

ON THE GENERALIZATION MYSTERY IN DEEP LEARNING

SATRAJIT CHATTERJEE AND PIOTR ZIELINSKI

ABSTRACT. The generalization mystery in deep learning is the following: Why do over-parameterized neural networks trained with gradient descent (GD) generalize well on real datasets even though they are capable of fitting random datasets of comparable size? Furthermore, from among all solutions that fit the training data, how does GD find one that generalizes well (when such a well-generalizing solution exists)?

We argue that the answer to both questions lies in the interaction of the gradients of different examples during training. Intuitively, if the per-example gradient vectors are well-aligned, that is, if they are *coherent*, then one may expect GD to be (algorithmically) stable, and hence generalize well. We formalize this argument with an easy to compute and interpretable metric for coherence, and show that the metric takes on very different values on real and random datasets for several common vision networks. The theory also explains a number of other phenomena in deep learning, such as why some examples are reliably learned earlier than others, why early stopping works, and why it is possible to learn from noisy labels. Moreover, since the theory provides a causal explanation of how GD finds a well-generalizing solution when one exists, it motivates a class of simple modifications to GD based on robust averaging of per-example gradients that attenuate memorization and improve generalization.

Generalization in deep learning is an extremely broad phenomenon, and therefore, it requires an equally general explanation. We conclude with a survey of alternative lines of attack on this problem, and argue that the proposed approach is the most viable one on this basis.

CONTENTS

1. Introduction	2
2. The Theory, Informally	3
3. An Illustrative Example	6
4. Metrics to Quantify Coherence	8
5. Bounding the Generalization Gap with α	11
6. Measuring α on Real and Random Datasets	12
7. From Measurement to Control: Suppressing Weak Descent Directions	16
8. Why are Some Examples (Reliably) Learned Earlier?	21
9. Learning With Noisy Labels	24
10. Depth, Feedback Loops, and Signal Amplification	25
11. What Should a Theory of Generalization Look Like?	31
12. Comparison with Other Theories and Explanations	33
13. Discussion and Directions for Future Work	39
A. Mathematical Properties of α	42
B. Comparison of α with Other Metrics	46
C. Proof of The Generalization Theorem	48
D. Methods to Measure α	56
E. Measuring α on Additional Datasets and Architectures	57
F. The Evolution of Coherence	61

G. Experimental Details of Easy and Hard Examples	63
H. The Under-Parameterized Case: A Preliminary Look	65
I. Additional Data	66
References	67

1. INTRODUCTION

In spite of the tremendous practical success of deep learning, we do not yet have a good understanding of why it works. Deep neural networks used in practice are over-parameterized, that is, they have many more parameters than the number of examples that are used to train them, and conventional wisdom holds that such over-parameterized models should not generalize well, but yet they do.¹

Although this gap in our understanding has been long known (see, for example, Bartlett [1996] and Neyshabur et al. [2014]), the problem was sharpened in an influential paper by Zhang et al. [2017] who showed that typical neural networks trained with stochastic gradient descent, which is the usual training method, can easily memorize a random dataset of the same size as the (real) dataset that they were designed for. They argued that this simple experimental observation poses a challenge to all known explanations of generalization in deep learning, and called for a “rethinking” of our approach. This led to a large effort in the community to better understand why neural networks generalize. However, although our understanding of deep learning has greatly improved as a result of this effort, to date, there does not appear to be an satisfactory explanation (see Zhang et al. [2021] for a detailed review).

Our goal in this work is to answer the broad question raised by the observations of Zhang et al., namely,

Why do neural networks generalize well in practice when they have sufficient capacity to memorize their training set?

Specifically, we want to answer the following questions:

- Q1.** Since by simply changing the dataset in an over-parameterized setting (for example, from real labels to random labels), we can obtain very different generalization performance, what property of the dataset controls the generalization gap (assuming, of course, that architecture, learning rate, size of training set, etc. are held fixed)? We stress that our interest is in the *gap*, that is the difference between training and test loss (and not in the test loss *per se*).
- Q2.** Why does gradient descent not simply memorize real training data as it does random training data? That is, from among all the models that fit the training set perfectly in an over-parameterized setting, how does gradient descent find one that generalizes well to unseen data when such a model exists? This property is often called the *implicit bias* of gradient descent.²

¹No less an authority than von Neumann is reported to have said, “With four parameters I can fit an elephant, and with five I can make him wiggle his trunk” [Mayer et al., 2010].

²It is an implicit bias because even in the absence of explicit regularizers (such as weight decay or drop out), gradient descent still generalizes well (where possible).

Most previous attempts to make progress on these questions have been based on the notion of uniform convergence [Vapnik and Chervonenkis, 1971], the primary theoretical tool in classical learning theory. However, a recent paper by Nagarajan and Kolter [2019b] provides a strong argument for why any method based on uniform convergence is unlikely to provide an explanation to the mystery.

In this work, we study the generalization mystery from the less-explored, but arguably more general, perspective of algorithmic stability [Devroye and Wagner, 1979, Bousquet and Elisseeff, 2002, Hardt et al., 2016]. We propose that the answer to both the questions lies in the interaction of the gradients of the examples during training. (The interaction of per-example gradients has not received much attention in the literature, with notable exceptions being Yin et al. [2018], Fort et al. [2020], Sankararaman et al. [2020], Liu et al. [2020c], He and Su [2020], Mehta et al. [2021] which we discuss later.) Specifically, we argue,

- A1. Gradient descent in an over-parameterized setting generalizes well when the gradients of different examples (during training) are similar, that is, when there is *coherence*.**
- A2. When there is coherence, the dynamics of gradient descent leads to models that are *stable*, that is, to models that do not depend too much on any one training example, and, as is well known, stable models generalize well.**

In this paper, we advance a theory along these lines, which we call the theory of **Coherent Gradients**.³

2. THE THEORY, INFORMALLY

The motivation for our theory comes from the observation that the ability to memorize random datasets, and yet generalize well on real datasets is not unique to deep neural networks, but also seen, for example, with decision trees (and random forests). But there is no generalization mystery there: a typical tree construction algorithm splits the training set recursively into similar subsets based on input features. If no similarity is found, eventually, each example is put into its own leaf to achieve good training accuracy, but, of course, at the cost of poor generalization. Thus, trees that achieve good training accuracy on a randomized dataset are larger than those on a real dataset (for example, see Expt. 5 in [Chatterjee and Mishchenko, 2020]).

We propose that something similar happens with neural networks trained with gradient descent (GD):

Gradient descent exploits patterns common across training examples during the fitting process, and if there are no common patterns to exploit, then the examples are fitted on a “case-by-case” basis.

Intuitively, this provides a *uniform* explanation of memorization and generalization: If a dataset is such that examples are fitted on a case-by-case basis, then we expect poor generalization (it corresponds to memorization), whereas, if there are common patterns that can be exploited to fit the data, then we should expect good generalization.

So how does gradient descent exploit common patterns to fit the training data (when such common patterns exist)? Since the only interaction between examples in gradient descent is in the parameter update step, the mechanism for commonality exploitation—if it exists—must be there. Let η be the learning rate and let $g_i(w_t)$ be the gradient for the i th training example at the point w_t in parameter space. Furthermore, for now, assume that the $g_i(w_t)$ all have the same scale (we

³Some preliminary results appeared in Chatterjee [2020], Zielinski et al. [2020], Chatterjee and Zielinski [2020].

discuss the consequences of relaxing this shortly). Now, consider the parameter update at step t of gradient descent:⁴

$$w_{t+1} \equiv w_t - \eta \frac{1}{m} \sum_{i=1}^m g_i(w_t)$$

If $g(w_t)$ denotes the average gradient, that is, $g(w_t) \equiv \frac{1}{m} \sum_{i=1}^m g_i(w_t)$, we observe the following:

- (1) The average gradient $g(w_t)$ is stronger in directions (components) where the per-example gradients are “similar,” and reinforce each other; and weaker in other directions where they are different, and do not add up (or perhaps cancel each other).
- (2) Since network parameters w_{t+1} are updated proportionally to gradients, those parameters that correspond to stronger gradient directions change more.

Therefore, the parameter changes of the network are biased towards those that benefit multiple examples.

When per-example gradients reinforce each other, we say they are “coherent,” and use the term *coherence* to informally refer to the similarity of per-example gradients (either in aggregate across entire gradients or across specific components of the gradients).

Note that it is possible that there are *no* directions or components where different per-example gradients add up, that is, there is no coherence. That case would correspond to fitting each example independently, that is, on a case-by-case basis. In other words, reducing loss on one example fails to reduce the loss on any other example. Intuitively, one might expect this to be the case for random datasets, and only possible when the learning problem is over-parameterized: the dimensionality of the tangent space (that is, the dimension of the gradient vectors) is greater than the number of training examples.

We can reason about the generalization of the above process through the notion of algorithmic stability [Bousquet and Elisseeff, 2002, Devroye and Wagner, 1979]. A learning process is *stable* if a replacing one example in the training sample by another example (from the example distribution) does not change the learned model too much. It can be shown that models obtained through stable processes generalize well.

Strong directions in the average gradient are stable since in those directions multiple examples support or reinforce each other. In particular, the absence or presence of a single example in the training set does not impact descent down a strong direction since other examples contribute to it anyway. Therefore, by stability theory, the corresponding parameter updates should generalize well, that is, they would lead to lower loss on unseen examples as well.

On the other hand, weak directions in the average gradient are unstable, since they are due to a few or even single examples. In the latter case, for example, the absence of the corresponding example in the training set would prevent descent down that direction, and therefore, the corresponding parameter updates would not generalize well.

With this observation, we can reason inductively about the stability of GD: since the initial values of the parameters do not depend on the training data, the initial function mapping examples to their gradients is stable. Now, if all parameter updates are due to strong gradient directions, then stability is preserved. However, if some parameter updates are due to weak gradient directions, then stability is diminished. Now, since stability (suitably formalized) is equivalent to generalization [Shalev-Shwartz et al., 2010], this allows us to see how generalization may degrade as training

⁴For simplicity of exposition, here we only consider the full-batch case and no batch normalization. For the stochastic case, we have to consider the expected gradient, but the argument carries over, as we shall see shortly in the formal development.

progresses.⁵ In summary:

When there is coherence, the dynamics of gradient descent, particularly the use of the average gradient, leads to models that are stable, that is, to models that do not depend too much on any one training example.

Of course, in general, there can be *no* dataset-independent guarantee of stability. We can understand the connection between coherence of the dataset and stability intuitively, by considering two extreme cases. If, on the one hand, the gradients of all the examples are pointing in the same direction (that is, we have perfect coherence), then the replacing one training example by another does not matter, and the loss on the replaced example should still decrease, even though it was not used for training. On the other hand, if the gradients of all the examples are all orthogonal to each other (we have low coherence), then replacing an example with another eliminates the descent down the gradient direction of the replaced example, and we should not expect loss to reduce on that example. In other words, in the presence of high coherence, training on the original training set, and a perturbed version should not differ too much.

This insight suggests a new modification to gradient descent to improve generalization, and also explains some existing regularization techniques:

- We can make gradient descent more stable by eliminating or attenuating weak directions by combining per-example gradients using robust mean estimation techniques (such as median) instead of simple averages.
- We can view ℓ^2 regularization (weight decay) as an attempt to eliminate movement in weak gradient directions by having a “background force” that pushes each parameter to a data-independent default value (zero). It is only in the case of strong directions that this background force is overcome, and parameters updated in a data-dependent manner.
- Earlier we assumed that all the $g_i(w_t)$ all have the same scale. But that is not always true. For example, during training, some examples get fitted earlier than others, and so their gradients become negligible. Now, fewer examples dominate the average gradient $g(w_t)$, and this leads to overfitting since stability is degraded. This may be viewed as a justification for early stopping as a regularizer.

