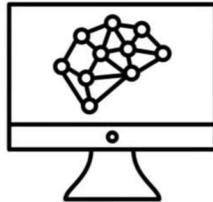


Estimating Example Difficulty using Variance of Gradients

Chirag Agarwal*, Sara Hooker*

Interpretability in current models?

I predict the suspect is guilty.



Artificial intelligence (AI)

How do you justify your judgment?



Judge

Justice system



Medical diagnosis

Motivation for Studying Relative Importance

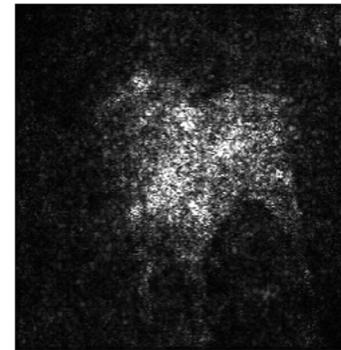
- o Automatically surfaces high priority examples for human inspection
- o Intra and Inter class understanding
- o Model bias

Gradient

Saliency maps can be noisy and ‘uninterpretable’



$$y = f(x)$$



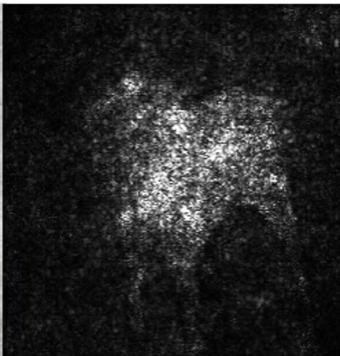
$$S = \nabla_x f(x)$$

Different variants of Gradients

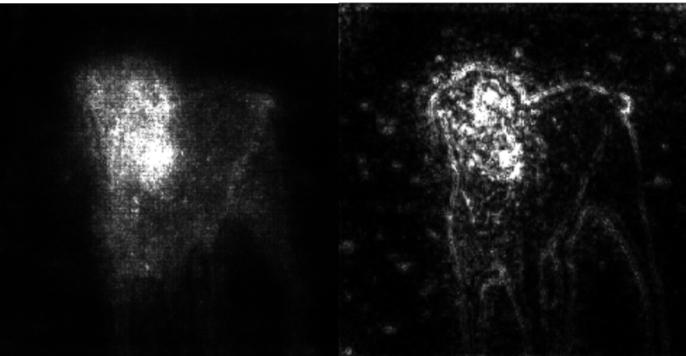
Gradient



SmoothGrad



GuidedBackProp



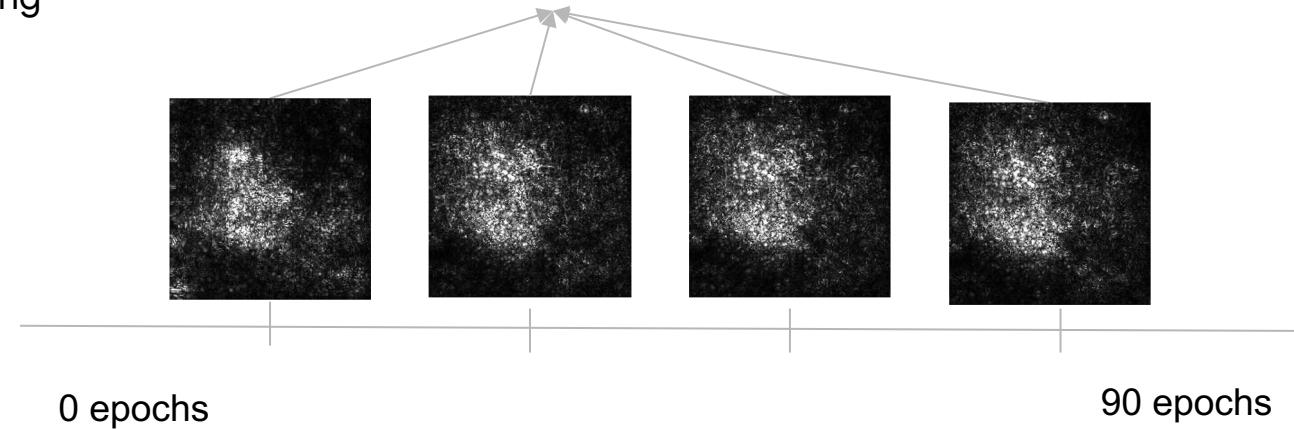
IntergratedGrad



Variance Of Gradients

$$VOG_i = \sqrt{\left(\frac{1}{K} \sum_{t=1}^K \frac{1}{N} (s_{ti} - \mu_i)^2 \right)}$$

