

# Predicting Deep Neural Network Generalization with Perturbation Response Curves

Yair Schiff, Brian Quanz,  
Payel Das, and Pin-Yu Chen

NeurIPS 2021 Poster Session



# Problem statement and background: Predicting neural network generalization

- There is a gap in the literature for an **efficient and intuitive measure that can predict generalization of a deep neural network**
- Predicting Generalization in Deep Learning (PGDL) NeurIPS 2020 encouraged participants to provide ***complexity* measures calculated from network weights and training data to predict generalization gaps**

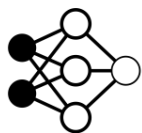


# Core idea



## User inputs

1. Trained model



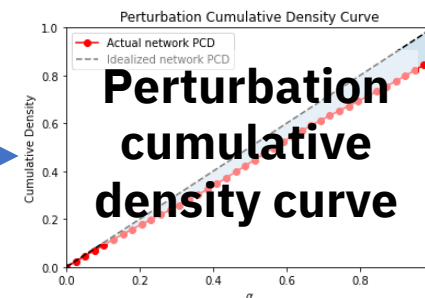
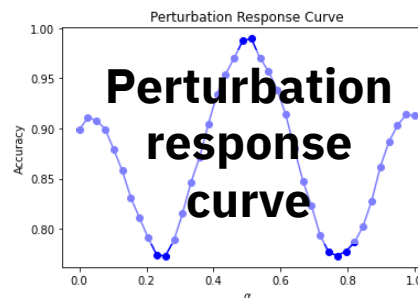
3. Parametric transformation

$\mathcal{T}_\alpha$

2. Training data



## Our proposed framework



**Output:**  
Gi-score  
&  
Pal-score



# Core idea



User inputs

1. Trained model



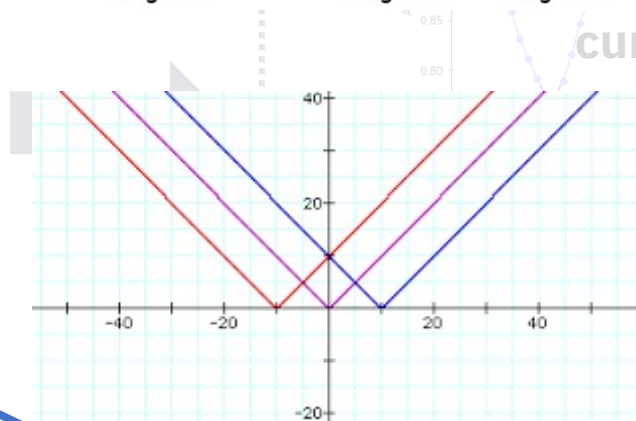
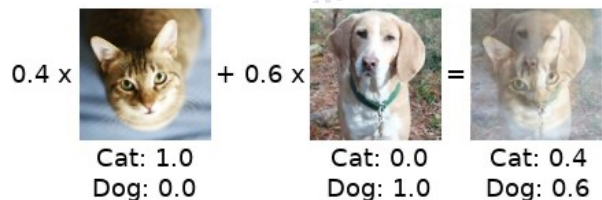
2. Training data