The properties of gradient descent as an optimizer do not play an important role in our generalization theory. We model gradient descent as a simple, “combinatorial” optimizer, a kind of greedy search with some hill-climbing thrown in (due to sampling and finite step size). Therefore, we are concerned less about the quality of solutions reached at the end of gradient descent, but more about staying “feasible” at all times during the search. In our context, feasibility means being able to generalize; and this naturally leads us to look at the transition dynamics of gradient descent to see if that preserves generalizability.

A notable feature of this approach is that it does not rely on any special property of the final solution obtained by gradient descent. Therefore, (1) it allows us to reason in an *uniform* manner about generalization with respect to early stopping and generalization at convergence,⁶ and (2) it

⁵Based on this insight, we shall see later how a simple modification to GD to suppress the weak gradient directions can dramatically reduce overfitting.

⁶Since stopping early typically leads to smaller generalization gap, the nature of the solutions of GD (for example, stationary points, the limit cycles of SGD at equilibrium) cannot be the ultimate explanation for generalization. In fact, think of the extreme case where no GD step is taken, and we have zero generalization gap!

applies *uniformly* to convex and non-convex problems (and uniformly across architectures without requiring simplifying assumptions such as homogeneous activations, infinite widths, etc.).

3. AN ILLUSTRATIVE EXAMPLE

The main ideas of the theory described above can be illustrated with a simple thought experiment where we fit an over-parameterized linear model to two toy datasets, one of which is “real” and the other is “random.”

Example 1 (“Real” and “Random”). Consider the task of fitting the linear model

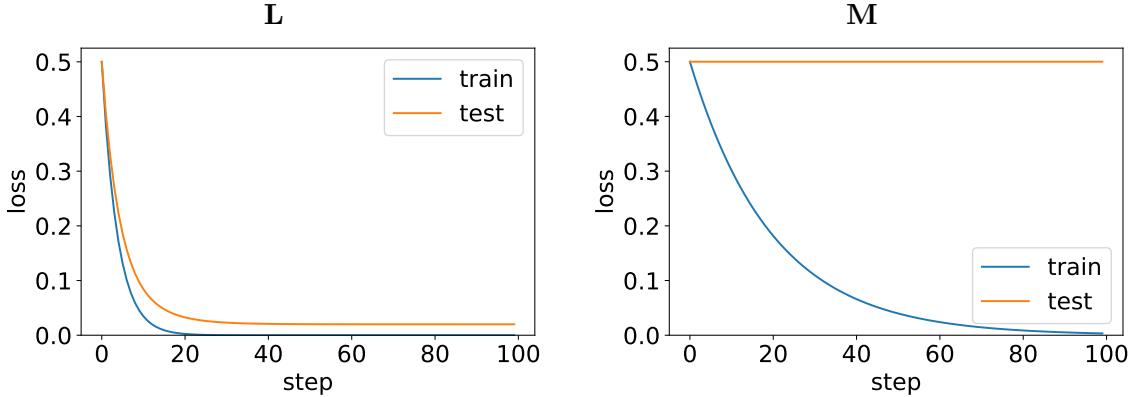
$$y = w \cdot x = \sum_{i=1}^6 w^{(i)} x^{(i)}$$

under the usual square loss $\ell(w) \equiv \frac{1}{2}(y - w \cdot x)^2$ to each of the following datasets:

L (“real”)			M (“random”)		
i	x_i	y_i	i	x_i	y_i
1	$\langle 1, 0, 0, 0, 0, 1 \rangle$	1	1	$\langle 1, 0, 0, 0, 0, 0 \rangle$	1
2	$\langle 0, -1, 0, 0, 0, -1 \rangle$	-1	2	$\langle 0, -1, 0, 0, 0, 0 \rangle$	-1
3	$\langle 0, 0, -1, 0, 0, -1 \rangle$	-1	3	$\langle 0, 0, -1, 0, 0, 0 \rangle$	-1
4	$\langle 0, 0, 0, 1, 0, 1 \rangle$	1	4	$\langle 0, 0, 0, 1, 0, 0 \rangle$	1
5	$\langle 0, 0, 0, 0, -1, -1 \rangle$	-1	5	$\langle 0, 0, 0, 0, -1, 0 \rangle$	-1

In each dataset, the first four examples, that is, (x_i, y_i) for $i \in [4]$ are used for training whereas the fifth example (x_5, y_5) is held out for evaluating generalization. Observe that the first 5 input features (shown in pink) are the same in both **L** and **M**, and furthermore, that they are each only predictive for one example (they are “idiosyncratic”). The last feature however is different in the two datasets. In **L** (shown in blue), it is predictive across all examples (“common” feature), whereas in **M** (shown in black), it is uninformative. We can think of **L** as a “real” dataset where there is something that can be learned by a linear model, whereas **M** is a “random” dataset that may at best be memorized.

Now, if we apply gradient descent (GD) to these problems (starting from $w = 0$), we see the following training and test curves:



In this simple setup, we see the same phenomena as in deep learning: (1) the model has sufficient capacity to memorize random data (**M**), yet it generalizes on real data (**L**), that is, the generalization is dataset dependent; and (2) from among all the models that fit **L** (which includes the solution for **M**), GD finds one that generalizes well (instead of just memorizing). However, in contrast to

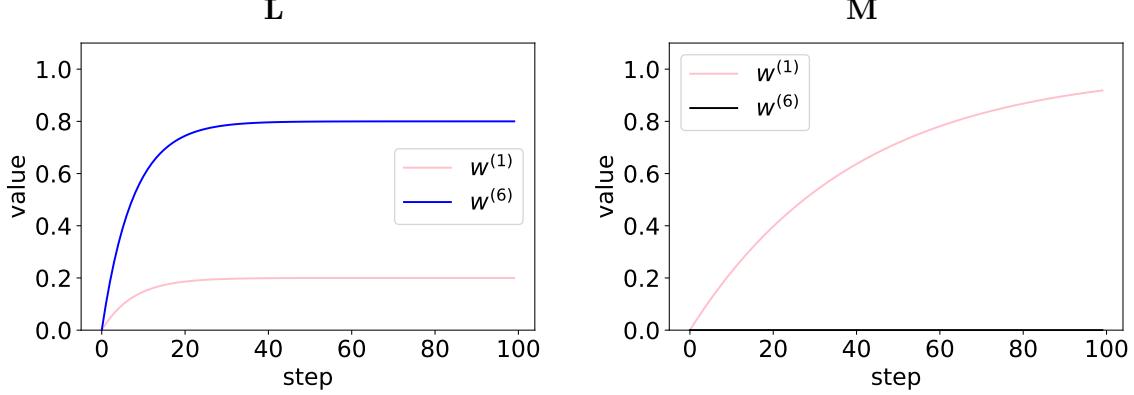
deep models, it is much easier to understand intuitively what is going on in this setup, particularly, from the point of view of the dynamics of GD.

First, consider the per-example gradients $g_i(w)$ and the average training gradient $g(w)$ at a point w on the trajectory of gradient descent starting from 0:

L		M	
i	$g_i(w)$	i	$g_i(w)$
1	$r(w) \langle 1, 0, 0, 0, 0, 1 \rangle$	1	$r(w) \langle 1, 0, 0, 0, 0, 0 \rangle$
2	$r(w) \langle 0, 1, 0, 0, 0, 1 \rangle$	2	$r(w) \langle 0, 1, 0, 0, 0, 0 \rangle$
3	$r(w) \langle 0, 0, 1, 0, 0, 1 \rangle$	3	$r(w) \langle 0, 0, 1, 0, 0, 0 \rangle$
4	$r(w) \langle 0, 0, 0, 1, 0, 1 \rangle$	4	$r(w) \langle 0, 0, 0, 1, 0, 0 \rangle$
$g(w)$	$r(w) \langle \frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4}, 0, 1 \rangle$	$g(w)$	$r(w) \langle \frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4}, 0, 0 \rangle$

where $r(w)$ is a scalar.⁷ Observe that **M** lacks coherence: the per-example gradients do not reinforce each other, that is, they do not add up in any component. Therefore, there are no strong directions in the average gradient. In contrast, in **L**, the per-example gradients are coherent: they share a *common* component (shown in blue) which “adds up” in the average gradient $g(w)$, which consequently, is **4 times** stronger in that component than in the other (idiosyncratic) components (shown in pink).

Finally, since parameter changes in GD are proportional to the average gradient $g(w)$, the parameter $w^{(6)}$ corresponding to the common input feature in **L** changes much more than any of the ones for the idiosyncratic features (for example, $w^{(1)}$)⁸ during training. In contrast, in **M**, only the weights for the idiosyncratic features are used to fit the data. This is confirmed by looking at the trajectories of $w^{(1)}$ and $w^{(6)}$ during training:



To summarize, in this simple example of “real” and “random” datasets that replicates the generalization mystery of deep learning, we see that:

- (1) The difference in the generalization gap between the two datasets can be understood in terms of the difference in the similarity of the per-example gradients, that is, in terms of the difference in coherence.
- (2) In the case of the “random” dataset, the training data was fit on a case-by-case basis (that is, memorized). Although the same would also have been an optimal solution for the

⁷ $r(w)$ is equal to $|y_i - w \cdot x_i|$ for $i \in [4]$, a quantity that is independent of i for w in the trajectory of GD due to the symmetries in this problem.

⁸From symmetry, $w^{(2)}, w^{(3)}$, and $w^{(4)}$ are the same as $w^{(1)}$.

problem of fitting the “real” training data, gradient descent produced a different solution by exploiting the commonality between training examples, as expressed in their gradients.

Finally, let us consider a modification to gradient descent where instead of simply averaging the per-example gradients, we take the component-wise **median** gradient $\tilde{g}(w)$:

L	M
<i>i</i>	<i>i</i>
1	1
2	2
3	3
4	4
$\tilde{g}(w)$	$\tilde{g}(w)$
$r(w) \langle 1, 0, 0, 0, 0, 1 \rangle$	$r(w) \langle 1, 0, 0, 0, 0, 0 \rangle$
$r(w) \langle 0, 1, 0, 0, 0, 1 \rangle$	$r(w) \langle 0, 1, 0, 0, 0, 0 \rangle$
$r(w) \langle 0, 0, 1, 0, 0, 1 \rangle$	$r(w) \langle 0, 0, 1, 0, 0, 0 \rangle$
$r(w) \langle 0, 0, 0, 1, 0, 1 \rangle$	$r(w) \langle 0, 0, 0, 1, 0, 0 \rangle$
$\tilde{g}(w)$	$\tilde{g}(w)$
$r(w) \langle 0, 0, 0, 0, 0, 1 \rangle$	$r(w) \langle 0, 0, 0, 0, 0, 0 \rangle$

It is easy to see that performing gradient descent with the median gradient $\tilde{g}(w)$ increases the stability of gradient descent on both datasets by eliminating the weak gradient directions, and leads to zero generalization gap for both “real” and “random”.⁹

□

The conventional explanation of generalization in over-parameterized linear models depends on the observation that GD from $\mathbf{0}$ finds the solution with minimum ℓ^2 norm (that is, the solution found by the pseudo-inverse method). In comparison, the explanation above, based on the alignment of per-example gradients and stability, is simpler and more direct. Crucially, it generalizes in a straightforward manner to deep networks where GD does not always find the minimum ℓ^2 norm solution. Thus, at a fundamental level, the Coherent Gradients approach allows us to decouple optimization and generalization; so even if we cannot say anything about the result of the optimization process after GD, we can say something about whether the solution generalizes or not by analyzing the per-example gradients along the way.

4. METRICS TO QUANTIFY COHERENCE

So far we have argued, informally, that if the gradients of different training examples are similar and reinforce each other (that is, if they are coherent), then the model produced by gradient descent is expected to generalize well; and we have illustrated this argument with a simple thought experiment. But in order to test this explanation in practical settings, such as the experiments of Zhang et al. [2017], we need to quantify the notion of coherence.

