

# ABSTRACTION MECHANISMS PREDICT GENERALIZATION IN DEEP NEURAL NETWORKS

ALEX GAIN<sup>1</sup>

HAVA SIEGELMANN<sup>2</sup>

1. THE JOHNS HOPKINS UNIVERSITY

2. UNIVERSITY OF MASSACHUSETTS AMHERST

- We developed metrics inspired by neuroscience that are predictive of generalization and training dynamics.
- Key to our metrics is relating input complexity to deep representations.
- This provides a new direction to approaching generalization with potential tertiary applications.

- DNNs require deeper representations in order to classify harder examples. [1]
- We formalize this as, for input  $x$  and network  $net$ :
  - ▶  $\alpha(x) \triangleq$  difficulty of input  $\in \mathbb{R}$
  - ▶  $\beta(x, net) \triangleq$  use of deep representations  $\in \mathbb{R}$
- I.e. in DNNs, for the distribution trained on:  
 $\implies \alpha \propto \beta$

- To what extent is  $\alpha \propto \beta$ ?  
     $\implies$  Let  $\rho_{\alpha,\beta}$  be a measure of  $\alpha \propto \beta$ .
- **Hypothesis:** Higher  $\rho_{\alpha,\beta}$ 's correspond to “better” networks.
- **Question:** Can  $\rho_{\alpha,\beta}$  predict generalization? (Assuming good definitions of  $\alpha$  and  $\beta$ )

- Are there good definitions of  $\alpha, \beta$  for the brain?
- In [5], neuroscience work shows that, for the brain,  $\rho_{\alpha, \beta}$  is large:
  - ▶  $\alpha(T)$  = abstractness of task  $T$ , as measured by mechanical turk surveys.
  - ▶  $\beta(T)$  = use of “deeper” neurons, where “deeper” corresponds to distance from sensory cortices.

# DEFINING $\alpha$

- Back to our formalization of the result from [1]:
  - ▶  $\alpha(x) \triangleq$  difficulty of input
  - ▶  $\beta(x, \text{net}) \triangleq$  use of deep representations
- Empirically, we show  $\alpha(x) \propto$  input complexity, as measured by:
  1. Compression ratio of  $x$  calculated via compression algorithms (as done in [4]).
  2. Shannon entropy estimation via histogram binning of feature values of  $x$ :

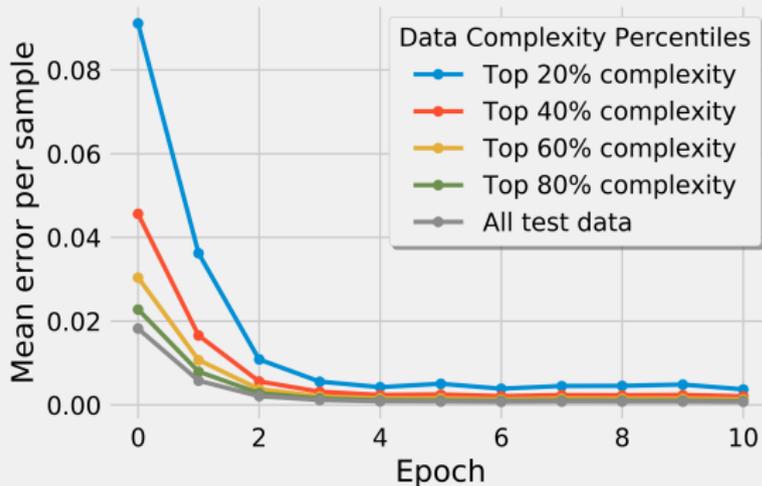
$$-\sum_{i=1}^N p_{bin_i} \log p_{bin_i}$$

for  $N$  bins where  $p_{bin_i}$  is proportional to the frequency of features in the range defined by bin  $i$ .

- Both give qualitatively similar results. We use method 2 for its ease of implementation.

Q: Does input complexity correspond to classification difficulty?

### Test Error for Different Complexity Scores



**Figure:** Mean test error versus training time for varying percentiles of input complexity. The model is an MLP trained on MNIST.

- Next, we define  $\beta(x)$  as:
  - ▶ The linear regressed slope on points  $(d, z_d(x))$  where  $d = 1, \dots, L$  for an  $L$  layer network, where  $z_d(x)$  is the sum of activation values for layer  $d$  of the network.
- $\beta$  as defined will be proportional to the use of deeper neurons – it is a coarse measure of use of depth.
- Note that  $\beta$  does *not* model the distribution of activation values versus depth (which is highly non-linear!). It is just a coarse measure that is practical for our purposes.

