

In Search of Probeable Generalization Measures

Jonathan Jaegerman^{1*}, Khalil Damouni^{1*}, Mahdi S. Hosseini^{2*}, Konstantinos N. Plataniotis¹

¹The Edward S. Rogers Sr. Department of Electrical & Computer Engineering, University of Toronto

²The Department of Electrical and Computer Engineering, University of New Brunswick

<https://github.com/mahdihosseini/GenProb>

Abstract—Understanding the generalization behaviour of deep neural networks is a topic of recent interest that has driven the production of many studies, notably the development and evaluation of generalization “explainability” measures that quantify model generalization ability. Generalization measures have also proven useful in the development of powerful layer-wise model tuning and optimization algorithms, though these algorithms require specific kinds of generalization measures which can probe individual layers. The purpose of this paper is to explore the neglected subtopic of probeable generalization measures; to establish firm ground for further investigations, and to inspire and guide the development of novel model tuning and optimization algorithms. We evaluate and compare measures, demonstrating effectiveness and robustness across model variations, dataset complexities, training hyperparameters, and training stages. We also introduce a new dataset of trained models and performance metrics, GenProb, for testing generalization measures, model tuning algorithms and optimization algorithms.

Index Terms—generalization, generalization measures, probeable generalization measures, eXplainability measures, complexity measures, quality metrics

I. INTRODUCTION

Deep learning has proven very successful this last decade, demonstrating time and time again its ability to generalize feature recognition from train data to test data. Despite all the attention, the underlying mechanisms in deep models that promote generalization are still open questions [1]. A number of papers attempt to consolidate the intuition behind deep learning generalization by developing “explainable” measures that attempt to quantify the generalization ability of a given model and dataset [2]–[4].

Generalization measures provide understanding of the learning mechanisms involved in the training of a given network using certain optimizers, datasets, and hyperparameters. Understanding and identifying trends in the relationship between these metrics and the network’s generalization gap or test accuracy provides methods of optimally selecting hyperparameter configuration, topologies, and other training parameters. Furthermore, generalization measures can be implemented for neural architecture search (NAS), hyperparameter optimization (HPO) and training optimization [5]–[9]. Other methods optimize deep models solely via holistic evaluations of performance, neglecting the discrepancies in training quality between layers, and foregoing the additional tuning required for a better solution.

Many generalization measures have proven effective and robust, but are not practical for implementation of model

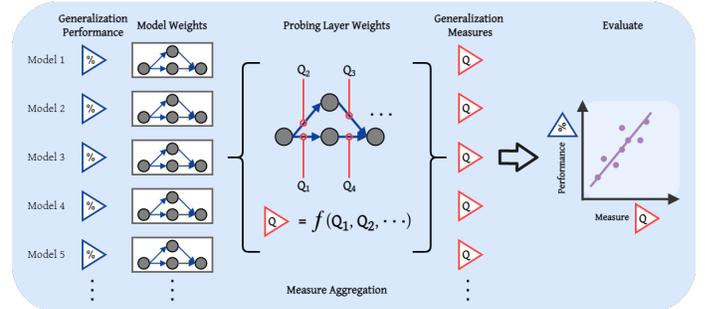


Fig. 1: Probeable Generalization Measure Evaluation Pipeline

tuning and optimization algorithms. It is optimal that a generalization measure can probe individual layers (a probeable measure), in order to tune models at a layer level (e.g. channel sizes, weight update rules, individual learning rates). The most successful measures are implemented from PAC-Bayes bounds and margin distributions, and are not probeable [2], [3], [10], [11]. Our objective in this paper is to investigate probeable generalization measures for model tuning and optimization applications, and to better understand generalization mechanisms in deep models. Given a dataset of trained models and their performance metrics, we evaluate probeable generalization metrics by the pipeline depicted in Fig. 1.

The following lists the contributions in this paper:

- We describe and investigate four probeable generalization measures that can be directly measured from individual layers of a deep network, without the need for additional training pipeline.
- We test these measures on the NATS-Bench dataset [12], and demonstrate their effectiveness for explaining the generalization of models of varying channel sizes.
- We study the evolution of these measures during training to understand their effectiveness at different stages of training.
- We introduce a new dataset—dubbed GenProb—to test probeable measures on models of varying channel sizes and training hyperparameters.
- We evaluate and compare the measures with GenProb, demonstrating effectiveness and robustness.

