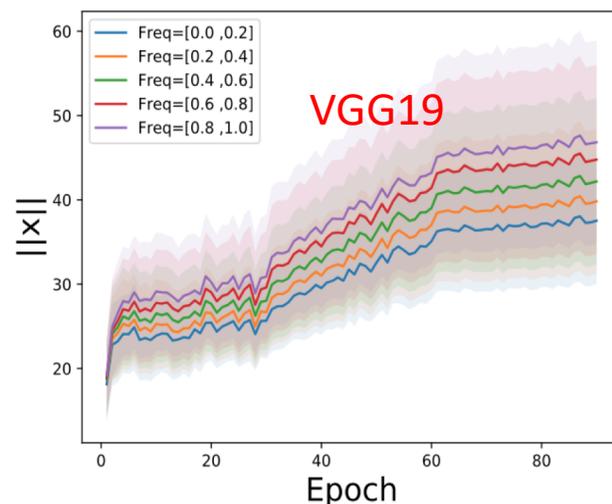
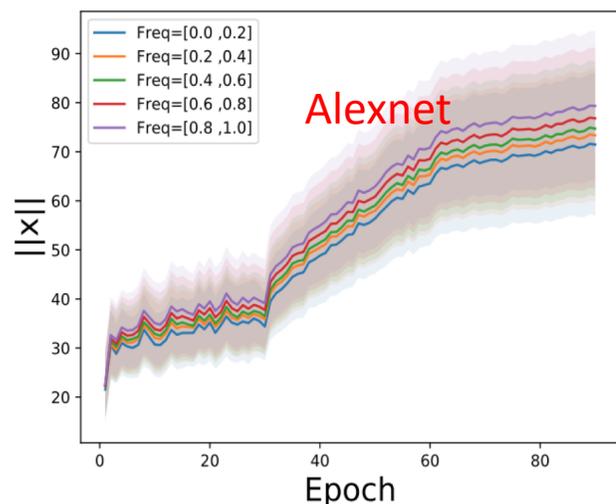
# Discovery 4

AVH's correlation with human selection frequency holds across models throughout training



# Discovery 5

The norm's correlation with human selection frequency is not consistent



# Conjecture on training dynamic of CNNs

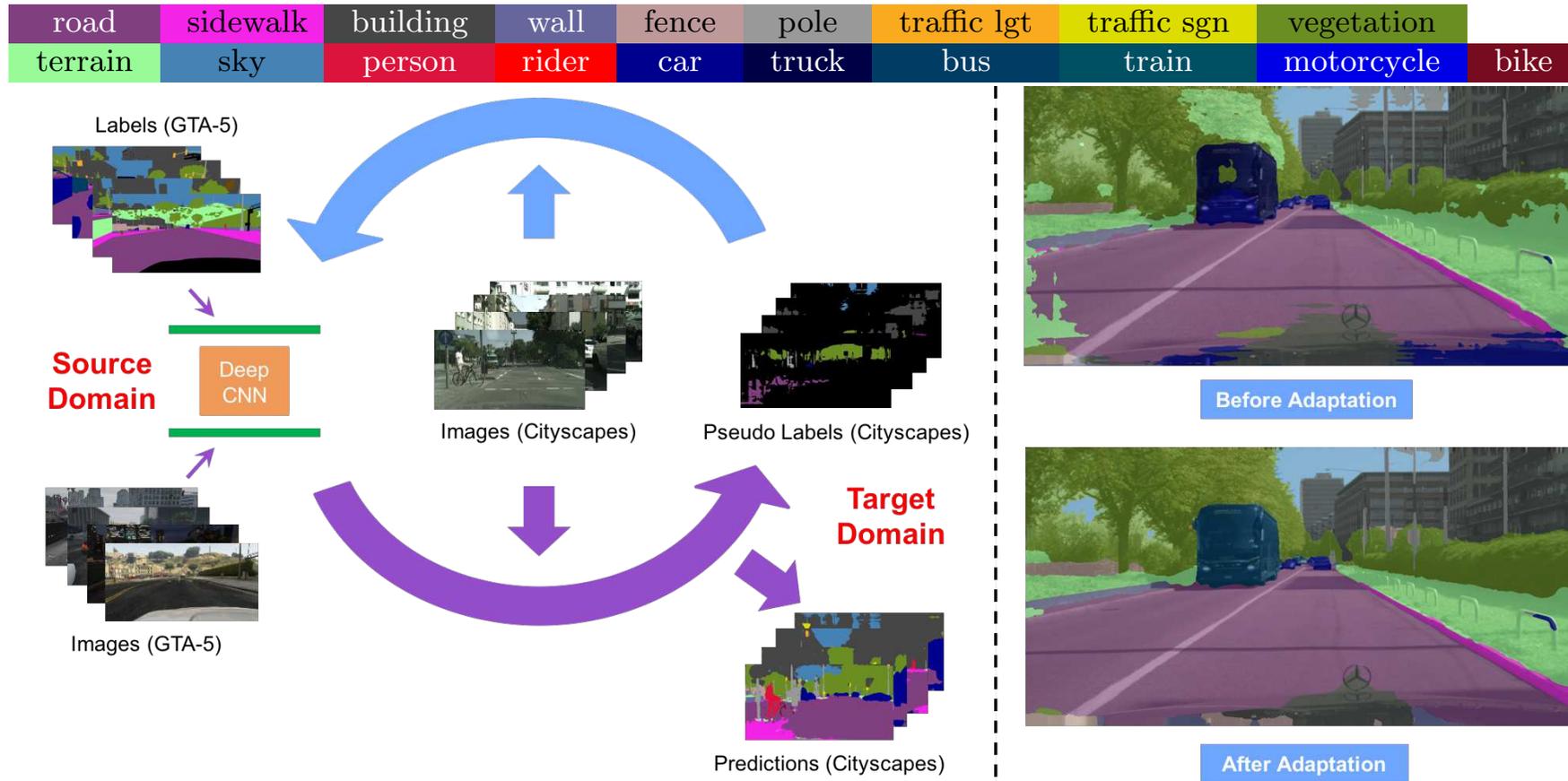
- Softmax cross-entropy loss will first optimize the angles among different classes while the norm will fluctuate and increase very slowly.
- The angles become more stable and change very slowly while the norm increases rapidly.
- Easy examples: the angles get decreased enough for correct classification, the softmax cross-entropy loss can be well minimized by increasing the norm.
- Hard examples: the plateau is caused by unable to decrease the angle to correctly classify examples or increase the norms otherwise hurting loss.