# COGNITIVE NEURAL ACTIVATION METRIC (CNA)

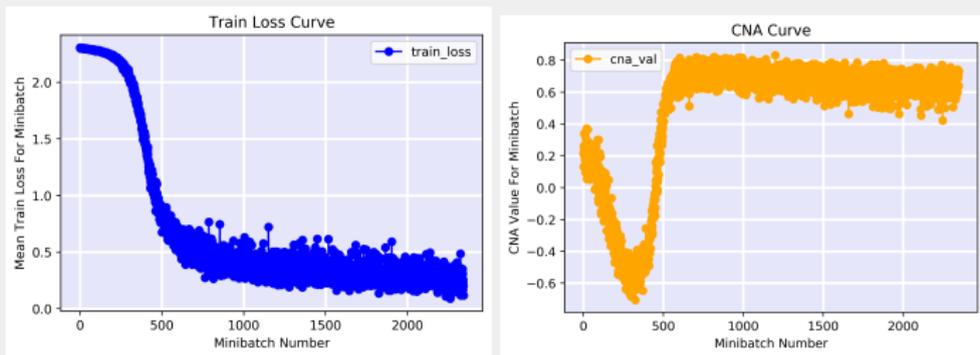
- The CNA (our metric for test performance) is defined for batch  $X$  and network  $net$ :

$$corr(\alpha(x), \beta(x, net) \mid x \in X)$$

where *corr* is Pearson correlation (though other correlations can be used).

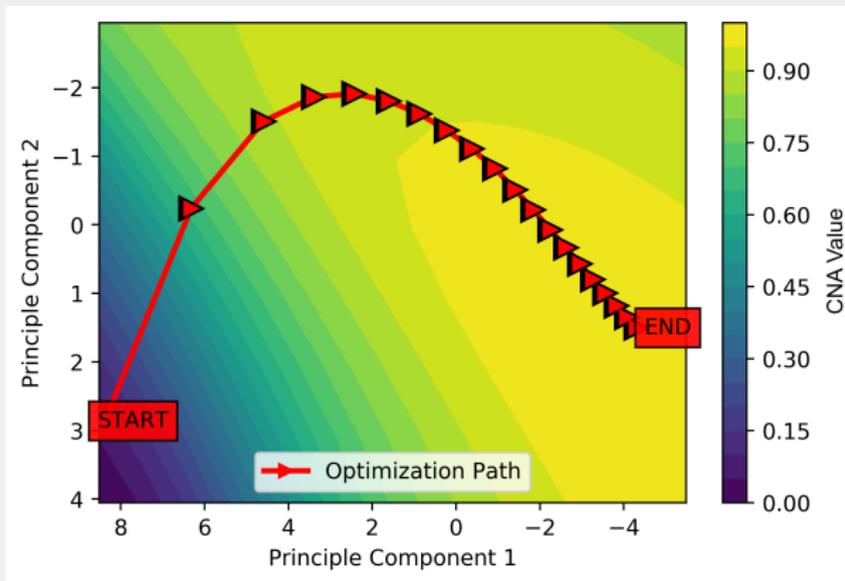
- Intuitively, well-performing networks should show higher CNA values.

# CNA RESULTS



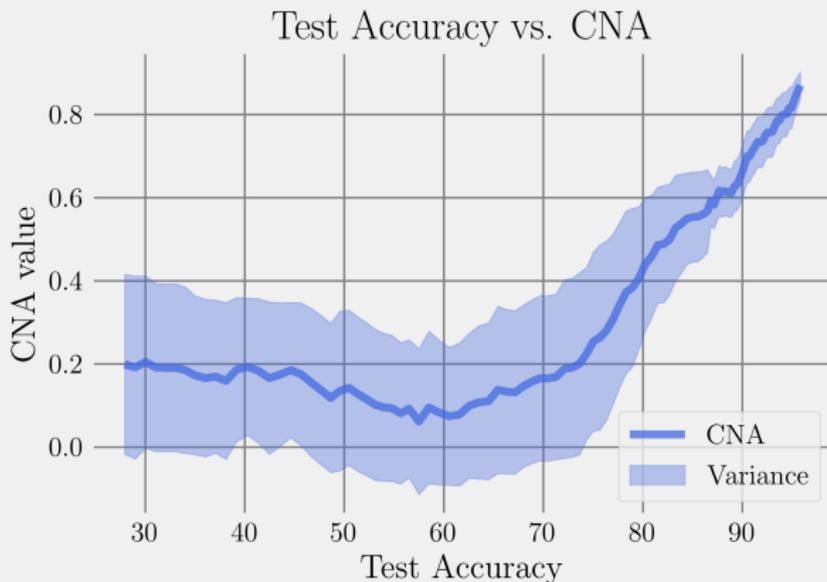
**Figure:** A simple MLP trained on MNIST with the training loss values (Left) and CNA values (Right) shown over training time. It is clear that the CNA shows a high correlation with training loss, with inflection points of both curves occurring at roughly the same timestep.