II. RELATED WORKS

With a recent surge in interest for understanding the generalization performance of deep models, studies have produced

many new complexity measures. Response to input perturbation measures [10], [11], [13], and PAC-Bayes sharpness and flatness measures [2], [10], [14] have demonstrated much success. Notably, the first and second place holders [10], [11] of the recent NeurIPS 2020: Predicting Generalization in Deep Learning competition [15], developed perturbation and sharpness measures. These measures have also been implemented in NAS, HPO and training optimization algorithms [6]–[8].

Studies also define measures from deep network margin distributions [16]–[18]. Some other measures are derived from model weights themselves, such as gradient signal-to-noise-ratio (GSNR) [19] and norm-based measures [5], [20]–[22]. Shallow networks are often trained onto the aforementioned measures for accurate generalization performance predictions, at the cost of ability to generalize predictions across different datasets and model structures [23]–[26]. Some large scale studies of generalization measures provide vigorous evaluations and comparisons of common and best performing measures, across different datasets and model structures [2]–[4].

These last three are the only studies that include probeable measures, however they are tested on sets of fully trained models with drastic variations; irrelevant for most model optimization algorithms. Probeable measures are implemented in [9] to reduce the computational cost and increase performance of NAS algorithms, and in [27] to adaptively optimize deep networks for greater generalization at no additional computational overhead.

III. METRIC DEFINITIONS

We define several measures that quantify the quality of a trained model (i.e. quality metrics or complexity measures) and describe its generalization ability. These quality metrics are probeable on individual layers of deep network, and quantify the contribution of each layer as a holistic measure for network representation, unlike other popular and successful measures elaborated above. The measures are thus explainable as they indicate how well training is optimized across layers of a deep network. The overview of all chosen measures is presented in Table I.

Stable quality (SQ) refers to the stability of encoding in a deep layer that is calculated with the relative ratio of stable rank and condition number of a layer defined in [27]. Stable rank encodes the space expansion under the matrix mapping of

TABLE I: Quality Metric Definitions

Metric	Formulation	Aggregation
Stable Quality	$\arctan s(\mathbf{W}_i)/\kappa(\mathbf{W}_i)$	$\prod_{i=1}^d (Q_{SQ}(\mathbf{W}_i))^{1/d}$
Effective Rank	$\sum_{i=1}^{n'} \bar{\sigma}_k(\mathbf{W}_i) \log(\bar{\sigma}_k(\mathbf{W}_i))$	$\log(\sqrt{\sum_{i=1}^d Q_{ER}(\mathbf{W}_i)^2/d})$
Spectral Norm	$\max(\boldsymbol{\sigma}(\mathbf{W}_i))$	$\log \sqrt{d(\prod_{i=1}^d \ \mathbf{W}_i\ _2^2)^{1/d}}$
Frobenius Norm	$\sqrt{\sum_{j=1}^m \sum_{k=1}^n w_{ijk} ^2}$	$\log \sqrt{d(\prod_{i=1}^d \ \mathbf{W}_i\ _F^2)^{1/d}}$

Here d is the depth of the model, n is the maximum rank of the weight matrix, n' is the rank of the weight matrix, $\boldsymbol{\sigma}(\mathbf{W}_i)$ is the vector of singular values of matrix \mathbf{W}_i with individual values $\sigma_k(\mathbf{W}_i)$ and normalized singular values $\bar{\sigma}_k(\mathbf{W}_i)$ and w_{ijk} is the value at row j and column k of matrix \mathbf{W}_i . Four dimensional weight tensors (such as convolutions) are first unfolded along their input and output channels for computation.

the layer, and condition number indicates the numerical sensitivity of the mapping layer. Altogether the measure introduces a quality measure of the layer as an autoencoder.

Effective rank (E) refers to the dimension of the output space of the transformation operated by a deep layer that is calculated with the Shannon entropy of the normalized singular values of a layer as defined in [28].

Frobenius norm (F) refers to the magnitude of a deep layer that is calculated with the sum of the squared values of a weight tensor. Frobenius norm is also calculated with the sum of the squared singular values of a layer.

Spectral norm (S) refers to the maximum magnitude of mapping by a transformation operated by a layer that is calculated as the maximum singular value of a weight tensor.