An obvious metric to quantify the coherence of the per-example gradients is their average pairwise dot product. Since this has a nice connection to the loss function, we start by reviewing the connection, and also set up some notation in the process. Formally, let $\mathcal{D}(z)$ denote the distribution of examples z from a finite set Z .¹⁰ For a network with d trainable parameters, let $\ell_z(w)$ be the loss for an example $z \sim \mathcal{D}$ for a parameter vector $w \in \mathbb{R}^d$. For the learning problem, we are interested in minimizing the expected loss $\ell(w) \equiv \mathbb{E}_{z \sim \mathcal{D}}[\ell_z(w)]$. Let $g_z \equiv [\nabla \ell_z](w)$ denote the gradient of the loss on example z , and $g \equiv [\nabla \ell](w)$ denote the average gradient. From linearity, we have,

$$g = \mathbb{E}_{z \sim \mathcal{D}} [g_z]$$

⁹Note that in the case of “random,” this reduction in the gap is achieved by preventing the training set from being memorized.

¹⁰We assume finiteness for mathematical simplicity since it does not affect generality for practical applications.

Now, suppose we take a small descent step $h = -\eta g$ (where $\eta > 0$ is the learning rate). From the Taylor expansion of ℓ around w , we have,

$$\ell(w + h) - \ell(w) \approx g \cdot h = -\eta g \cdot g = -\eta \mathbb{E}_{z \sim \mathcal{D}} [g_z] \cdot \mathbb{E}_{z \sim \mathcal{D}} [g_z] = -\eta \mathbb{E}_{z \sim \mathcal{D}, z' \sim \mathcal{D}} [g_z \cdot g_{z'}] \quad (1)$$

where the last equality can be checked with a direct computation. Thus, the expected pairwise dot product is equal to the change in loss divided by the learning rate, that is, the “instantaneous” change in loss.

Example 2 (Perfect similarity v/s pairwise orthogonal). Consider a sample with m examples z_i where $1 \leq i \leq m$. Let g_i be the gradient of z_i and further that $\|g_i\| = \|u\|$ for some u . Let $1 \leq j \leq m$. If all the g_i are the same, that is, if coherence or similarity is maximum, then $\mathbb{E}[g_i \cdot g_j] = \|u\|^2$. However, if they are pairwise orthogonal, i.e., $g_i \cdot g_j = 0$ for $i \neq j$, then $\mathbb{E}[g_i \cdot g_j] = (1/m)\|u\|^2$.

Observe that the loss reduces m times faster in the case of maximum coherence than when the gradients are pairwise orthogonal. \square

As the above example illustrates, the average expected dot product can vary significantly depending on the similarity of the gradients, and so could be used as a measure of coherence. However, since it has no natural scale—just rescaling the loss changes its value—it is difficult to interpret in an experimental setting. For example, it is not immediately clear if say, a value of 17.5 for the expected dot product indicates good or bad coherence.

Now, there is a natural scaling factor that can be used to normalize the expected dot product of per-example gradients. Consider the Taylor expansion of each individual loss ℓ_z around w when we take a small step h_z down its gradient g_z :

$$\ell_z(w + h_z) - \ell_z(w) \approx g_z \cdot h_z = -\eta g_z \cdot g_z$$

Taking expectations over z we get,

$$\mathbb{E}_{z \sim \mathcal{D}} [\ell_z(w + h_z) - \ell_z(w)] = -\eta \mathbb{E}_{z \sim \mathcal{D}} [g_z \cdot g_z] \quad (2)$$

The quantity in (2) has a simple interpretation: It is the expected reduction in the per-example loss ℓ_z if we could take different steps for different examples. In that sense, it is an *idealized reduction in loss* that real gradient descent cannot usually attain. As might be expected intuitively, it is a bound on the quantity in (1) and is tight when all the per-example gradients are identical. Thus, it serves as a natural scaling factor for the expected dot product, and we obtain a normalized metric for coherence, called α , from (1) and (2):

$$\alpha \equiv \frac{\ell(w + h) - \ell(w)}{\mathbb{E}_{z \sim \mathcal{D}} [\ell_z(w + h_z) - \ell_z(w)]} = \frac{\mathbb{E}_{z \sim \mathcal{D}, z' \sim \mathcal{D}} [g_z \cdot g_{z'}]}{\mathbb{E}_{z \sim \mathcal{D}} [g_z \cdot g_z]} = \frac{\mathbb{E}_{z \sim \mathcal{D}} [g_z] \cdot \mathbb{E}_{z \sim \mathcal{D}} [g_z]}{\mathbb{E}_{z \sim \mathcal{D}} [g_z \cdot g_z]} \quad (3)$$

From the discussion above, it is easy to see that

$$0 \leq \alpha \leq 1.$$

Coherence is 1 when all per-example gradients are identical, and 0 when the expected gradient is $\mathbf{0}$, that is, when training converges.¹¹ This is formalized as Theorem 2 in Appendix A.

¹¹If all examples are fit at the end of training, the denominator vanishes. However, in that case, the numerator is also zero, and we define the coherence to be 0. This is also consistent with adding a small positive epsilon to the denominator when computing coherence to avoid division by zero. This choice can also be justified with a continuity argument (see Appendix A).

Example 3 (Orthogonal Sample and Orthogonal Limit). If we have a sample of m examples whose gradients g_i ($1 \leq i \leq m$) are pairwise orthogonal (an “orthogonal” sample), then,

$$\alpha = \frac{(1/m) \mathbb{E}_i [g_i \cdot g_i]}{\mathbb{E}_i [g_i \cdot g_i]} = \frac{1}{m}$$

Thus, for an orthogonal sample, α is independent of the actual gradients, and we call $1/m$ the *orthogonal limit* for a sample of size m , and denote it by α_m^\perp . \square

In our experiments, we measure coherence on a sample, as we typically do not have access to the underlying distribution \mathcal{D} . However, as one might expect, α as measured from a sample is a biased estimator of the true (that is, the distributional) α . To distinguish between the two, we use α_S to denote the estimate of α obtained from a sample S . Since the sample S is often clear from the context, we commonly also use α_m instead of α_S where m is the size of S (to preserve the distinction with α and to remind ourselves of the sample size dependence).¹²

In the case of α estimated from a sample, the coherence of a (possibly fictitious¹³) orthogonal sample provides a convenient yardstick to judge the coherence of any other sample of the same size. So given a sample of size m , rather than show α_m , we show the ratio $\alpha_m/\alpha_m^\perp = m\alpha$. This quantity has a simple intuitive interpretation:

The metric α_m/α_m^\perp for coherence measures how many examples on average (including itself) a given example helps at a given step of gradient descent.

This intuition may be justified by considering different kinds of possible training samples:

- For an orthogonal sample, $\alpha_m = 1/m$ and, therefore, $\alpha_m/\alpha_m^\perp = 1$. In this case, each example in the sample is fitted independently of the others: a step in a direction of any given example’s gradient does not affect the loss on any other example. Each example “helps” only itself and does not positively or negatively affect the other examples in the sample.
- For a sample with perfect coherence (all the examples are identical), $\alpha_m = 1$, and we have, $\alpha/\alpha_m^\perp = m$. Here, each example in the sample “helps” all the other examples in the sample during gradient descent.
- At the end of training, when the expected gradient of the sample is $\mathbf{0}$, either the loss on a given example cannot be improved, or improving it comes at the expense of worsening the loss on other examples in the sample. Thus, an example “helps” no examples (including itself) on average, and $\alpha_m = \alpha_m/\alpha_m^\perp = 0$.
- For the $d \ll m$ (under-parameterized) case, consider α_m/α_m^\perp for a sample that comprises $k \gg 1$ copies of d orthogonal gradients living in a d -dimensional tangent space (that is, $m = kd \gg d$). Now, α_m of this replicated sample is also $1/d$ (since replicating a sample does not change the empirical distribution) and thus, its α_m/α_m^\perp is $(1/d)/(1/kd) = k$, which agrees with the intuition that each example in the replicated sample helps k other examples.

Example 4 (“Real” and “Random”). In Example 1, consider the training sample for \mathbf{L} ($m = 4$). It is easy to check with a direct computation that $\alpha_m = 5/8$ and $\alpha_m/\alpha_m^\perp = 2.5$, that is, each example

¹²Typically, α_S overestimates α . A sample S of size 1 provides an extreme example of this, since in that case, α_S is 1 regardless of α . Another example is provided by a distribution \mathcal{D} where all the per-example gradients have the same norm, and it is not too hard to see that $\mathbb{E}_{S \sim \mathcal{D}^m}[\alpha_S] \geq \alpha$. This is also confirmed by experiments (for example, see Figure 11 in Appendix D).

¹³If the dimension d of the tangent space is less than m , that is, we have more examples than parameters, then an orthogonal sample of size m cannot actually be constructed.

helps 2.5 examples. Intuitively, 2.5 may be understood as an example helping reduce the loss on itself 100%, and reducing the loss on 3 other examples by 50% due to the common component (the other half comes from their own idiosyncratic components).

The training sample for \mathbf{M} is pairwise orthogonal, and therefore, $\alpha_m/\alpha_m^\perp = 1$. \square

One can imagine many different metrics for coherence. We have seen two so far: the expected pairwise dot product and α . As we will see in the next section, another metric that arises naturally is

$$\mathbb{E}_{z \sim \mathcal{D}, z' \sim \mathcal{D}} [\|g_z - g_{z'}\|], \quad (4)$$

and there are other metrics in the literature such as stiffness [Fort et al., 2020] and GSNR [Liu et al., 2020c]. Furthermore, coherence metrics may be computed at the level of the network (that is, with entire per-example gradient vectors), or at the level of layers or even individual parameters (that is, with different projections of the per-example gradients).

Compared to the other metrics, α (and by extension, α_m/α_m^\perp) has certain advantages: It is simultaneously (1) easy to compute at scale, (2) more interpretable, and (3) has convenient theoretical properties that, among other things, allows us to provide a qualitative bound on the generalization gap. However, α is not a perfect metric and has some significant limitations as we shall see.

5. BOUNDING THE GENERALIZATION GAP WITH α

We can now formalize the argument in Section 2 by using α as a metric of coherence. We build on the work of Hardt et al. [2016] and Kuzborskij and Lampert [2018] who showed that (small-batch) *stochastic* gradient descent is stable since each training example is looked at so rarely, that it cannot have much influence on the final model if the training is not run for too long (or the learning rate is decayed quickly enough). Obviously, such a dataset-independent argument for stability is inapplicable to the generalization mystery since it rules out memorization.¹⁴ In contrast, we provide a dataset-dependent argument for stability may be summarized as follows:

$$\text{coherence} \implies \text{stability} \implies \text{generalization}$$

In the reverse direction, although it is known that generalization implies stability, generalization or stability do not imply coherence. There may be good generalization in spite of low coherence simply by virtue of having many training examples (for example, if the learning problem is under-parameterized) or by training for a short time.

In this context, our main result is a bound on the expected generalization gap, that is, the expected difference between training and test loss over all samples of size m from \mathcal{D} (denoted by $\text{gap}(\mathcal{D}, m)$) in terms of α . In its most general form, where we make no assumption on the learning rate schedule, it is as follows:

Theorem 1 (Generalization Theorem). *If (stochastic) gradient descent is run for T steps on a training set consisting of m examples drawn from a distribution \mathcal{D} , we have,*

$$|\text{gap}(\mathcal{D}, m)| \leq \frac{L^2}{m} \sum_{t=1}^T [\eta_k \beta]_{k=t+1}^T \cdot \eta_t \cdot \sqrt{2(1 - \alpha(w_{t-1}))} \quad (5)$$

¹⁴In the light of this, the experiments of Zhang et al. [2017] may be seen as demonstrating that in practice we run SGD long enough that it is no longer (unconditionally) stable.

where $\alpha(w)$ denotes the coherence at a point w in the parameter space, w_t the parameter values seen during gradient descent, η_t is the step size at step t ,

$$[\eta_k \beta]_{k=t_0}^{t_1} \equiv \prod_{k=t_0}^{t_1} (1 + \eta_k \beta),$$

and L and β are certain Lipschitz constants.