# CNA RESULTS



**Figure:** Network optimization trajectory (red curve) over training time visualized in 2D via PCA. The CNA gradient (contour background) is approximated via sampling from the principle component space and using the inverse PCA transformation.

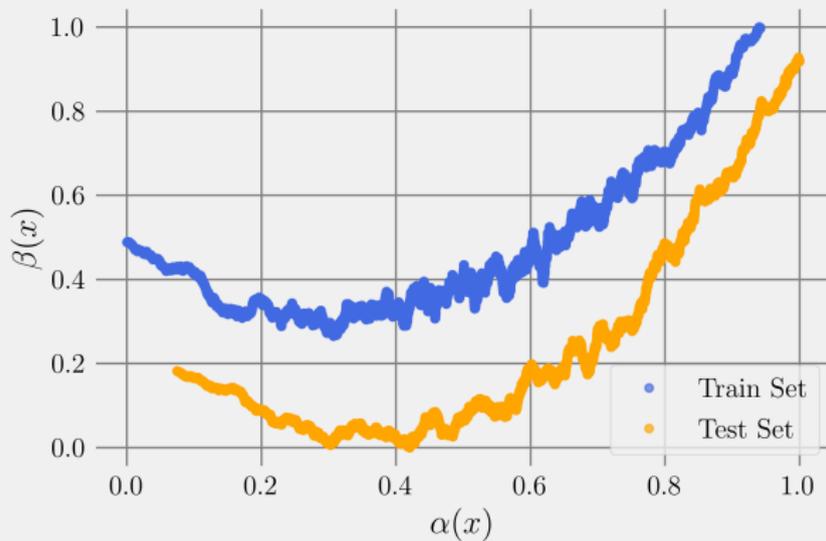
# CNA RESULTS



**Figure:** CNA vs Test Accuracy (moving average) for close to 200 network instances. Datasets: MNIST, FashionMNIST, CIFAR-10, CIFAR-100, ImageNet-32. Architectures: MLP, VGG-18, ResNet-18, ResNet-101.

- We develop another formulation focused on measuring the generalization gap.
  - ▶ *Generalization gap: difference in performance between training and test set for a given network and data distribution.*

# CNA-MARGIN: PREDICTING THE GENERALIZATION GAP



**Figure:**  $(\alpha, \beta)$ -curves for an MLP trained on SVHN, with the train set in blue and the test set in orange. The gap between the two curves corresponds to the gap between the train and test set performance for the network.

# CNA-MARGIN: PREDICTING THE GENERALIZATION GAP

- We define the CNA-Area (CNA-A) as the “area” between the two curves shown:

$$CNA_{\mathcal{A}}(X_{\text{train}}, X_{\text{test}}) \triangleq \max_{P \in \mathcal{S}} A(P)$$

where  $A(P)$  denotes the area of a polygon  $P$ , and  $\mathcal{S}$  denotes the set of polygons that can be inscribed between the curves.

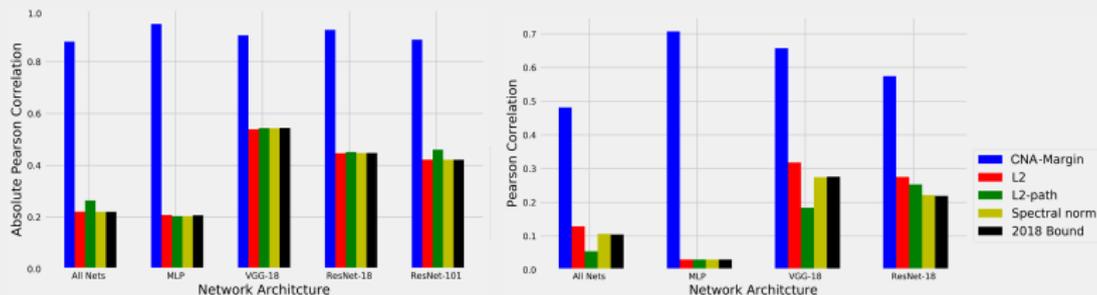
- The CNA-Margin (CNA-M) is simply CNA-A multiplied by the training margin  $\gamma_{\text{margin}}$  (see [3] for definition):

$$CNA_{\mathcal{M}}(X_{\text{train}}, X_{\text{test}}) \triangleq \gamma_{\text{margin}} \cdot CNA_{\mathcal{A}}(X_{\text{train}}, X_{\text{test}})$$