3. Parametric transformation

$\mathcal{T}_\alpha$

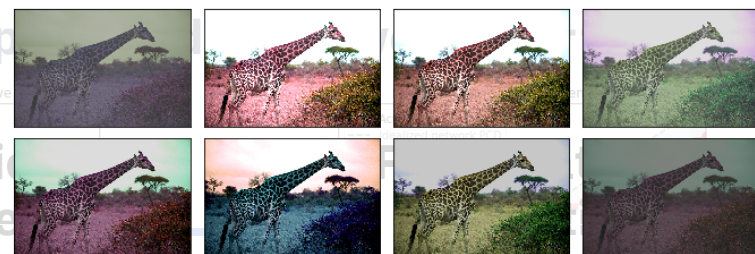
## Mixup



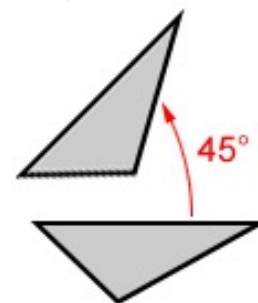
## Translation



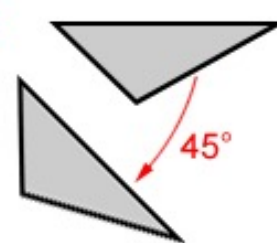
## Color jitter



density curve



## Rotation



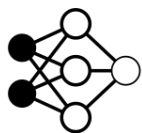


# Core idea



## User inputs

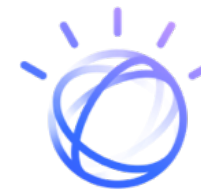
1. Trained model



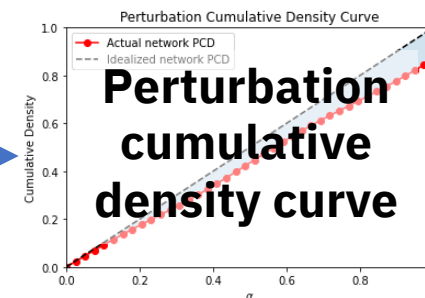
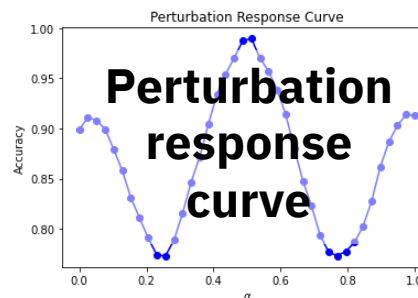
3. Parametric transformation