Proof. Please see Appendix C for the full formal statement and the proof. \square

Like Hardt et al. [2016], our bound depends on the length of training (shorter the training, better the generalization), and the size of the training set (more the examples, better the generalization).¹⁵ However, unlike Hardt et al. [2016], the bound also depends on the coherence as measured by α during training (greater the coherence, better the generalization). Due to the presence of the “expansion” term $[\eta_k \beta]_{k=t+1}^T$, an interesting qualitative aspect of this dependence is that high coherence early on in training is better than high coherence later on. Also in contrast to Hardt et al. [2016], our bound applies in an uniform manner to the stochastic and the full-batch case in the general non-convex setting.¹⁶

We emphasize that as with most theoretical generalization bounds in deep learning, our bound is extremely loose, and is therefore only useful in a qualitative sense. This is not only due to the Lipschitz constants. The coherence term based on α only plays a strong role when it is close to 1, but as we shall see in the next section, α on real datasets is quite far from 1. This is because α , being an average over the entire network, is a rather blunt instrument, and conjecture that a tighter bound may be obtained in terms of layer-wise coherences, or perhaps better yet, in terms of “minimum coherence cuts” of the network.

The proof of the theorem follows the iterative framework introduced in Hardt et al. [2016]. We analyze the evolution of two models under stochastic gradient descent, one trained on the original training set, and another trained on a perturbed training set where the i th training example (z_i) is replaced with a different one from the data distribution (call it z'_i). When a minibatch with the i th example is encountered, the two models diverge (further) due to the different gradients, and this divergence (in expectation) is bounded by the quantity in (4). This quantity, in turn, is bounded by α as follows (see Lemma 6 in Appendix A):

$$\mathbb{E}_{z \sim \mathcal{D}, z' \sim \mathcal{D}} [\|g_z - g_{z'}\|] \leq \sqrt{2(1-\alpha) \mathbb{E}_{z \sim \mathcal{D}} [g_z \cdot g_z]}.$$

6. MEASURING α ON REAL AND RANDOM DATASETS

In this section, we use the α_m/α_m^\perp metric developed in Section 4 to investigate the generalization mystery. We can measure α_m at any point in training by computing it directly using (3) on a given sample.¹⁷ The sample could be either from the training set (giving us “training” coherence), or the test set (“test” coherence).

One complication is the use of batch normalization [Ioffe and Szegedy, 2015] in most practical networks, since in that case, per-example gradients are no longer well-defined. This is a problem not just with α , but with any metric that is based on per-example gradients. However, α has an interesting mathematical property that allows us to impute it using the coherence of *mini-batch*

¹⁵In fact, when m is large, even with random data, we can have good generalization since we are interested in the gap, and not just test loss—fitting the training set gets harder with larger m . Although natural, this inverse dependence on m does not usually hold for bounds based on uniform convergence as was pointed out by Nagarajan and Kolter [2019b].

¹⁶In fact, since our analysis is in expectation, the size of the minibatch drops out of the bound.

¹⁷Observe that α can be computed efficiently in an online fashion by keeping a running sum for $\mathbb{E}[g_z]$, and another for $\mathbb{E}[g_z \cdot g_z]$.

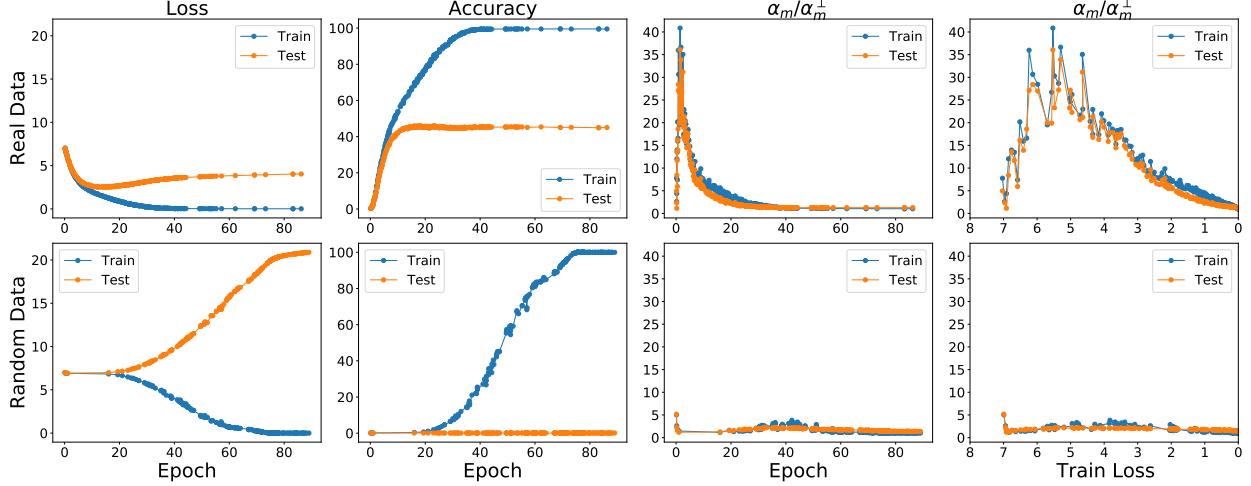


FIGURE 1. An experiment in the spirit of Zhang et al. [2017] illustrating the generalization mystery in deep learning. We train two ResNet-50 models, one on ImageNet with original labels (“real”, top row), and another on ImageNet with images replaced by Gaussian noise (“random”, bottom row) using vanilla SGD and no explicit regularization. As the loss and accuracy curves (first two columns) show, the network has sufficient capacity to memorize the training data, yet, it generalizes in one case and not the other. We believe that the reason for this difference in behavior can be found by analyzing the similarity between per-example gradients during training, that is, *coherence*. Using the α/α_m^\perp metric for coherence (last two columns), we see that in the case of real data, the per-example gradients are much more similar, and each example helps reduce the loss on many other examples, as compared to the random case.

gradients. We discuss this in greater detail in Appendix D. In what follows, we use this imputation method for networks with batch normalization.

In our main experiment, along the lines of Zhang et al. [2017], we train a ResNet-50 network on two datasets, one real (ImageNet) and the other random (“ImageNet” with random input pixels). Both networks are trained with vanilla SGD (that is, no momentum) with a constant learning rate of 0.5 and batch size of 4096. We do not use any explicit regularizers such as weight decay or dropout. We estimate test and train α during training using the test set (which has 50K samples) and a training sample (a random sample of 50K training examples). We take snapshots of the network during training each time the (minibatch) training loss falls by 1% from the previous low watermark, and compute loss, accuracy, and α_m/α_m^\perp on these snapshots using the test and training samples.

The resulting training curves are shown in the first two columns of Figure 1. As expected, both networks achieve zero training loss and near perfect training accuracy, but only the network trained on (real) ImageNet shows non-trivial generalization to the test set.

Coherence as measured by α_m/α_m^\perp ($m = 50,000$) is shown in the third and fourth columns of Figure 1. The third column shows coherence as a function of the number of epochs trained, whereas the fourth column shows it in terms of the training loss. Since the realized rate of learning is

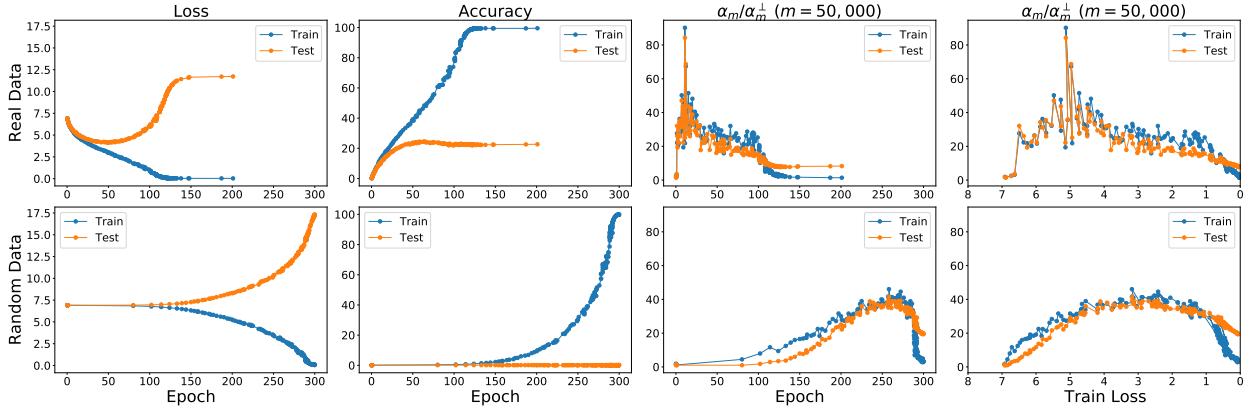


FIGURE 2. Unlike the situation with ResNet-50 (Figure 1), with AlexNet we find that the peak coherence for random data (second row) as measured by α_m/α_m^\perp can be surprisingly high, even though it happens much later in training, and is lower than that of real (first row). Although this appears to be a contradiction to the theory, it is not; it is a limitation of the metric. α_m/α_m^\perp in this plot is a measure of coherence over the entire network (that is, over entire per-example gradients), and is therefore an average quantity. A closer look at the layer-by-layer values of α_m/α_m^\perp as shown in Figure 3 reveals, once again, a significant difference between real and random data.

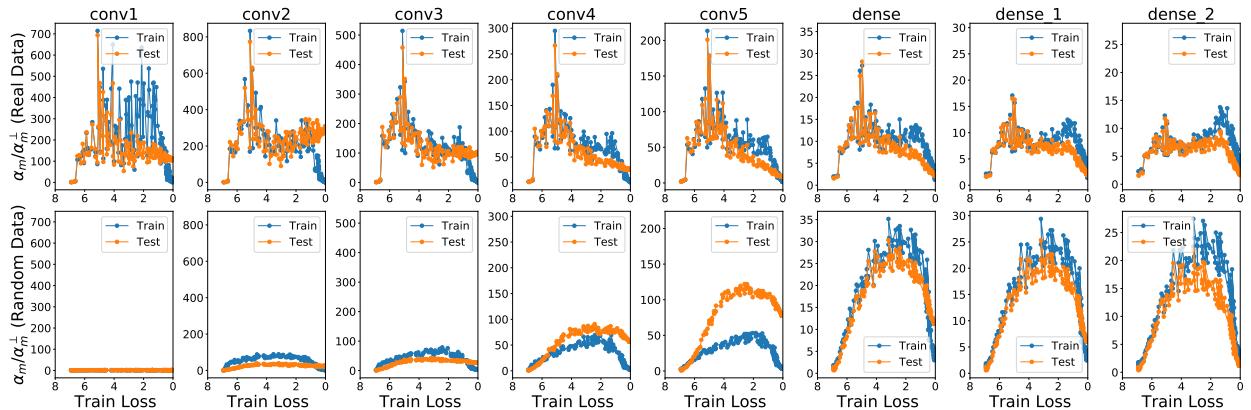


FIGURE 3. A layer-by-layer breakdown of α_m/α_m^\perp for AlexNet from Figure 2 shows that on random data (second row), α_m/α_m^\perp is indeed close to 1 and much lower than that of real data (first row) for the first few layers. For the higher (dense) layers, coherence is comparable between real and random, though note the difference in scale of α_m/α_m^\perp between the convolutional and dense layer plots.

different for the two datasets (the real data is learned in fewer epochs), plotting coherence against loss allows us to compare across the two datasets more easily.