The formulations of these quality metrics are presented in Table I. With these layer-level quality measures, we aim to aggregate across model layers for a single meaningful model-level quality metric. For the norm-based measures we borrow the best performing aggregation method for each layer-level measure from [3]. For stable quality and effective rank, we inspire ourselves from literature aggregation methods, test all variations, and choose a single best performing method [2], [3].

The notation convention used in Table I and hereinafter to represent different quality metrics is: Q_M^{AGG} where aggregation $AGG \in \{L2 = \text{depth-normalized L2 norm}, p = \text{depth-normalized product}\}$ and metric $M \in \{SQ = \text{stable quality}, E = \text{effective rank}, F = \text{Frobenius norm}, S = \text{stable norm}\}$.

Low-rank factorization (LRF) is a preprocessing technique we also employ, aiming to increase the consistency of quality metrics across stages of training. By LRF, the low-rank component of a weight matrix, which represents the useful information, is extracted from raw weights, stripping the weights of residual noise from random initialization. We compute the factorization of our weight matrices by means of EVBMF defined in [29]. A wide hat $\hat{\cdot}$ can be included in the notation as \hat{Q}_M^{AGG} to indicate preprocessing of weights by low rank factorization.

IV. PRELIMINARY EXPERIMENTS

A. Quality Metrics and Varying Network Architecture

NATS-Bench is a dataset of trained deep neural networks for benchmarking neural architecture search algorithms [12]. A set of various related model architectures are trained and have their weights and performances saved to generate NATS-Bench, a mapping of model architectures and weights to generalization performance. NATS-Bench models are designed with a standard architecture skeleton as described in [12], and vary layer operations or channel sizes, at different locations in the skeleton. We can test the quality metrics with NATS-Bench by computing them on the provided trained weights, and correlating them with model generalization gap and test accuracy.

The NATS-Bench dataset consists of two subsets: the topology search space contains data from a set of models of varying layer operations, and the size search space contains data from

a set of models of varying layer channel sizes. The topology search space features a set of five operations (Table II) for six layers (15,625 total permutations), and the size search space features a set of eight channel sizes (Table II) for five layers (32,768 total permutations). The models are trained on CIFAR10, CIFAR100 and ImageNet-16-120 datasets [30]. Model weights were saved at 12 epochs, and at completion. Training was optimized with Nesterov momentum SGD for cross-entropy loss, with L2 weight decay and cosine annealing (see Table II for hyperparameter selection). Data augmentation included random flips with probability of 0.5, 32x32 (16x16 for ImageNet-16-120) random crops with 4 pixel padding, and RGB channel normalization.

Quality metrics are computed with and without use of the Low-Rank Factorization (LRF) on each set of trained model weights, then the Spearman correlations of these measures, with test accuracy and generalization gap, are computed to generate Fig. 2. Spearman rank-order correlation describes the quality of a relationship between two variables as an arbitrary non-parametric monotonic function. Spearman correlation therefore fits our purpose of evaluating generalization measures, as a generalization measure only needs to rank the relative performance of competing models.

Several measures show very high correlations with test accuracy and generalization gap in the NATS size search space, reaching to and above 0.9, as shown in Fig. 2(a). Most raw quality metrics show a lot of potential, though use of LRF preprocessing of weights does not provide a consistent improvement. A trend is not observed across different epochs; the patterns for 12 epochs can't be found in the 90 and 200 epoch alternative, indicating measure inconsistency across stages of training. Trends only hold across datasets in the case of NATS size search space for test accuracy, also indicating measure inconsistency across datasets. The patterns from the size search space do not appear for the topology search space; as illustrated in Fig. 2(b), correlations don't pass 0.5, and many of the previous peaks are now troughs. The measures are not valuable across models of changing topologies, as they are sensitive to and vary with the types of layers included in the model.

The NATS-Bench search spaces feature drastic changes between models which yield stochastic results and inhibit meaningful analysis; however, the high correlations under the size search space motivate further investigation. Notably for HPO and training optimization algorithms, it would be

TABLE II: NATS-Bench Hyperparameter Variation Summary

Hyperparameter	Size Search Space	Topology Search Space
Learning Rate	0.1 \rightarrow 0 (cosine)	
Weight Decay	5e-4	
Batch Size	256	
Epochs	12, 90	12, 200
Channel Size Variations	8, 16, 24, 32, 40, 48, 56, 64	-
Layer Operation Variations	-	zeroize, skip, 1x1 conv, 3x3 conv, average-pool

valuable to use a search space with variations in training hyperparameters instead. A last remark is that LRF boosts correlations for certain configurations (e.g. stable quality in the topology search space), which warrants further investigation.