$\mathcal{T}_\alpha$

2. Training data



## Our proposed framework



**Output:**  
Gi-score  
&  
Pal-score

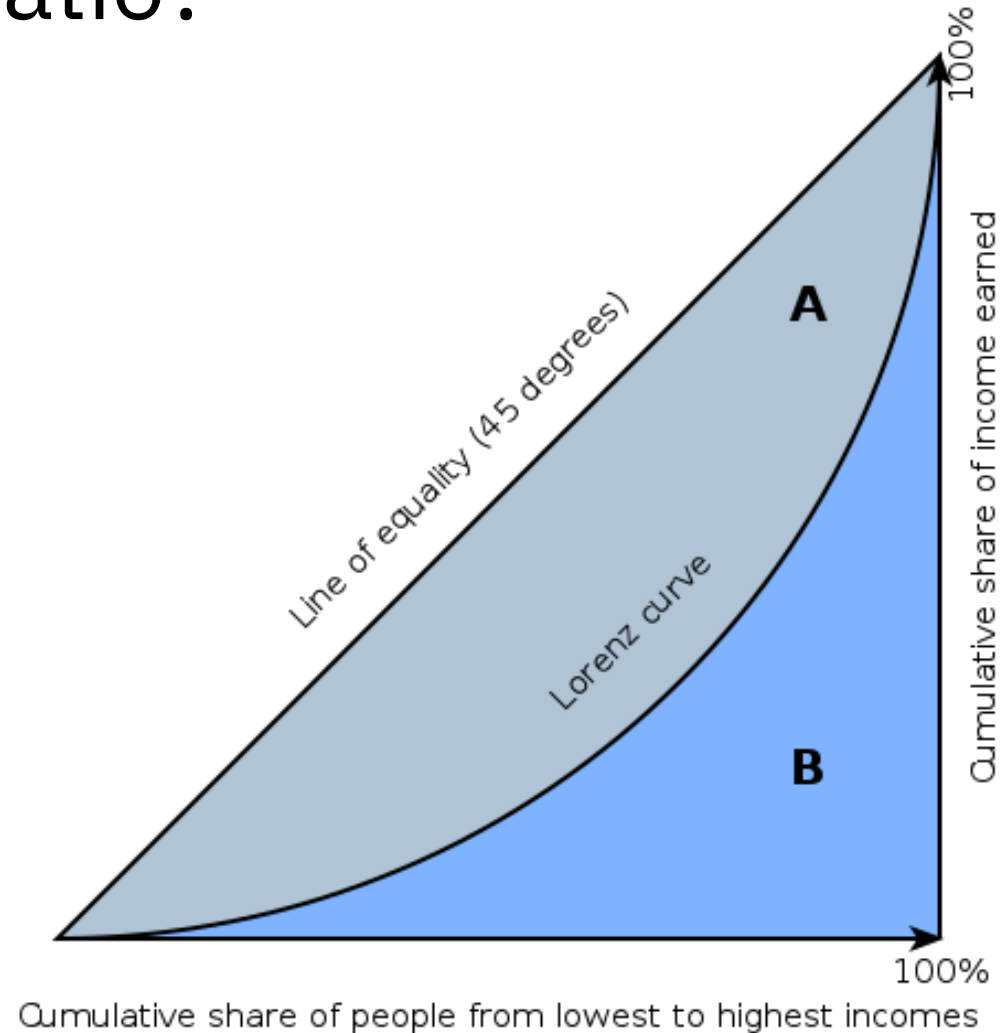


# Use cases and value

- This framework will be useful for **data scientists** and **machine learning practitioners**
- Useful for **predicting generalization and robustness** and defining new **regularization** approaches
- Our work can therefore serve as a **model selection criterion**, similar to  $R^2$  and other related statistics



# Interlude: What are the Gini coefficient and Palma ratio?





# Detailed description: Step 1 – Calculate Perturbation Response curve

---

## Algorithm 1: Building Perturbation Response (PR) Curve

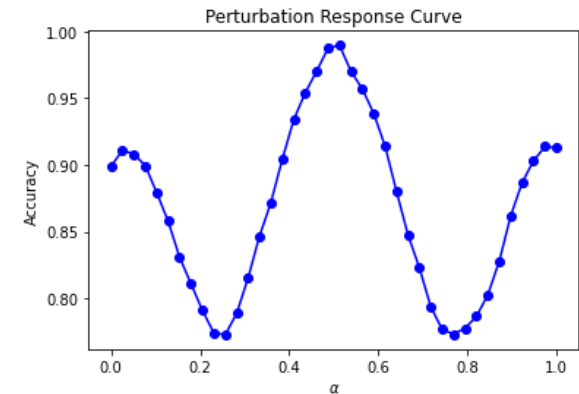
---

**Inputs:** Trained model  $f$ ; Dataset  $\mathcal{D}$ ; Perturbation  $\mathcal{T}_\alpha$ ; Min perturbation magnitude  $\alpha_{\min}$ ; Max perturbation magnitude  $\alpha_{\max}$ ; Number of perturbation magnitudes to measure  $n_p$ ; Layer at which to apply the perturbation  $\ell$ ; number of batches to sample  $n_b$ ; batch size  $b_s$

**Output:** PR Curve: Arrays of regularly spaced perturbation magnitudes ranging from  $\alpha_{\min}$  to  $\alpha_{\max}$  of length  $n_p$   $[\alpha_{\min}, \alpha_{\max}][n_p]$  and accuracy array at each perturbation magnitude of length  $n_p$   $\mathcal{A}_\alpha[n_p]$

```
for  $i \leftarrow 0$  to  $n_p - 1$  do
   $\alpha_i \leftarrow [\alpha_{\min}, \alpha_{\max}][i]$ 
  Shuffle  $\mathcal{D}$ 
  for  $k \leftarrow 0$  to  $n_b - 1$  do
     $\mathcal{D}_{sample} \leftarrow \mathcal{D}[kb_s : (k+1)b_s]$  // batch  $k$  of  $\mathcal{D}$ 
     $\mathcal{A}_{\alpha_i}^{(\ell)}[k] \leftarrow$  batch accuracy under perturbation  $\mathcal{T}_{\alpha_i}$  (Equation 1)
   $\mathcal{A}_\alpha[i] \leftarrow \sum_k \mathcal{A}_{\alpha_i}^{(\ell)}[k] / n_b$ 
```

---



- Works with any trained model, image dataset, and parametric perturbation
- Can be applied at any depth of a neural network





# Detailed description: Step 2 – Calculate $G_i$ and Pal scores

---

**Algorithm 2:**  $G_i$ -Score computation given PR Curve for a model

---

**Inputs:** Arrays of perturbation magnitude  $\alpha[n]$  and accuracy  $\mathcal{A}_\alpha[n]$

**Output:**  $G_i$ -score  $g_i$

$a_t[0] \leftarrow 0$  // initialize 1st element of trapezoidal areas array with 0

**for**  $i \leftarrow 0$  to  $n - 2$  **do**

$a_t[i + 1] \leftarrow 0.5(\alpha[i + 1] - \alpha[i])(\mathcal{A}_\alpha[i] + \mathcal{A}_\alpha[i + 1])$