- This takes a similar form to sharpness-based generalization metrics, which are SOTA for generalization gap prediction [2], e.g. the L2 norm metric is proportional to:

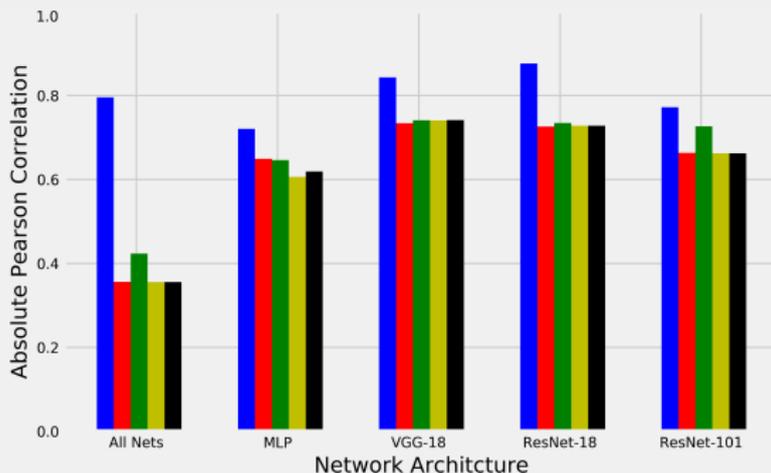
$$\frac{1}{\gamma_{\text{margin}}^2} \sum_k \|W_k\|_2$$

# CNA-M RESULTS



**Figure:** Generalization gap correlation for CNA-M and comparison generalization metrics. Left figure is w.r.t. standard datasets and a Gaussian noise dataset, right figure is w.r.t. standard datasets for varying degrees of shuffled labels (ranging from 10% to 50%).

# CNA-M RESULTS



**Figure:** Generalization gap correlation for CNA-M and comparison generalization metrics for standard datasets and training only (i.e. no Gaussian noise dataset included). Legend in previous slide.

# SUMMARY AND CONCLUSIONS

- CNA shows interesting properties w.r.t. training.
  - ⇒ Potentially useful as an unsupervised loss term and for understanding training dynamics.
- CNA-M shows excellent prediction of the generalization gap.
  - ⇒ Potentially can lead to insights into generalization not attainable by current directions.
- Tertiary contributions:
  - ▶ These nice results suggest studying DNNs under the  $\alpha - \beta$  framing in general can lead to more insights or applications.
  - ▶ We have formalized and provided direct empirical measures of “deeper representations are needed in order to classify harder examples.”

# REFERENCES



GAO HUANG, DANLU CHEN, TIANHONG LI, FELIX WU, LAURENS VAN DER MAATEN, AND KILIAN Q WEINBERGER.  
**MULTI-SCALE DENSE NETWORKS FOR RESOURCE EFFICIENT IMAGE CLASSIFICATION.**  
*arXiv preprint arXiv:1703.09844*, 2017.



YIDING JIANG, BEHNAMEH NEYSHABUR, HOSSEIN MOBAHI, DILIP KRISHNAN, AND SAMY BENGIO.  
**FANTASTIC GENERALIZATION MEASURES AND WHERE TO FIND THEM.**  
*arXiv preprint arXiv:1912.02178*, 2019.



BEHNAMEH NEYSHABUR, SRINADH BHOJANAPALLI, DAVID MCALLESTER, AND NATI SREBRO.  
**EXPLORING GENERALIZATION IN DEEP LEARNING.**  
In *Advances in Neural Information Processing Systems*, pages 5947–5956, 2017.



JOAN SERRÀ, DAVID ÁLVAREZ, VICENÇ GÓMEZ, OLGA SLIZOVSKAIA, JOSÉ F NÚÑEZ, AND JORDI LUQUE.  
**INPUT COMPLEXITY AND OUT-OF-DISTRIBUTION DETECTION WITH LIKELIHOOD-BASED GENERATIVE MODELS.**  
*arXiv preprint arXiv:1909.11480*, 2019.



P TAYLOR, JN HOBBS, J BURRONI, AND HT SIEGELMANN.  
**THE GLOBAL LANDSCAPE OF COGNITION: HIERARCHICAL AGGREGATION AS AN ORGANIZATIONAL PRINCIPLE OF HUMAN CORTICAL NETWORKS AND FUNCTIONS.**  
*Scientific reports*, 5(1):1–18, 2015.