B. Evolution of Quality Metrics during Training

We trained ResNet34 on CIFAR10 and CIFAR100 with AdaS [27] until completion and saved model weights at every epoch. The saved set of weights enables the observation of the evolution of quality metrics over training.

In Fig. 3(a) and Fig. 3(b) we observe high noise perturbation at initialization that quickly fades in the first 10 epochs, then rises subtly and plateaus at a low value. The residual noise extracted by LRF fades, as we'd expect of the random noise from initialization; it is replaced with learned structure, though it remains significant at later epochs with CIFAR100. All raw measures fail to mirror test accuracy and generalization gap at early epochs. Stable quality and effective rank measures begin mirroring test accuracy and generalization gap at around 10 epochs, and Frobenius and spectral norm measures only begin at the 60th epoch. The quality metrics do a better job of mirroring test accuracy and generalization gap at earlier



Fig. 2: Evolution of Generalization Measures during Training with and without LRF

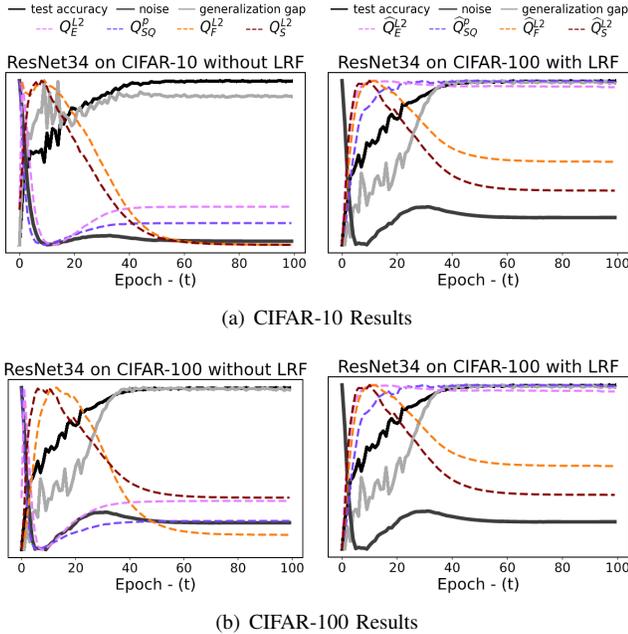


Fig. 3: Evolution of Quality Metrics during Training with LRF

epochs with LRF. Frobenius and spectral norm measures, however, still follow a distinct drop from epoch 10 to 60 with LRF. It intuitively follows that we can expect stable quality and effective rank measures to correlate better with test accuracy and generalization gap, notably with LRF preprocessing, through all stages of training.

V. RESULTS DEMONSTRATION

For a comprehensive display of experimental results please refer to the Supplementary Material.

A. Generating a Family of Trained Models

To test the effectiveness of the measures for tracking generalization performance at earlier stages of training, we train families of models with varied hyperparameter and channel size configurations, then save model weights and performances at each epoch. Models are trained for 70 epochs on CIFAR10 and CIFAR100 with various optimizers. We dub this dataset Generalization Dataset for Probeable Measures (GenProb).

The model architecture used is the same as that in NATS-Bench’s size search space [12], described in Table III. The convolutional blocks can be described as directed acyclic

TABLE III: GenProb Model Architecture

Block Index	Block Type	Output Shape
0	input	$32 \times 32 \times 3$
1	3×3 convolution	$32 \times 32 \times 8$
2	convolutional block	$32 \times 32 \times \{40, 48\}$
3	residual block	$18 \times 18 \times \{40, 48\}$
4	convolutional block	$18 \times 18 \times \{40, 48\}$
5	residual block	$9 \times 9 \times \{40, 48\}$
6	convolutional block	$9 \times 9 \times \{40, 48\}$
7	global average pooling	$1 \times 1 \times \{40, 48\}$
8	linear	$1 \times 1 \times \{10, 100\}$