**for**  $i \leftarrow 1$  to  $n - 1$  **do**

$a_t[i] \leftarrow a_t[i] + a_t[i - 1]$ . // cumulative sum

$d[i] = \alpha[i] - a_t[i], \forall i$

$g_i = 0$

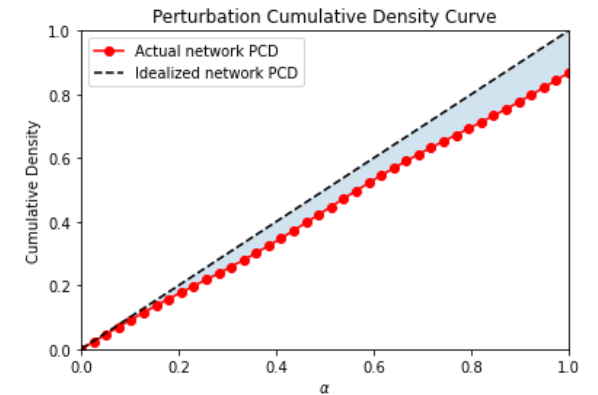
**for**  $i \leftarrow 0$  to  $n - 2$  **do**

$g_i \leftarrow g_i + 0.5(\alpha[i + 1] - \alpha[i])(d[i] + d[i + 1])$

$g_i \leftarrow g_i / (0.5\alpha[n - 1]^2)$  // Divide by area under line of equality