**PART**

**03**

**Applications**

# Self-training and Domain Adaptation



# AVH for Self-training and Domain Adaptation

Replace Softmax-based confidence with AVH-based one during sample selection:

$$\mathcal{AVC}(c|\mathbf{x}; \mathbf{w}) = \frac{\pi - \mathcal{A}(\mathbf{x}, \mathbf{w}_c)}{\sum_{k=1}^K (\pi - \mathcal{A}(\mathbf{x}, \mathbf{w}_k))}$$

Similarly, AVH-based pseudo label

$$\hat{y}_t^{(k)*} = \begin{cases} 1, & \text{if } k = \arg \max_k \left\{ \frac{\mathcal{AVC}(k|\mathbf{x}_t; \mathbf{w})}{\lambda_k} \right\} \\ & \text{and } \mathcal{AVC}(k|\mathbf{x}_t; \mathbf{w}) > \lambda_k \\ 0, & \text{otherwise} \end{cases}$$

# Main Results

Method	Aero	Bike	Bus	Car	Horse	Knife	Motor	Person	Plant	Skateboard	Train	Truck	Mean
Source [51]	55.1	53.3	61.9	59.1	80.6	17.9	79.7	31.2	81.0	26.5	73.5	8.5	52.4
MMD [42]	87.1	63.0	76.5	42.0	90.3	42.9	85.9	53.1	49.7	36.3	<b>85.8</b>	20.7	61.1
DANN [16]	81.9	77.7	82.8	44.3	81.2	29.5	65.1	28.6	51.9	54.6	82.8	7.8	57.4
ENT [19]	80.3	75.5	75.8	48.3	77.9	27.3	69.7	40.2	46.5	46.6	79.3	16.0	57.0
MCD [50]	87.0	60.9	<b>83.7</b>	64.0	88.9	79.6	84.7	76.9	88.6	40.3	83.0	25.8	71.9
ADR [51]	87.8	79.5	<b>83.7</b>	65.3	<b>92.3</b>	61.8	<b>88.9</b>	73.2	87.8	60.0	85.5	32.3	74.8
Source [65]	68.7	36.7	61.3	<b>70.4</b>	67.9	5.9	82.6	25.5	75.6	29.4	83.8	10.9	51.6
CBST [65]	87.2	78.8	56.5	55.4	85.1	79.2	83.8	77.7	82.8	<b>88.8</b>	69.0	<b>72.0</b>	76.4
CRST [65]	88.0	79.2	61.0	60.0	87.5	81.4	86.3	78.8	85.6	86.6	73.9	68.8	78.1
Proposed	<b>93.3</b>	<b>80.2</b>	78.9	60.9	88.4	<b>89.7</b>	<b>88.9</b>	<b>79.6</b>	<b>89.5</b>	86.8	81.5	60.0	<b>81.5</b>