Compute the average variance in gradients (VOG) for an image over training



Understand how feature importance forms over the course of training.



0 epochs

Early-stage training



90 epochs

Late Stage Training

Achille et al. 2017. Critical learning periods in deep neural networks. arXiv preprint arXiv:1711.08856.
Mangalam, K. and Prabhu, V. U. Do deep neural networks learn shallow learnable examples first. 2019.
Jiang et al. 2020.. Exploring the memorization-generalization continuum in deep learning. arXiv preprint arXiv:2002.03206

VOG computes a relative ranking of each class

Higher VOG score samples have cluttered backgrounds

Lowest VOG



Highest VOG



Lowest VOG



Highest VOG

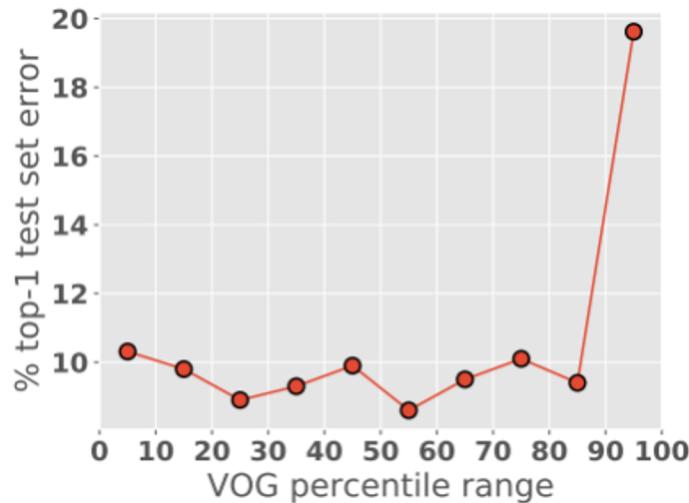


Early-stage training

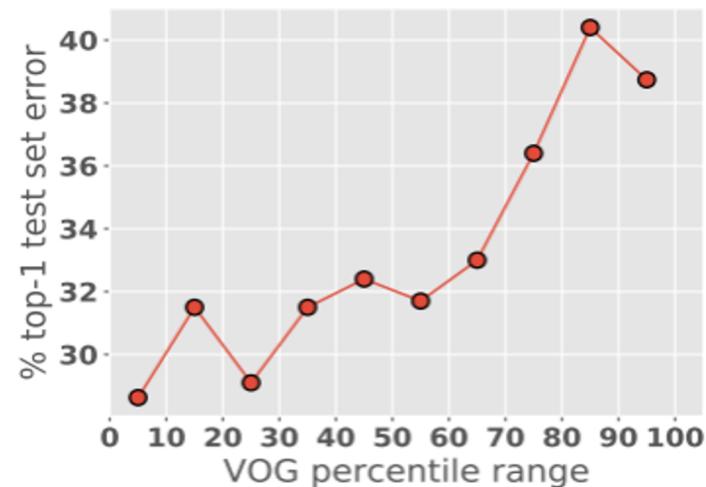
Late-stage training

VOG effectively discriminates easy & challenging examples

Mis-classification increases with an increase in VOG scores



CIFAR-10



CIFAR-100

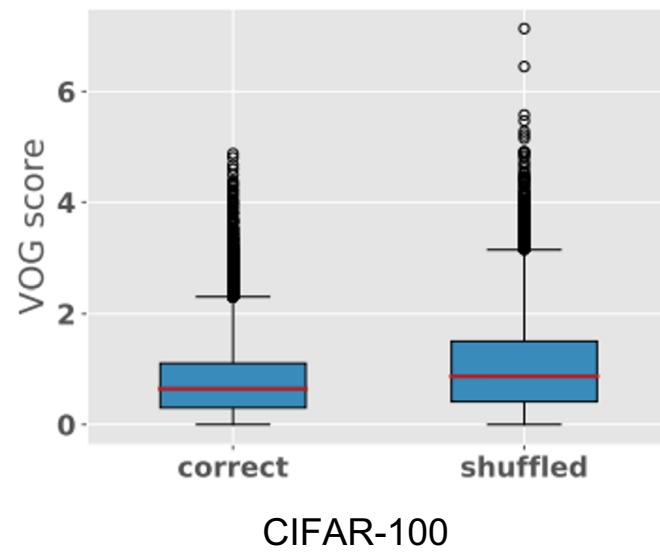