In the case of real data, we observe that α_m/α_m^\perp starts off low (around 1) in early training and then increases to a maximum (about 40) within the first few epochs and then returns to a low value again (around 1) at the end of training.¹⁸ In the plot of α_m/α_m^\perp against training loss, we see that when actual learning happens, that is, when the loss comes down, α_m/α_m^\perp stays around 20. In other words, when training with real labels, each training example in our set of 50K examples used to measure coherence helps many other examples.

In contrast, for random data, although the evolution of α_m/α_m^\perp is similar to that of real data, the actual values, particularly, the peak is very different. α_m/α_m^\perp starts off low (around 1), increases slightly (staying usually below 5), and then returns back to a low value (around 1). Therefore, each training example in the case of random data, helps only one or two other examples during training, that is, the 50K random examples used to estimate coherence are more or less orthogonal to each other.¹⁹

In summary,

With a ResNet-50 model on real ImageNet data, in a sample of 50K examples, each example helps tens of other examples during training, whereas on random data, each example only helps one or two others.

This provides evidence that the difference in generalization between real and random stems from a difference in similarity between the per-example gradients in the two cases, that is, from a difference in coherence.

While experiments with other architectures and datasets also show similar differences between real and random datasets (see Appendix E), there are cases when the coherence of random data *as measured by α_m/α_m^\perp over the entire network* can be surprisingly high for an extended period during training.²⁰ In our experiments, we found an extreme case of this when we replaced the ResNet-50 network in the previous experiment with an AlexNet network (learning rate of 0.01). The training curves and measurements of α_m/α_m^\perp in this case are shown in Figure 2. As we can see, unlike the ResNet-50 case, α_m/α_m^\perp reaches a value of 40 for $m = 50,000$. In other words, in a sample of 50K examples, at peak coherence, each random example helps 40 other examples!²¹

What is going on? An examination of the per-layer values of α_m/α_m^\perp provides some insight. These are shown in Figure 3. We see that for the first convolution layer (`conv1`) in the case of random—and *only* in that case— α_m/α_m^\perp is approximately 1 indicating that the per-example gradients in that layer are pairwise orthogonal (at least over the sample used to measure coherence).²² This indicates that the first layer plays an important role in “memorizing” the random data since each example is pushing the parameters of the layer in a different direction (orthogonal to the rest). This is not surprising since the images are comprised of random pixels.

¹⁸For now, we ignore the small differences in training and test coherence.

¹⁹We note here that very early in training, that is, the first few steps (not shown in Figure 1, but presented in Figure 16 instead), α_m/α_m^\perp can be very high even for random data due to imperfect initialization. All the training examples are coordinated in moving the network to a more reasonable point in parameter space. As may be expected from our theory, this movement generalizes well: the test loss decreases in concert with training loss in this period. Rapid changes to the network early in training is well documented (see, for example, the need for learning rate warmup in He et al. [2016] and Goyal et al. [2017]).

²⁰As we discussed earlier, coherence even for random can be high for a short period early on in training due to imperfections in initialization. But the difference here is *sustained* high coherence.

²¹That said, note that (1) even in this case, at its peak α_m/α_m^\perp for real is more than $2\times$ the peak for random; and (2) the high coherence of random occurs much later in training than that of real which possibly indicates the importance of the “expansion term” ($[\eta_k \beta]_{k=t+1}^T$) in the bound of Theorem 1 (see discussion in Section 5).

²²The difference in α_m/α_m^\perp in the first layer between real and random is also seen when the entire training set is used to measure α_m/α_m^\perp (Figure 19).

Now, the overall (network) α is a convex combination of the per-layer α s (see Theorem 7 in Appendix A). Since the fully connected layers have high coherence, overall α (as shown in Figure 2) can be high even though there is a layer with very low α (at the orthogonal limit). In other words,

As a measure of coherence, α_m/α_m^\perp over the whole network, being an average, is a blunt instrument, and therefore, a finer-grained analysis, for example, on a per-layer basis, is sometimes necessary.

An important open problem, therefore, is to devise a better metric for coherence that accounts for the structure of the network, and to use that metric to obtain a bound stronger than that in Theorem 1. Please also see the discussion in Section 10, and particularly, Example 9 for a closer look at this in the context of a simple, illustrative example.

Evolution of Coherence. Experiments across several architectures and datasets show a common pattern in how coherence as measured by α_m/α_m^\perp (or equivalently, α) changes during training. Ignoring the initial transient in the first few steps of training, coherence follows a broad parabolic pattern: It starts off at a low value, rises to a peak, and then comes back down to the orthogonal limit.²³ This happens regardless of whether the dataset is random or real, indicating that this is an optimization (as opposed to a generalization) effect. We discuss the reasons behind this in Appendix F.

7. FROM MEASUREMENT TO CONTROL: SUPPRESSING WEAK DESCENT DIRECTIONS

Weak directions in the average gradient are supported by few examples or perhaps even one, and therefore, are less stable than strong directions which are supported by many examples. Therefore, as discussed in Section 2, suppressing weak directions in the average gradient should lead to less overfitting and better generalization. Now, although existing regularization techniques such as weight decay, dropout, and early stopping may be viewed through this lens, the theory also suggests a new, more direct, regularization technique we call *winsorized gradient descent* (WGD).

In WGD, instead of updating each parameter with the average gradient as in gradient descent, we update it with a “winsorized” average where the most extreme values (outliers) are clipped. Formally, let $w_t^{(j)}$ represent the j th trainable parameter (that is, the j th component of the parameter vector w_t) at step t , and $g_i^{(j)}(w_t)$ the j th component of the gradient of the i th example at w_t . In normal gradient descent, we update $w_t^{(j)}$ as follows:

$$w_{t+1}^{(j)} = w_t^{(j)} - \eta \frac{1}{m} \sum_{i=1}^m g_i^{(j)}(w_t).$$

Now, let $c \in [0, 50]$ be a hyperparameter that controls the level of winsorization. Define $l^{(j)}$ to be the c th percentile of $g_i^{(j)}(w_t)$ (over the examples i). Likewise, let $u^{(j)}$ be the $(100 - c)$ th percentile of $g_i^{(j)}(w_t)$. The update rule with winsorized gradient descent is as follows:

$$w_{t+1}^{(j)} = w_t^{(j)} - \eta \frac{1}{m} \sum_{i=1}^m \text{clip}(g_i^{(j)}(w_t), l^{(j)}, u^{(j)})$$

where $\text{clip}(x, l, u) \equiv \min(\max(x, l), u)$.

²³In rare cases, such as a fully connected network on MNIST where the signal is strong and easy to find, coherence starts off high.

The modified update rule minimizes the effect of outliers in the per-example gradients on a per-coordinate basis. The value of c dictates what is an outlier. When $c = 0$, nothing is an outlier, and this corresponds to normal gradient descent, whereas when $c = 50$, all values other than the median are considered outliers. Thus, increasing c reduces variance and increases bias.

Example 5 (WGD applied to Example 1). Recall that using the component-wise median gradient in Example 1 instead of the average gradient reduced the generalization gap to zero for both “real” and “random.” We note that this corresponds to running WGD with $c = 50$. □

Although the modification for WGD is a conceptually simple change, it greatly increases the computational expense due to the need to compute and store per-example gradients for all the examples. The computational expense can be reduced by performing winsorized *stochastic* gradient descent (WSGD). This is a straight forward modification of SGD where the winsorization is only performed over the examples in the mini-batch rather than all examples in the training set.

WSGD on MNIST. We train a small fully-connected ReLU network with 3 hidden layers of 256 neurons each for 60,000 steps (100 epochs) with a fixed learning rate of 0.1 on MNIST and 4 variants with different amounts of label noise ϵ ranging from 25% to 100%.²⁴ We use WSGD with a batch size of 100 and vary c in $\{0, 1, 2, 4, 8\}$. Since we have 100 examples in each minibatch, the value of c immediately tells us how many outliers are clipped in each minibatch. For example, $c = 2$ means the 2 largest and 2 lowest values of the per-example gradient in a batch are clipped (independently for each trainable parameter in the network), and as always, $c = 0$ corresponds to unmodified SGD.

The resulting training and test curves shown in Figure 4. The columns correspond to different amounts of label noise and the rows to different amounts of winsorization. In addition to the training and test accuracies (ta and va , respectively), we show the level of overfit which is defined as $ta - va$.

As expected, as c increases, and the weaker directions are suppressed more, the extent of overfitting decreases. Furthermore, for larger values of c (for example, $c = 8$) the ability to fit the corrupted labels is severely impacted. The training accuracy stays below the accuracy that would be obtained if only the uncorrupted labels were learned (shown by dashed gray lines in the plot). The ability to memorize, that is, fit random labels (100% noise) is impacted more than the ability to fit real labels (0% noise): the effective learning rate (the rate at which training accuracy increases) is much lower for random than for real.

In summary,

If we suppress weak gradient directions by modifying SGD to use a robust average of per-example gradients that excludes outliers, generalization improves. This provides further direct evidence that weak directions are responsible for overfitting and memorization.

Finally, we notice that with a large amount of winsorization, there can be optimization instability (not to be confused with algorithmic stability which as we have seen is improved): training accuracy can fall after a certain point. We are not sure what causes the instability but conjecture that this is

²⁴A label noise of 25% means that the dataset is constructed by randomly assigning labels to 25% of the training examples in MNIST that are chosen uniformly at random. The remaining 75% have the correct labels as do the test examples. Since MNIST has 10 possible labels, this means that $75\% + (1/10) 25\% = 77.5\%$ of training examples have pristine (that is, correct) labels and 22.5% have corrupted labels.