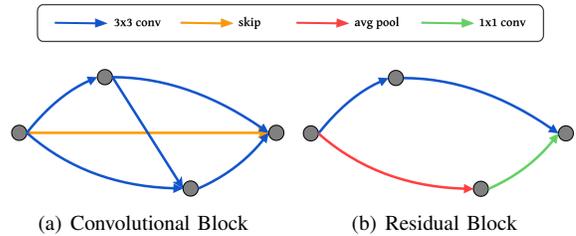


Fig. 4: GenProb Model Block Architecture

graphs with five nodes of activations, as depicted in Fig. 4(a). All nodes are ordered, and each node is connected to all nodes in front of it with a 3×3 convolution. The longest connection in the convolutional block (first node to last node) is replaced by a skip connection. The residual block is composed of a main path and a shortcut path, as illustrated in Fig. 4(b). The main path contains a 3×3 convolution (stride 2) followed by another 3×3 convolution (stride 1), and an average pooling layer (stride 2) followed by a 1×1 convolution in the shortcut path. All 3×3 convolutions are followed by batch normalizations, and all changes in channel sizes take place at the earliest layers possible in the blocks. All activation functions are ReLU, including the prediction layer. The output channel sizes of blocks 2-6 are varied, the output channel size of block 7 depends on block 6, and the output channel size of block 8 depends on the dataset.

The family of trained models is generated by varying initial learning rate, weight decay and channel sizes with the options listed in Table IV. Furthermore, input data is normalized channel-wise and augmented by random 32×32 crops with 4 pixels of padding and random horizontal flips of probability 0.5. Five repeated blocks in the model have their output channel sizes chosen, totaling $2^5 = 32$ permutations. From the set of all choices of hyperparameters we generate $5 \times 3 \times 32 = 480$ combinations, each yielding a unique trained model. We repeat with 7 different optimizers for both CIFAR10 and CIFAR100, training $480 \times 7 \times 2 = 6720$ unique configurations. Saving the model weights and performance at every epoch, we generate a total of $6720 \times 70 = 470400$ unique trained model files for our dataset.¹

B. Effective Rank, Test Accuracy and Generalization Gap

To visualize the relationship between the quality metrics and both generalization gap and test accuracy, we produce scatter plots of test accuracy and generalization gap over the quality metrics. Furthermore, by organizing these separate scatter plots relative to the quantity of training each set of models has undergone, we can study the evolution of the relationship. The quality metrics bias at earlier epochs due to residual noise from random initialization; we observe a lack of form in the effective rank scatter plots on models trained with Adam on CIFAR10 in Fig. 5(a) and Fig. 5(b) at earlier epochs.

¹GenProb will be released upon publication, alongside the GitHub repository.