**return**  $g_i$

---



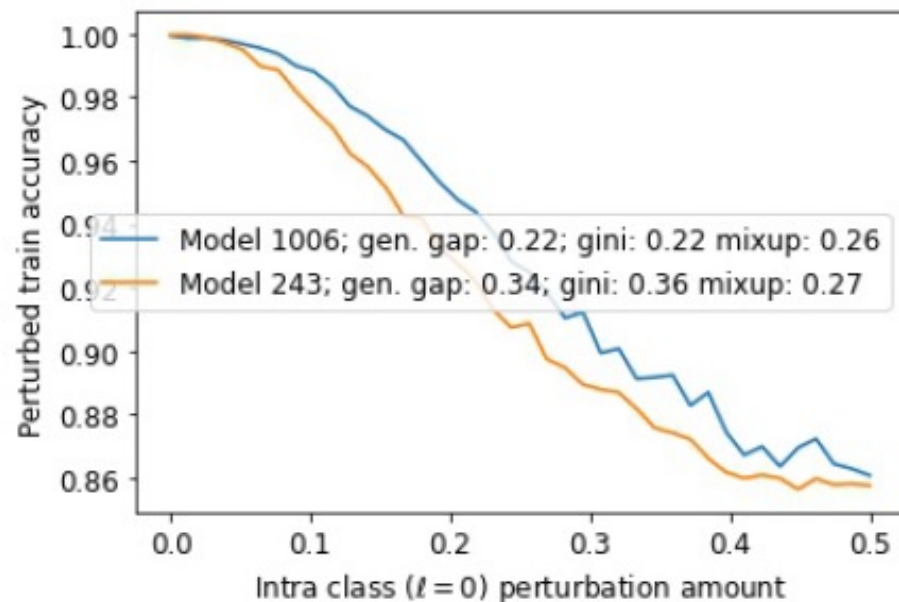


# PGDL results: We outperform winning team from PGDL competition in majority of tasks & overall

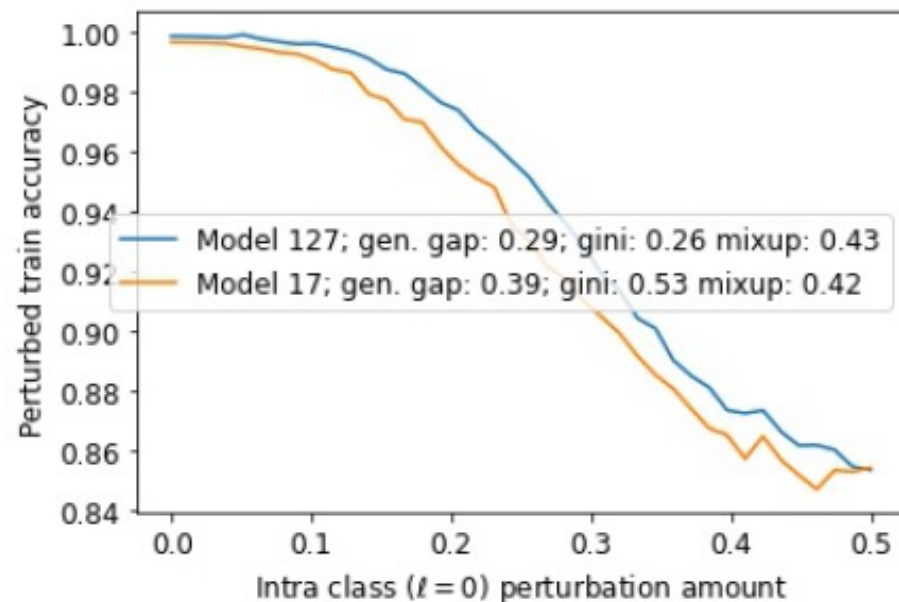
	CIFAR-10		SVHN	CINIC-10		Oxford Flowers	Oxford Pets	Fashion MNIST	<i>All Avg</i>
	VGG	NiN	NiN	Conv w/bn	Conv w/o bn	NiN	NiN	VGG	
<i>Single measures only</i>									
Gi inter $\ell=0$	3.03	<b>34.34</b>	26.58	21.01	6.96	33.05	<b>18.46*</b>	4.48	18.49
Gi inter $\ell=1$	<b>7.88</b>	22.59	12.17	12.58	8.39	7.52	4.68	<b>16.16*</b>	11.49
Pal inter $\ell=0$	3.14	26.39	24.25	21.11	6.37	29.62	15.96	4.21	16.38
Pal inter $\ell=1$	7.31	12.75	9.79	12.09	7.71	6.37	3.46	14.13	9.20
Gi intra $\ell=0$	0.84	30.54	<b>41.75*</b>	22.97	11.46	<b>42.44</b>	16.21	5.10	<b>21.41</b>
Gi intra $\ell=1$	0.22	17.18	10.96	9.50	12.43	6.92	3.60	5.55	8.29
Pal intra $\ell=0$	0.61	24.36	31.82	24.15	11.01	38.10	14.04	5.12	18.65
Pal intra $\ell=1$	0.44	10.34	13.48	8.68	11.09	5.88	3.02	6.25	7.40
Mixup	0.03	14.18	22.75	<b>30.30</b>	<b>19.51</b>	35.30	9.99	7.75	17.48
Mani. Mixup	2.24	2.88	12.11	4.23	4.84	0.03	0.13	0.19	3.33
<i>Combination measures</i>									
PCA Gi&Mix.	0.04	33.16	38.08	<b>33.76*</b>	<b>20.33*</b>	40.06	13.19	<b>10.30</b>	<b>23.62*</b>
Pal $\ell=0*\ell=1$	1.71	<b>35.77*</b>	<b>41.58</b>	25.14	9.50	38.92	<b>18.41</b>	5.61	22.08
Pal inter+intra	<b>24.84*</b>	29.70	14.04	1.64	3.45	14.84	2.13	4.89	11.94
DBI* <sup>1</sup> Mixup	0.00	25.86	32.05	31.79	15.92	<b>43.99*</b>	12.59	9.24	21.43