# Inner Metric

	TP Rate	AVH (avg)	Model Confidence	Norm $\ x\ $
CBST+AVH	0.844	0.118	0.961	20.84
CBST/CRST	0.848	0.117	0.976	21.28



**Examples chosen by  
AVH but not Softmax**

# AVH-based loss for Domain Generalization

AVH-based Loss:

$$\mathcal{L}_{AVH} = \sum_i \frac{\exp(s(\pi - \mathcal{A}(\mathbf{x}_i, \mathbf{w}_{y_i})))}{\sum_{k=1}^K \exp(s(\pi - \mathcal{A}(\mathbf{x}_i, \mathbf{w}_k)))}$$

Method	Painting	Cartoon	Photo	Sketch	Avg
AlexNet (Li et al., 2017)	62.86	66.97	89.50	57.51	69.21
MLDG (Li et al., 2018)	66.23	66.88	88.00	58.96	70.01
MetaReg (Balaji et al., 2018)	<b>69.82</b>	70.35	<b>91.07</b>	59.26	<b>72.62</b>
Feature-critic (Li et al., 2019)	64.89	<b>71.72</b>	89.94	<b>61.85</b>	72.10
Baseline CNN-9	66.46	67.88	89.70	51.72	68.94
CNN-9 + AVH	<b>71.56</b>	<b>69.25</b>	<b>89.93</b>	<b>60.86</b>	<b>72.90</b>

**PART**

**04**

**Conclusion**

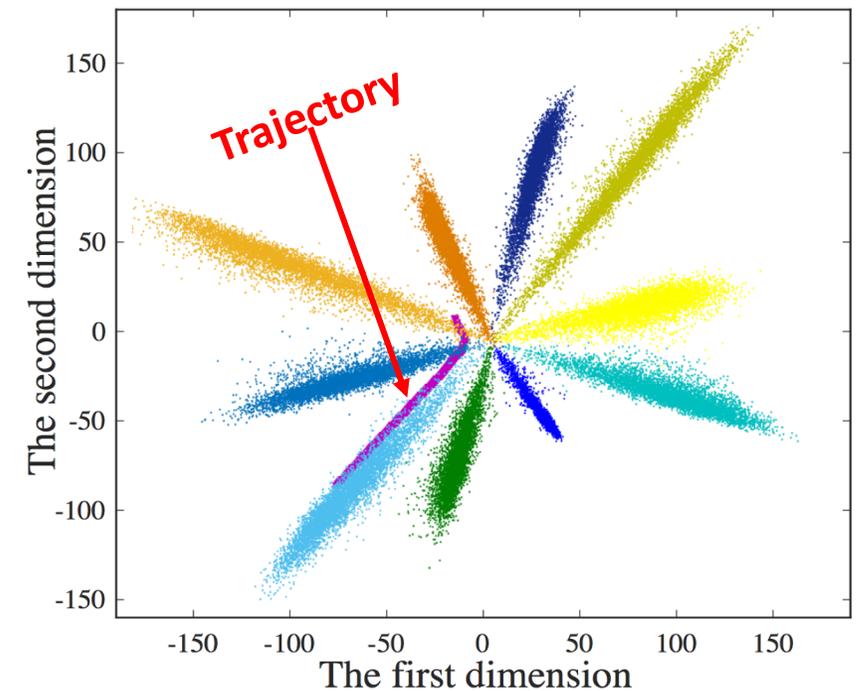
# Summary

- Propose AVH as a measure for visual hardness
- Validate that AVH has a statistically significant stronger correlation with human selection frequency
- Make observations on the dynamic evolution of AVH scores during ImageNet training
- Show the superiority of AVH with its application to self-training for unsupervised domain adaptation and domain generalization

# Discussions

- Connection to deep metric learning
- Connection to fairness in machine learning
- Connection to knowledge transfer and curriculum learning
- Uncertainty estimation (Aleatoric and Epistemic)
- Adversarial Example: A Counter Example?

Trajectory of an adversarial example switching from one class to another





# THANKS

Paper URL



Contact: [beidi.chen@rice.edu](mailto:beidi.chen@rice.edu)