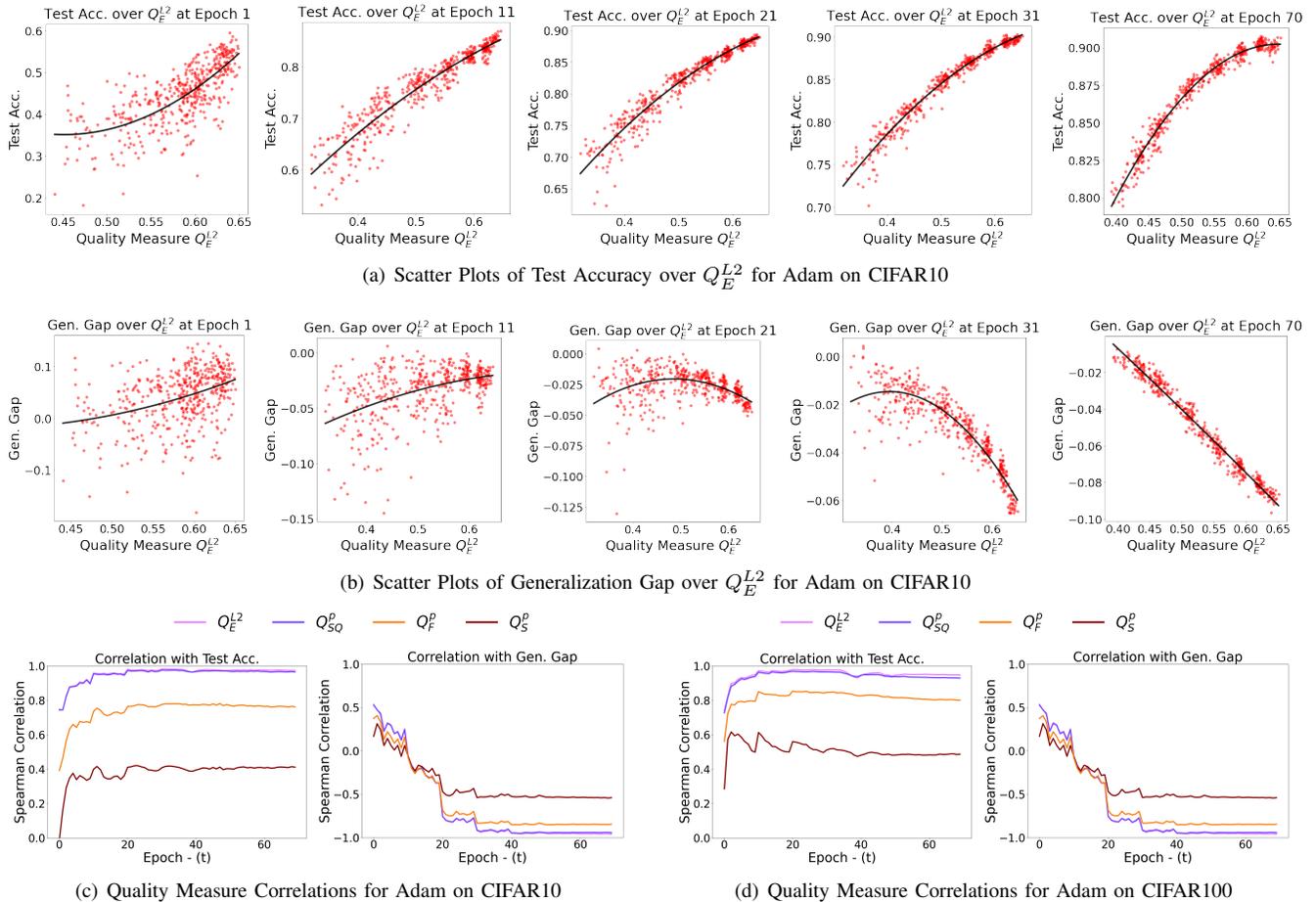


Fig. 5: Testing Quality Measures with GenProb

The quality metric evolves into clear, strong trends with test accuracy and generalization gap as training progresses and learned structure develops in the model weights. We observe similar behaviour with all three other metrics provided in the Supplementary Material (Fig. 1-16), and LRF only seems to weaken any trends, see Supplementary Material (Fig. 17-32).

We observe a clear linear relationship between the effective rank measure and generalization gap at later epochs, and a 2nd order relationship between the effective rank measure and test accuracy. The plateauing trend with test accuracy delineates a bound on test accuracy; maximizing effective rank above this bound would still increase generalization gap (linear trend) however, suggesting an increase in train accuracy

without changes in test accuracy. It is still evident that for a model trained on CIFAR10 with Adam, a greater effective rank indicates greater test accuracy, and a greater (negative) generalization gap. In fact, the relationships take form before model completion, as trends with test accuracy show as early as epoch 10, and trends with generalization gap show at epoch 40; this may be of interest for implementation of effective rank in NAS, HPO and optimization algorithms. We find similar timelines with all other quality metrics, see Supplementary Material (Fig. 1-32).

C. Quality Measure Correlation Evolution

By plotting the correlations of the quality metrics with test accuracy and generalization gap in Fig. 5(c) and Fig. 5(d), we can understand the relative progression of the effectiveness of these measures through different stages of training. As the aforementioned trends in Fig. 5(a) and Fig. 5(b) become more distinct, the corresponding correlations increase in magnitudes, some nearly up to 1. LRF preprocessing of weights only decreases the magnitude of correlation through all stages of training, see Supplementary Material (Fig. 41-42)

The large correlations indicate robustness to changes in training hyperparameters, and model channel sizes. Effective

TABLE IV: GenProb Hyperparameter Variation Summary

Hyperparameter	RMSGD	SGD	SAM	SGDP	Adam	AdaBound	AdamP
Learning Rate	0.02, 0.04, 0.06, 0.08, 0.10	2e-3, 4e-3, 6e-3, 8e-3, 0.01					
Weight Decay			1e-4, 5e-4, 1e-3				
Channel Sizes			40, 48				
Momentum			0.9				
Batch Size			256				
Epochs			70				
Scheduler Step Size	-			10			

rank and stable quality measures prove to be the most effective and robust generalization measures through all training phases and across dataset complexities. The Frobenius and spectral norm measures underperform consistently and are less consistent across dataset complexities, yielding lower correlations with CIFAR10. Also, as previously highlighted in the scatter plots, correlations with test accuracy plateau after 10 epochs and correlations with generalization gap plateau after 40 epochs.