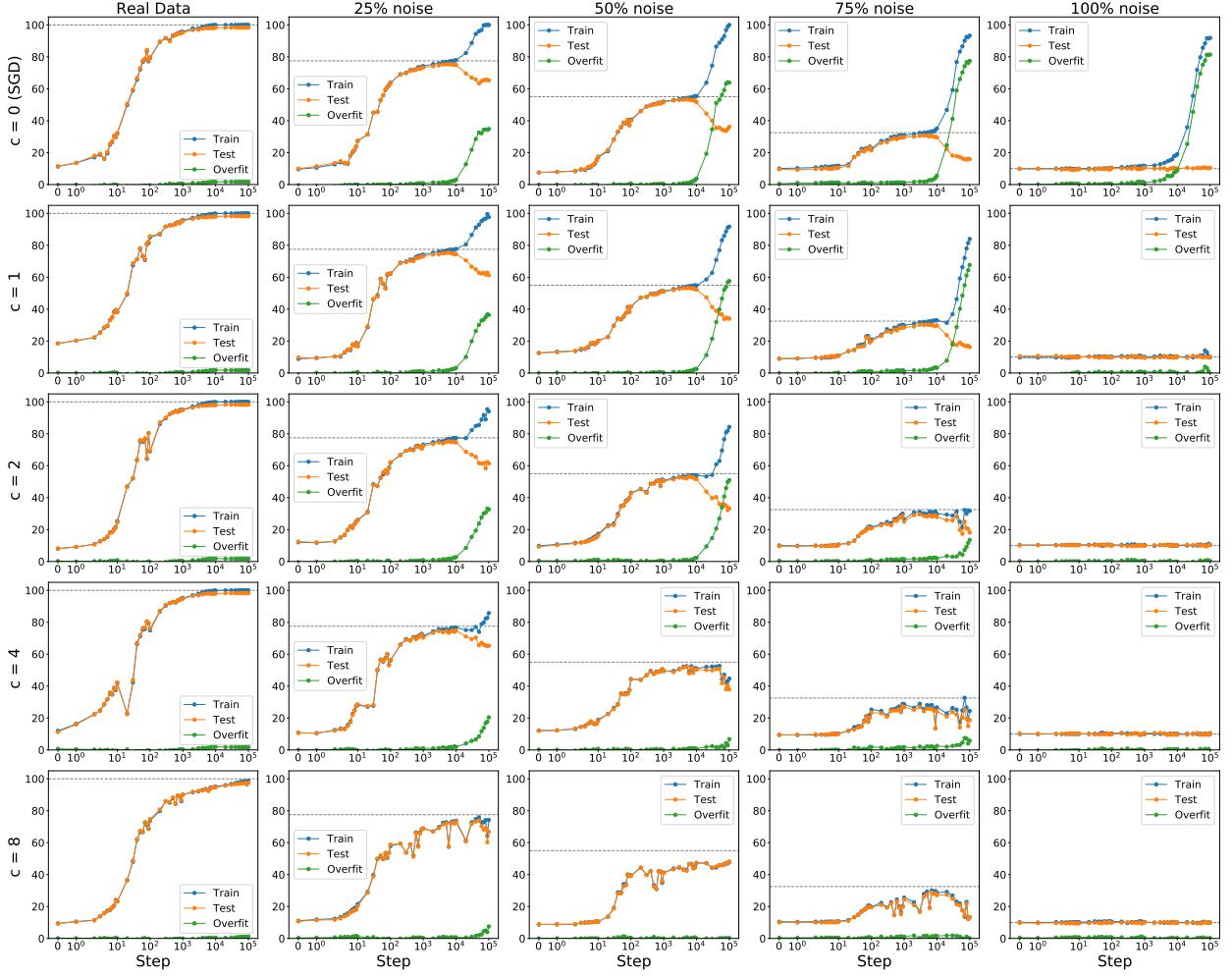


FIGURE 4. Generalization improves when weak directions in the average gradient are suppressed during gradient descent. Weak directions are suppressed by winsorization, that is, by clipping extreme per-examples gradients (independently for each coordinate of the gradient). The parameter c controls the level of winsorization and $c = 0$ corresponds to using the (usual) average gradient. We train a fully connected network on MNIST with varying amounts of label noise (see Figure 26 for a similar experiment with random pixels).

because of a strengthening of positive feedback loops in strong directions. Positive feedback loops are discussed in Section 10.

Scaling up. Winsorization requires computing and storing many per-example gradients.²⁵ This makes it impractical for large datasets such as ImageNet, as well as inapplicable to popular networks

²⁵For exact computation, the storage would depend on c and may be tractable for small c . This could be reduced further by employing approximations to compute streaming quantiles.

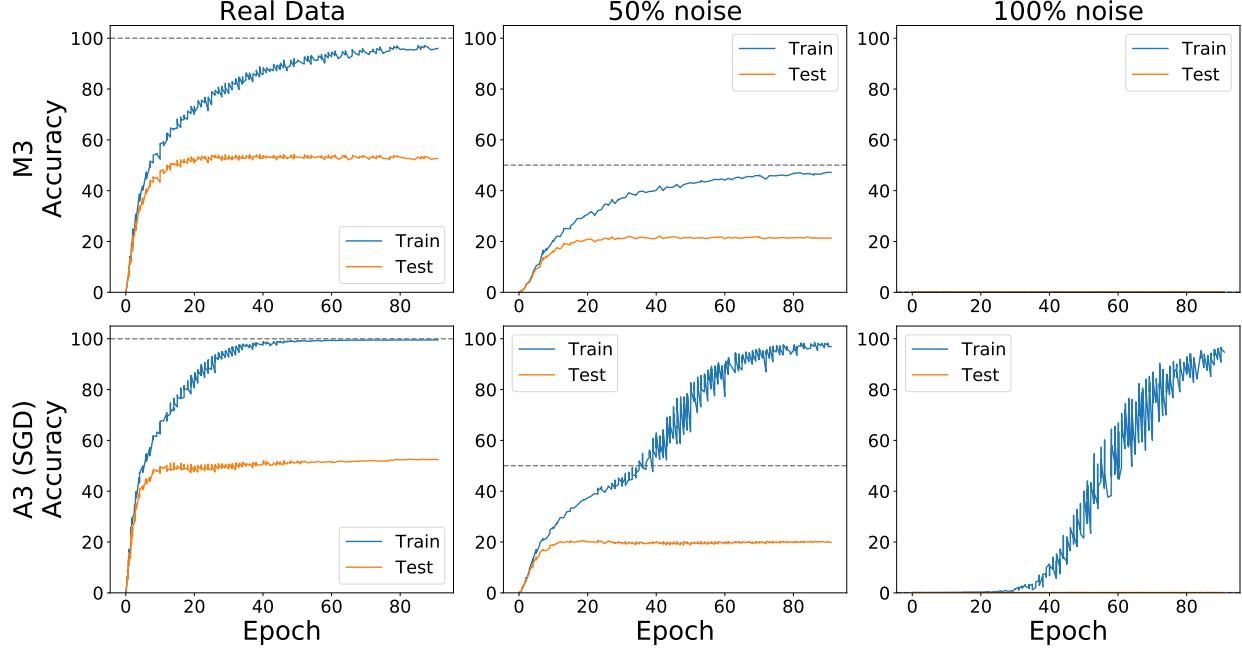


FIGURE 5. Taking the coordinate-wise median of 3 micro-batches (M3) suppresses weak gradient directions more than taking their coordinate-wise average (A3) (the same as ordinary SGD), leading to better generalization. Here, we train a ResNet-50 on 3 datasets derived from ImageNet by replacing some fraction of the training images with Gaussian noise. In the case of random data (100% noise), M3 prevents memorization, so the training and test curves both lie on the x-axis. In contrast, as expected, A3 (SGD) memorizes the training set. When only half the training images are replaced by noise (50% noise), M3 reaches a training accuracy near 50% suggesting that only the real images are learned. This is confirmed in Figure 6.

such as ResNet-50 which employ batch normalization as per-example gradients are not defined in that case.

An alternative to winsorization that addresses both these problems is the well known technique of taking median of means (see, for example, Minsker [2013]).²⁶ The main idea of the *median of means* algorithm is to divide the samples into k groups, computing the sample mean of each group, and then returning the median of these k means. The most obvious way to apply this idea to SGD is to divide a mini-batch into k groups of equal size. We compute the mean gradients of each group as usual, and then take their coordinate-wise median. The median is then used to update the weights of the network.²⁷

²⁶Median of means has the theoretical property that its deviation from the true mean is bounded above by $O(1/\sqrt{m})$ with high probability (where m is the number of samples). The sample mean satisfies this property only if the observations are Gaussian.

²⁷Since we are simply replacing the mini-batch gradient with a more robust alternative, this technique may be used in conjunction with other optimizers.

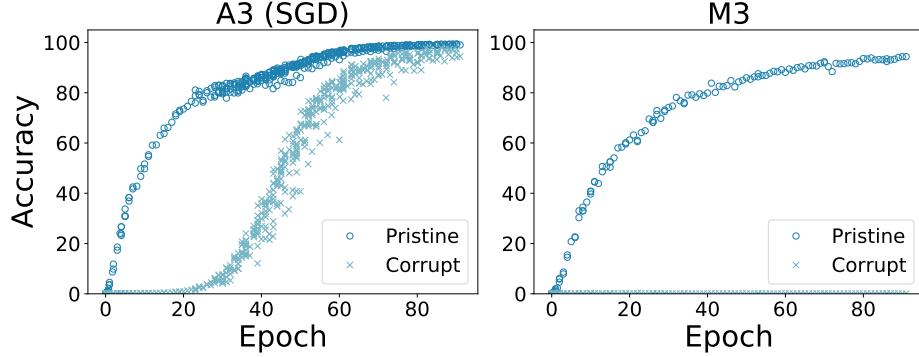


FIGURE 6. In the experiment of Figure 5, on the dataset where half the training images are replaced by noise (50% noise), M3 only learns the real (pristine) training images but does not learn the random (corrupt) images. In contrast, A3 (that is, SGD) learns both, though the corrupt examples are learned later in training. This is further evidence that the gradients of real data are well-aligned, and add up to strong directions in the average gradient; whereas those of random data are not well-aligned, and may be suppressed if weak directions are attenuated through robust averaging. This difference in the strength of the gradients in the two cases explains the observed difference in generalization between them.

Even though the algorithm is straightforward, its most efficient implementation (that is, where the k groups are large and processed in parallel) on modern hardware accelerators requires low-level changes to the stack to allow for a median-based aggregation instead of the mean. Therefore, in this work, we simply compute the mean gradient of each group as a separate *micro-batch* and only update the network weights with the median every k micro-batches, i.e., we process the groups serially.

In the serial implementation, $k = 3$ is a sweet spot. We have to remember only 2 previous micro-batches, and since $\text{median}(x_1, x_2, x_3) = \sum_i x_i - \min_i x_i - \max_i x_i$ (where $i \in \{1, 2, 3\}$), we can compute the median with simple operations. We call this M3 (median of 3 micro-batches).

Example 6 (M3 applied to Example 1). In this case, applying M3 with a micro-batch size of 1 (batch size of 3) completely suppresses the weak gradient directions (corresponding to the idiosyncratic signals) in both the “real” and “random” datasets, and leads to perfect generalization. □

Figure 5 shows the effect of M3 on 3 datasets: real ImageNet, ImageNet with half the images replaced by Gaussian noise (50% noise), and “ImageNet” with all the images replaced by Gaussian noise (100% noise). The 100% noise dataset is the same as the “random” dataset in Section 6. We run M3 with a micro-batch of 256 (that is, batch size of 3×256) for 90 epochs and a learning rate of 0.20. For comparison we also run SGD with an identical setup. Instead of the median of 3 micro-batches, in SGD we compute their average, and hence, in this context, we refer to SGD

as A3 to highlight that the *only* difference in the two cases is replacing the median operation with average.²⁸

We observe the following:

- (1) On real data, M3 only slightly slows down training after about 50% of the dataset has been learned, and almost reaches the 100% training accuracy achieved by A3. Also, M3 has a lower generalization gap (43.34%) compared to A3 (47.15%).
- (2) With 50% noise, M3 reaches a training accuracy less than 50% which is much lower than the 100% training accuracy of A3. Thus M3 has a much lower generalization gap (25.90%) compared to A3 (77.13%). Further investigation shows that M3 learns the real 50% of the training dataset (the “pristine” examples) well, but does not learn the remaining 50% of the training set that is Gaussian noise (the “corrupt” examples). This is, of course, in contrast to A3 which learns both almost equally well. See Figure 6.
- (3) Finally, in the case of 100% noise, M3 fails to learn any of the training data at all in sharp contrast to A3.

This provides further evidence that (a) weak directions are responsible for memorization, and (b) suppressing weak directions improves generalization (via improving stability).

The choice of hyper-parameters of M3 requires care. Larger the size of the micro-batch, less effective is the median operation in suppressing weak directions (consider the extreme case when the micro-batch is the entire training set and taking the median serves no purpose). However, smaller the micro-batch, less reliable are the batch statistics (for batch normalization), and higher the effective level of winsorization (consider a micro-batch size of 1 where the effective winsorization level now is $c = 50$). The learning rate has to be set in conjunction with the micro-batch size, and for some learning rates we found the training to be unstable (see Figure 27). This is similar to the optimization stability seen with winsorization previously.

8. WHY ARE SOME EXAMPLES (RELIABLY) LEARNED EARLIER?

In a study on memorization in shallow fully-connected networks and small convolutional networks on MNIST and CIFAR-10, Arpit et al. [2017] discovered that for real datasets, starting from different random initializations, many examples are consistently classified correctly or incorrectly after one epoch of training. They call these *easy* and *hard* examples respectively. They hypothesize that this variability of difficulty in real data “is because the easier examples are explained by some simple patterns, which are reliably learned within the first epoch of training.”

But that begs the question what makes a pattern “simple” and why are such patterns reliably learned early? The theory of Coherent Gradients provides a refinement of their hypothesis:

Easy examples are those examples that have a lot in common with other examples where commonality is measured by the alignment of the gradients of the examples.