# Example PR Curve Pairs: Mixup $\alpha=0.5$ not enough to differentiate them



(a) SVHN



(b) CIFAR-10 NiN

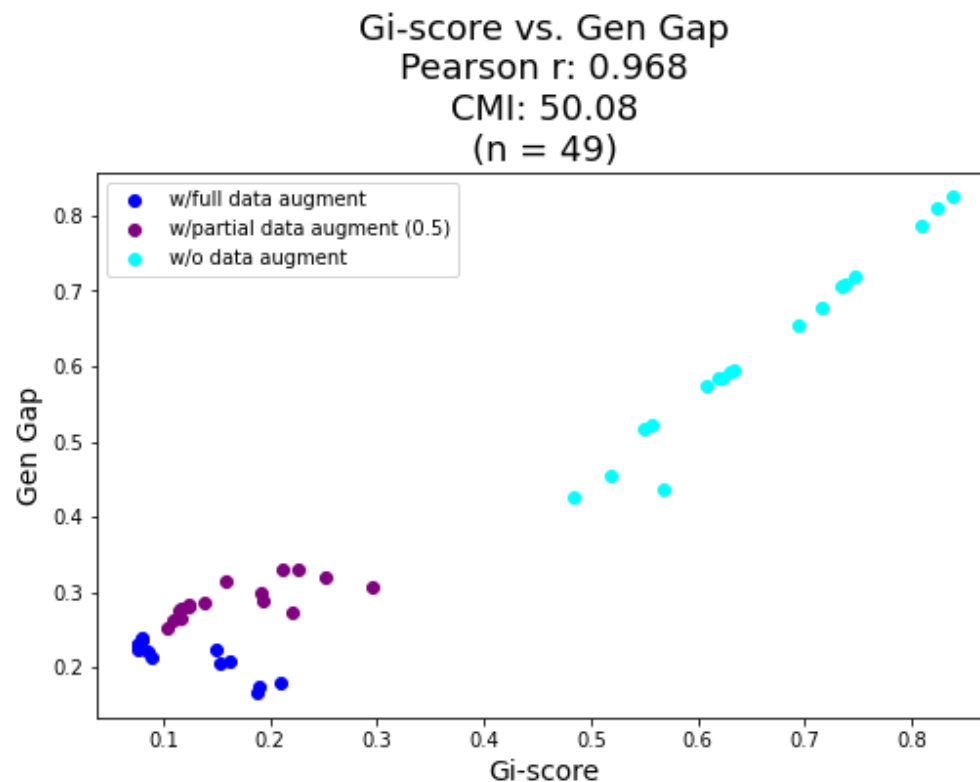
Figure 2: Examples of PR Curves with normalized scores and gen. gaps for 2 different models showing different performance fall-off captured by Gi intra score, but mixup scores roughly the same.



# Measuring invariance: experimental setup

	CIFAR-10	SVHN
Rotation	(-180, 179)	(-90, 90)
Horizontal translation	(-0.5, 0.5)	(-0.5, 0.5)
Vertical translation	(-0.5, 0.5)	(-0.5, 0.5)
Color jittering	(-0.25, 0.25)	(-0.25, 0.25)

Table 2: Perturbation minimum and maximum magnitudes by perturbation type and dataset. Minimum and maximums are displayed in each cell as an ordered pair.





# Measuring invariance:

We accurately predict degree of invariance across perturbation types

	CIFAR-10		SVHN	
	Resnet	VGG	Resnet	VGG
<i>Rotation</i>	( $n = 34$ )	( $n = 93$ )	( $n = 49$ )	( $n = 142$ )
Acc. on augmented train subset	27.99	<b>16.99</b>	47.24	42.97
Mean acc. on PR curve	27.61	15.61	48.14	44.05
Gi-score	<b>41.54</b>	15.29	<b>54.11</b>	<b>46.11</b>
<i>Horizontal translation</i>	( $n = 36$ )	( $n = 112$ )	( $n = 50$ )	( $n = 143$ )
Acc. on augmented train subset	41.79	33.48	29.49	24.20
Mean acc. on PR curve	45.03	33.00	29.88	24.28
Gi-score	<b>50.07</b>	<b>34.31</b>	<b>34.56</b>	<b>25.94</b>
<i>Vertical translation</i>	( $n = 36$ )	( $n = 107$ )	( $n = 49$ )	( $n = 141$ )
Acc. on augmented train subset	26.79	35.68	51.88	50.98
Mean acc. on PR curve	26.55	37.33	52.39	52.22
Gi-score	<b>34.85</b>	<b>39.07</b>	<b>59.02</b>	<b>52.83</b>
<i>Color-jittering</i>	( $n = 44$ )	( $n = 130$ )	( $n = 49$ )	( $n = 143$ )
Acc. on augmented train subset	35.77	28.44	43.08	<b>30.07</b>
Mean acc. on PR curve	39.12	29.37	43.26	28.67
Gi-score	<b>44.63</b>	<b>30.79</b>	<b>50.08</b>	29.45



# Summary

- We propose a **flexible framework** that provides high quality prediction of a trained neural network's generalization capability
- We provide multiple **new and efficient neural network generalization predictors**: Gi-score, Pal-score, and their combinations
- Our work **can be used with any parametric transformation** to compare the degree to which a network is invariant to that transformation

**Thank you!**