VI. CONCLUSION

In this work we investigated "explainable" generalization measures that can probe individual layers of a deep neural network. We tested four quality metrics, that are calculated from layer weight tensors, over spaces of similar model architectures for NAS (NATS-Bench). We analyzed the evolution of the quality metrics during training, then produced a dataset for more meaningful analysis and testing related to HPO and training optimization. We introduced GenProb, a dataset of 470,400 trained models and their performance metrics, distinguished by hyperparameter and channel size variations. We demonstrated effective rank and stable quality measure effectiveness and robustness to variations in training hyperparameters, channels sizes, dataset complexities and stages of training. Furthermore, we investigate the amount of training required to produce meaningful generalization measures from model weights, and the shape of the relationships of these measures with test accuracy and generalization gap.

We hope this work inspires and guides the production of novel NAS, HPO and training optimization algorithms that leverage probeable generalization measures to maximize model generalization performance at a layer level. We also aim to motivate further development of probeable generalization measures, for which GenProb will prove a useful tool. Our investigated measures don't prove robust to drastic variations in model architecture, and as such may not be suitable for all NAS algorithms. We also only consider simple convolutional neural networks and two datasets (CIFAR10 and CIFAR100); we hope that future investigations will experiment with other types of model, deeper models, and other datasets.

REFERENCES

- [1] H. Minyoung, M. Hossein, Z. Richard, C. Brian, A. Pulkit, and I. Phillip, "The low-rank simplicity bias in deep networks," 2021.
- [2] Y. Jiang, B. Neyshabur, H. Mobahi, D. Krishnan, and S. Bengio, "Fantastic generalization measures and where to find them," 2019.
- [3] G. K. Dziugaite, A. Drouin, B. Neal, N. Rajkumar, E. Caballero, L. Wang, I. Mitliagkas, and D. M. Roy, "In search of robust measures of generalization," 2021.
- [4] B. Neyshabur, S. Bhojanapalli, D. McAllester, and N. Srebro, "Exploring generalization in deep learning," 2017.
- [5] B. Neyshabur, R. R. Salakhutdinov, and N. Srebro, "Path-sgd: Path-normalized optimization in deep neural networks," in *Advances in Neural Information Processing Systems*, C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett, Eds., vol. 28. Curran Associates, Inc., 2015. [Online]. Available: <https://proceedings.neurips.cc/paper/2015/file/ea32c96f620053cf442ad32258076b9-Paper.pdf>
- [6] P. Foret, A. Kleiner, H. Mobahi, and B. Neyshabur, "Sharpness-aware minimization for efficiently improving generalization," *CoRR*, vol. abs/2010.01412, 2020. [Online]. Available: <https://arxiv.org/abs/2010.01412>
- [7] J. J. Cheria, A. G. Taube, R. T. McGibbon, P. Angelikopoulos, G. Blanc, M. Snarski, D. D. Richman, J. L. Klepeis, and D. E. Shaw, "Efficient hyperparameter optimization by way of pac-bayes bound minimization," 2020.
- [8] K. Kandasamy, W. Neiswanger, J. Schneider, B. Póczos, and E. P. Xing, "Neural architecture search with bayesian optimisation and optimal transport," *CoRR*, vol. abs/1802.07191, 2018. [Online]. Available: <http://arxiv.org/abs/1802.07191>
- [9] M. S. Abdelfattah, A. Mehrotra, Łukasz Dudziak, and N. D. Lane, "Zero-cost proxies for lightweight nas," 2021.
- [10] P. Natekar and M. Sharma, "Representation based complexity measures for predicting generalization in deep learning," 12 2020.
- [11] S. A. K. D. Kashyap, and N. Subramanyam, "Robustness to augmentations as a generalization metric," 2021.