Under this assumption, it is easy to see why an easy example is learned sooner reliably: most gradient steps benefit it.

Note that this hypothesis is more nuanced than the conjecture in Arpit et al. [2017]. First, the difficulty of an example is not simply a property of that example (whether it has simple patterns or

²⁸In practice, SGD with a batch size of 3×256 is subtly different from A3 with a micro-batch of 256 in how the low-level computation is distributed among the accelerators and in the resulting statistics for batch normalization. Although that difference is immaterial, we use A3 here to eliminate *all* differences other than the method of aggregation.

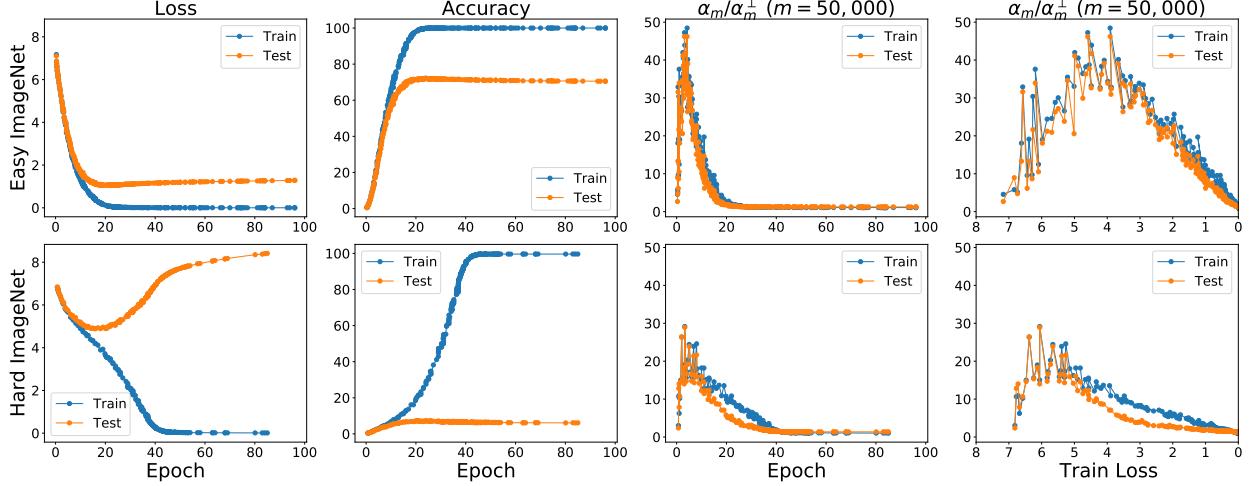


FIGURE 7. Coherence as measured by α_m/α_m^\perp is higher for examples that are learned early in ImageNet training (Easy ImageNet) than for those that are learned later (Hard ImageNet). As expected, the higher coherence translates into better generalization. The network used is a ResNet-50.

not), but depends on the relationship of that example to others in the training set (what it shares with other examples).

Second, the dynamics of training, including initialization, can determine the difficulty of examples. Consequently, it can accommodate the observed phenomenon of adversarial initialization [Liao et al., 2018, Liu et al., 2019] where examples that are easy to learn with random initialization become significantly harder to learn with a different, specially constructed initialization, and the generalization performance of the network suffers. Any notion of simplicity of patterns intrinsic to an example alone cannot explain adversarial initialization since the dataset remains the same, and therefore, so do the patterns in the data.

In order to test the above hypothesis, we create two datasets with 500K training and 50K test examples from the standard ImageNet training set. These datasets consists of examples that a ResNet-50 reliably learns early (“Easy ImageNet”) or late (“Hard ImageNet”). See Appendix G for more details on this construction. Next, as in Section 6, we measure loss, accuracy, and coherence during training on both datasets.

The results are shown in Figure 7. We observe that the coherence for Easy ImageNet is significantly higher than that of Hard ImageNet, and as might be expected from the theory, the generalization gap of Easy ImageNet is smaller since the average gradient in that case is more stable.²⁹

In other words, since easy examples, as a group have more in common with each other (than with the hard examples, or the hard examples have amongst themselves), (1) a model is less impacted by the presence or absence of a single easy example in the training set, leading to good generalization on easy examples; and (2) easy examples have a stronger (focused) effect on the average gradient

²⁹This is also seen in an *in situ* measurement of coherence of easy and hard examples during regular ImageNet training, but due to a technicality involving batch normalization statistics, that measurement is less reliable. See Appendix G for details.

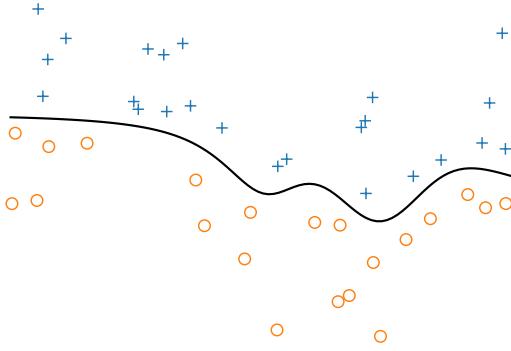


FIGURE 8. If gradient descent enumerates hypotheses of increasing complexity then the examples learned early, that is, the easy examples, should be the ones far away from the decision boundary since they can be separated by simpler hypotheses. However, a model learned from the hard examples then, should generalize well to easy examples, but that is not what we find on ImageNet. A model trained on hard ImageNet examples has only a 17% accuracy on easy ImageNet examples (in contrast to a 71% accuracy obtained by a model trained on other easy examples).

direction than the hard examples (which are diffuse) leading to easy examples being learned faster. Thus there is a correlation between how quickly examples are learned and their generalization.³⁰

Although the above experiment shows that easy examples have more in common with each other, it does not rule out that nonetheless there is something intrinsic to each easy example that accelerates learning. For example, it may be that easy examples have larger gradients than hard examples. To investigate this, we train a network on a *single* example from Easy ImageNet or Hard ImageNet at a time and measure the number of steps it takes to fit the example. We find that there is no statistically significant difference between the two datasets in this regard. It is only when sufficiently many easy examples come together (as in the earlier experiment), that they train faster than hard examples. This may be seen as further evidence that it is the relationship between examples rather than something intrinsic to an example that controls difficulty. Please see Appendix G for more details.

Does gradient descent explore functions of increasing complexity? One intuitive explanation of generalization is that GD somehow explores candidate hypotheses of increasing “complexity” during training. Although this is backed by some observational studies (for example, Arpit et al. [2017], Nakkiran et al. [2019], Rahaman et al. [2019]), we are not aware of any causal mechanism that has been proposed to explain this sequential exploration.³¹

Our theory suggests one. As described above, as training progresses, more and more examples get fit, starting with examples whose gradients are well-aligned with those of others, and ending with those examples whose gradients are more idiosyncratic. Thus, one may observe the function implemented by the network to get more and more complicated in the course of training simply because it fits (interpolates through) more and more training examples (points). An interesting

³⁰This may be seen as an additional justification for early stopping to the one in Section 2.

³¹For example, there is no causal intervention similar to suppressing weak gradient directions.

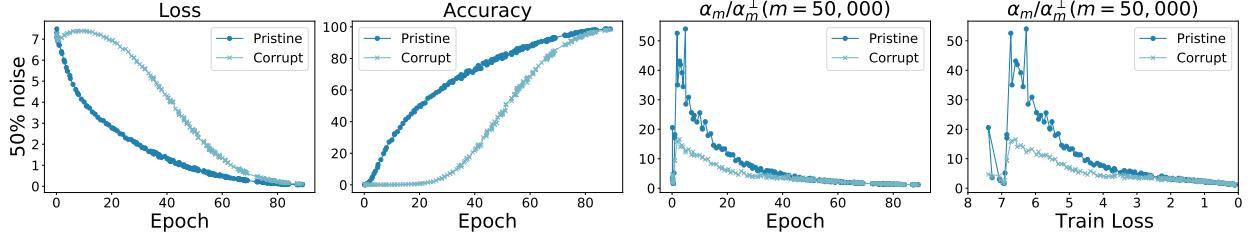


FIGURE 9. Pristine examples, that is examples with correct labels, show higher coherence than the corrupt examples, and consequently are learned much faster. Here, we train a Resnet-50 on ImageNet where half the images in the training set have randomly assigned labels (that is, ImageNet with 50% label noise).

question for future work is if this argument can be formalized to explain some of the specific empirical observations in the literature (such as ReLU networks learning low frequency functions first [Rahaman et al., 2019] or learning low descriptive-complexity functions first [Valle-Perez et al., 2019]).

However, one may wonder if there is some *other* causal mechanism that causes SGD to explore functions in increasing order of complexity? A simple extension of the previous experiment with Easy and Hard ImageNet suggests not. To motivate the extension, consider a thought experiment motivated by the 2D classification problem shown in Figure 8. If SGD were to somehow “try” simpler hypothesis early on, then the easy examples would be the ones far from the decision boundary (which could be easily separated by say, linear classifiers), and the hard ones would be ones close to the boundary (which would need more complex curves to separate). Therefore, based on the intuition provided by Figure 8, one might expect that the decision boundary learned from the hard examples, would generalize well to the easy examples.

But this is not what we see with real data in high dimensions: The model trained on Hard ImageNet has an accuracy of only 17% on the Easy ImageNet test set. This is much lower than the 71% test accuracy of the model trained on Easy ImageNet (as seen in Figure 7). In other words,

Examples learned late by gradient descent (that is, the hard examples), by themselves, are insufficient to define the decision boundary to the same extent that early (or easy) examples are.

This observation leads us to a view of deep models that is closer in spirit to kernel regressors [Nadaraya, 1964, Watson, 1964], but, of course, with a “learned kernel,” and is consistent with the intuition that in high dimensions, learning (almost) always amounts to extrapolation [Balestriero et al., 2021].

9. LEARNING WITH NOISY LABELS

The theory of Coherent Gradients provides insight on why it is possible to learn in the presence of label noise [Rolnick et al., 2017]. As Figure 9 shows, similar to real and random data, coherence as measured by α_m/α_m^\perp is greater for examples with correct labels (the pristine examples) than it is for examples with incorrect labels (corrupt examples). Consequently, the strong directions in the average gradient correspond to the pristine examples, and they get learned before the corrupt examples, mirroring the situation with easy and hard examples discussed in the previous section.

We can thus explain the empirical success of early stopping when learning with noisy labels and the empirical observation that the loss of pristine examples falls faster than that of corrupt

examples (see, for example, “small-loss” trick in Song et al. [2020, Section III E]). Furthermore, even if early stopping is not used, and training is continued until convergence, the gradients of corrupt examples are not strong enough to “interfere” with those of the pristine examples which reinforce each other.³²

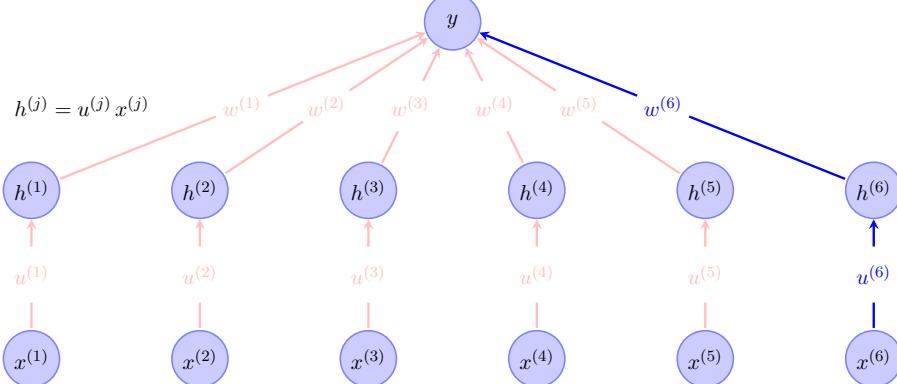
10. DEPTH, FEEDBACK LOOPS, AND SIGNAL AMPLIFICATION

Gradient descent through the lens of Coherent Gradients is a mechanism for soft feature selection, or equivalently, for signal amplification. The over-parameterized linear regression model in Example 1 provides a good illustration of this. When training the model on dataset \mathbf{L} (“real”), we see that the weight of the common signal grows more rapidly since increasing it simultaneously reduces the loss on all training examples.