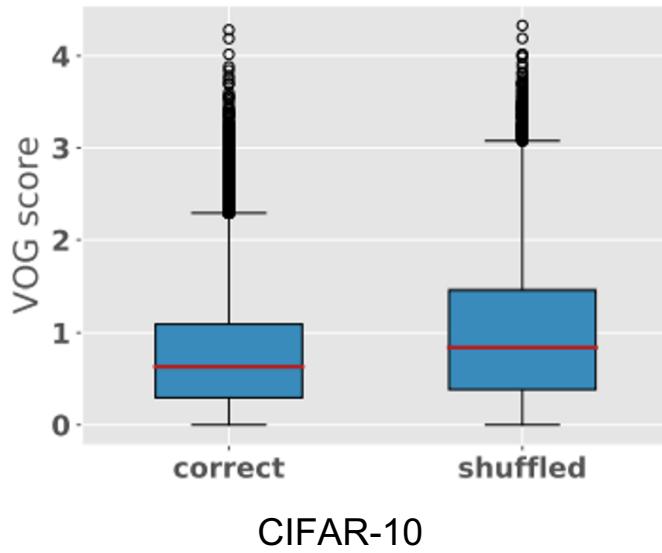
Memorization experiment

- o Overparameterized networks achieve zero training error by memorizing examples*
- o Replace 20% of all labels in the training set with random shuffled labels
- o Re-train the model from random initialization and compute VOG scores

*Zhang et al. Understanding deep learning requires rethinking generalization. arXiv preprint arXiv:1611.03530, 2016.

VOG distribution

Higher mean (red line) and spread for shuffled data

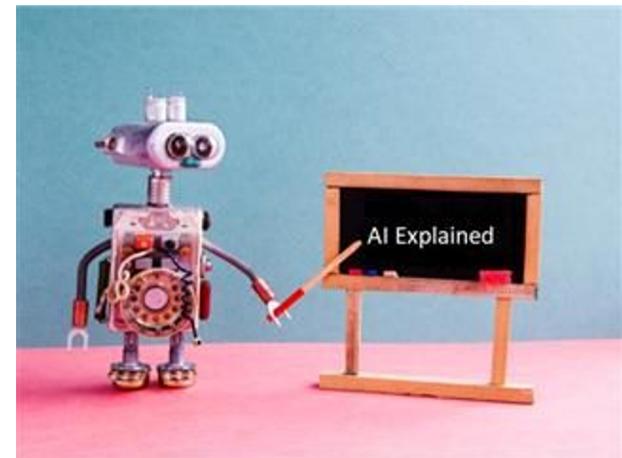


Future directions

- o Use VOG score to analyze complex datasets
- o Investigate various dataset and model bias
- o Leverage VOG to aid in curriculum learning

Final takeaways

- o An interpretability tool for ranking training and test examples
- o No need for modifying the architecture
- o VOG aids in clustering images with distinct visual properties



Paper: <https://arxiv.org/pdf/2008.11600.pdf>
chiragagarwall12@gmail.com
 @_cagarwal
shooker@google.com
@sarahookr