- [12] X. Dong, L. Liu, K. Musial, and B. Gabrys, "Nats-bench: Benchmarking nas algorithms for architecture topology and size," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, p. 1–1, 2021. [Online]. Available: <http://dx.doi.org/10.1109/TPAMI.2021.3054824>
- [13] Y. Schiff, B. Quanz, P. Das, and P. Chen, "Predicting deep neural network generalization with perturbation response curves," *CoRR*, vol. abs/2106.04765, 2021. [Online]. Available: <https://arxiv.org/abs/2106.04765>
- [14] G. K. Dziugaite and D. M. Roy, "Computing nonvacuous generalization bounds for deep (stochastic) neural networks with many more parameters than training data," *CoRR*, vol. abs/1703.11008, 2017. [Online]. Available: <https://arxiv.org/abs/1703.11008>
- [15] Y. Jiang, P. Foret, S. Yak, D. M. Roy, H. Mobahi, G. K. Dziugaite, S. Bengio, S. Gunasekar, I. Guyon, and B. Neyshabur, "Neurips 2020 competition: Predicting generalization in deep learning," 2020.
- [16] C. Chuang, Y. Mroueh, K. H. Greenewald, A. Torralba, and S. Jegelka, "Measuring generalization with optimal transport," *CoRR*, vol. abs/2106.03314, 2021. [Online]. Available: <https://arxiv.org/abs/2106.03314>
- [17] S. Arora, R. Ge, B. Neyshabur, and Y. Zhang, "Stronger generalization bounds for deep nets via a compression approach," 2018.
- [18] P. Bartlett, D. J. Foster, and M. Telgarsky, "Spectrally-normalized margin bounds for neural networks," 2017.
- [19] J. Liu, G. Jiang, Y. Bai, T. Chen, and H. Wang, "Understanding why neural networks generalize well through gsnr of parameters," 2020.
- [20] V. Nagarajan and J. Z. Kolter, "Generalization in deep networks: The role of distance from initialization," *CoRR*, vol. abs/1901.01672, 2019. [Online]. Available: <http://arxiv.org/abs/1901.01672>
- [21] T. Liang, T. A. Poggio, A. Rakhlin, and J. Stokes, "Fisher-rao metric, geometry, and complexity of neural networks," *CoRR*, vol. abs/1711.01530, 2017. [Online]. Available: <http://arxiv.org/abs/1711.01530>
- [22] R. Novak, Y. Bahri, D. A. Abolafia, J. Pennington, and J. Sohl-Dickstein, "Sensitivity and generalization in neural networks: an empirical study," 2018.
- [23] Y. Jiang, D. Krishnan, H. Mobahi, and S. Bengio, "Predicting the generalization gap in deep networks with margin distributions," in *International Conference on Learning Representations*, 2019. [Online]. Available: <https://openreview.net/forum?id=HJIQfnCqKX>
- [24] S. Yak, J. Gonzalvo, and H. Mazzawi, "Towards task and architecture-independent generalization gap predictors," 2019.
- [25] C. Corneanu, M. Madadi, S. Escalera, and A. Martinez, "Computing the testing error without a testing set," 2020.
- [26] T. Unterthiner, D. Keysers, S. Gelly, O. Bousquet, and I. Tolstikhin, "Predicting neural network accuracy from weights," 2021.
- [27] M. S. Hosseini and K. N. Plataniotis, "Adas: Adaptive scheduling of stochastic gradients," *CoRR*, vol. abs/2006.06587, 2020. [Online]. Available: <https://arxiv.org/abs/2006.06587>
- [28] O. Roy and M. Vetterli, "The effective rank: A measure of effective dimensionality," 01 2007.
- [29] S. Nakajima, M. Sugiyama, and S. Babacan, "Global solution of fully-observed variational bayesian matrix factorization is column-wise independent," in *Advances in Neural Information Processing Systems*, J. Shawe-Taylor, R. Zemel, P. Bartlett, F. Pereira, and K. Q. Weinberger, Eds., vol. 24. Curran Associates, Inc., 2011. [Online]. Available: <https://proceedings.neurips.cc/paper/2011/file/b73ce398c39f506af761d2277d853a92-Paper.pdf>
- [30] P. Chrabaszcz, I. Loshchilov, and F. Hutter, "A downsampled variant of imagenet as an alternative to the cifar datasets," 2017.