In this section, we shall see how the signal amplification effect becomes stronger when the model has multiple layers, due to the interaction between the parameters belonging to the different layers. We illustrate this with two simple examples that build on Example 1 to add depth.

Example 7 (Adding Depth to Linear Regression). Instead of fitting the a model of the form $y = w \cdot x$ as is the case in linear regression, consider fitting the following “deep” model:

$$y = (w \odot u) \cdot x = \sum_{j=1}^6 w^{(j)} h^{(j)} = \sum_{j=1}^6 w^{(j)} u^{(j)} x^{(j)}$$



under the square loss $\ell(w) \equiv \frac{1}{2}(y - (w \odot u) \cdot x)^2$ to the dataset \mathbf{L} (“real”) from Example 1:

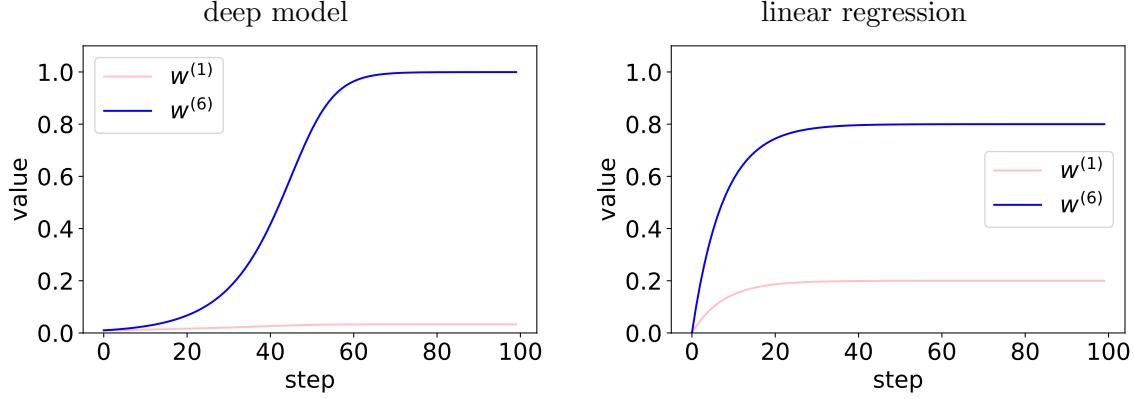
i	x_i	y_i
1	$\langle 1, 0, 0, 0, 0, 1 \rangle$	1
2	$\langle 0, -1, 0, 0, 0, -1 \rangle$	-1
3	$\langle 0, 0, -1, 0, 0, -1 \rangle$	-1
4	$\langle 0, 0, 0, 1, 0, 1 \rangle$	1
5	$\langle 0, 0, 0, 0, -1, -1 \rangle$	-1

Now start gradient descent (GD) with a constant learning rate of 0.1 at $w^{(j)} = u^{(j)} = 0.01$ for $j \in [6]$.³³ As in the case of linear regression, the parameters corresponding to the common feature $x^{(6)}$, that is, $u^{(6)}$ and $w^{(6)}$ grow more quickly than those corresponding to the idiosyncratic features

³²However, there may still be some coherence amongst incorrectly labeled examples which may cause unintended generalization to the test set. This is manifested as a degradation in test performance as the corrupt examples are learned, and can be seen in experiments (see, for example, Chatterjee [2020, Figure 1 (b)], Zielinski et al. [2020, Figure 1 (a)], and in Arpit et al. [2017, Figures 7 and 8 (b)]).

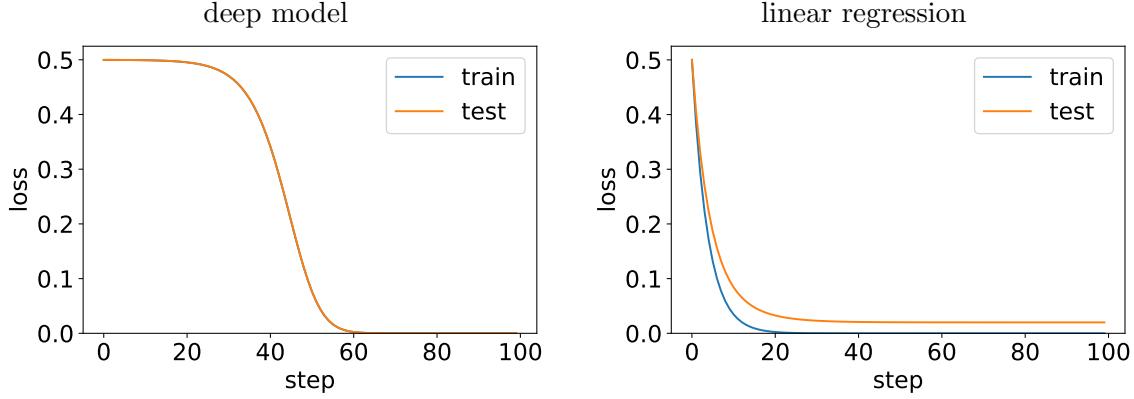
³³Unlike linear regression, we use a near-zero initialization for the deep model rather than a zero initialization since the latter is a stationary point.

$x^{(1)}, \dots, x^{(4)}$. But what is striking is the extent of the relative difference. This can be seen by comparing the plots of their values in the two cases:³⁴



In fact, for the deep model, the difference between $w^{(1)}$ and $w^{(6)}$ understates the importance of the corresponding features since the weight of feature i in the model is given by $u^{(j)}w^{(j)}$ (which, in this case, is $(w^{(j)})^2$ due to symmetry).

The improvement in amplification efficiency is also reflected in the generalization gap (using the 5th example in \mathbf{L} as the test example) which is negligible for the deep model in comparison to that of linear regression:



This naturally leads to the following question:

Why is the signal amplification so much greater in the case of the deep model compared to linear regression?

The answer lies in the interaction of the weights across the two layers. Consider the components of the gradient corresponding to $u^{(j)}$ and $w^{(j)}$ for the i th example. These are, respectively,

$$\frac{\partial \ell_i}{\partial u^{(j)}} = r_i(u, w) x^{(j)} w^{(j)} \quad \text{and} \quad \frac{\partial \ell_i}{\partial w^{(j)}} = r_i(u, w) x^{(j)} u^{(j)}$$

³⁴We only show $w^{(1)}$ and $w^{(6)}$ since, along the trajectory of GD, from symmetry, $u^{(j)} = w^{(j)}$ for all $j \in [6]$ in the deep model and $w^{(2)}, w^{(3)}$, and $w^{(4)}$ are equal to $w^{(1)}$ in both models.

where $r_i(u, w) \equiv y_i - (w \odot u) \cdot x_i$ is the residual. Observe that the component of the gradient corresponding to $u^{(j)}$ is directly proportional to $w^{(j)}$ and vice-versa. By slight abuse of terminology we call this a *positive feedback loop* between $u^{(j)}$ and $w^{(j)}$.³⁵

Now consider the gradient $g_i(u, w)$ for the i th example, and the average gradient $g(u, w)$:³⁶

i	$g_i(u, w) = \langle \frac{\partial \ell_i}{\partial u^{(1)}}, \frac{\partial \ell_i}{\partial u^{(2)}}, \frac{\partial \ell_i}{\partial u^{(3)}}, \frac{\partial \ell_i}{\partial u^{(4)}}, \frac{\partial \ell_i}{\partial u^{(5)}}, \frac{\partial \ell_i}{\partial u^{(6)}}, \frac{\partial \ell_i}{\partial w^{(1)}}, \frac{\partial \ell_i}{\partial w^{(2)}}, \frac{\partial \ell_i}{\partial w^{(3)}}, \frac{\partial \ell_i}{\partial w^{(4)}}, \frac{\partial \ell_i}{\partial w^{(5)}}, \frac{\partial \ell_i}{\partial w^{(6)}} \rangle$
1	$r(u, w) \langle w^{(1)}, 0, 0, 0, 0, w^{(6)}, u^{(1)}, 0, 0, 0, 0, u^{(6)} \rangle$
2	$r(u, w) \langle 0, w^{(2)}, 0, 0, 0, w^{(6)}, 0, u^{(2)}, 0, 0, 0, u^{(6)} \rangle$
3	$r(u, w) \langle 0, 0, w^{(3)}, 0, 0, w^{(6)}, 0, 0, u^{(3)}, 0, 0, u^{(6)} \rangle$
4	$r(u, w) \langle 0, 0, 0, w^{(4)}, 0, w^{(6)}, 0, 0, 0, u^{(4)}, 0, u^{(6)} \rangle$
$g(u, w)$	$r(u, w) \langle \frac{1}{4}w^{(1)}, \frac{1}{4}w^{(2)}, \frac{1}{4}w^{(3)}, \frac{1}{4}w^{(4)}, 0, w^{(6)}, \frac{1}{4}u^{(1)}, \frac{1}{4}u^{(2)}, \frac{1}{4}u^{(3)}, \frac{1}{4}u^{(4)}, 0, u^{(6)} \rangle$

We see that, just as in linear regression with dataset **L** (Example 1), the per-example gradients in the common component (shown in blue) reinforce each other, and add up in the average gradient $g(u, w)$, but the idiosyncratic components (shown in pink) do not.³⁷

However, in contrast to linear regression where the *relative* difference in the gradient directions was constant during training (1 v/s $\frac{1}{4}$), in the deep model, the relative difference depends on u and w (for example, $w^{(6)}$ v/s $\frac{1}{4}w^{(1)}$ and $u^{(6)}$ v/s $\frac{1}{4}u^{(1)}$). This combined with the positive feedback loop between the $u^{(i)}$ and the corresponding $w^{(i)}$ leads to the relative difference between the gradient directions *exponentially increasing* as training progress.

As a result, the relative importance of the common signal increases much faster in the deep model than in the linear regression case leading to a winner-take-all situation. In summary,

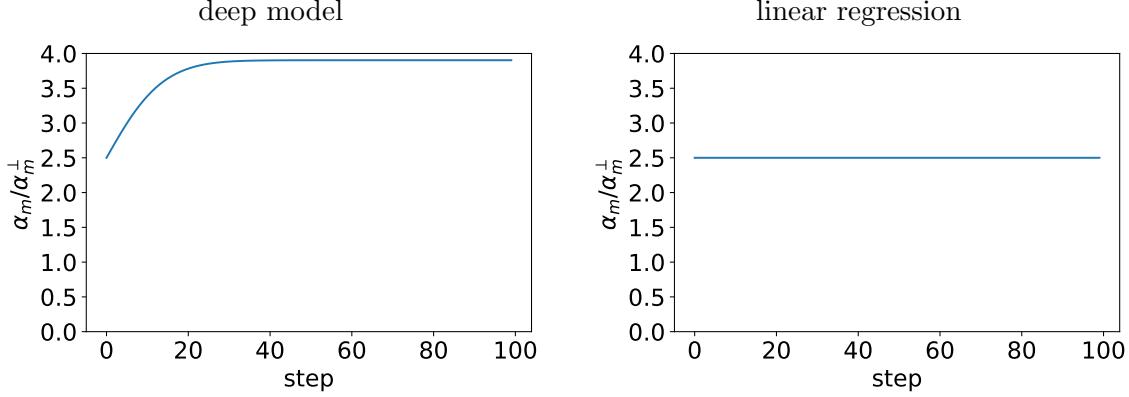
In the deep model, in contrast to linear regression, the relative difference between the common and idiosyncratic gradient directions gets exponentially amplified during training due to positive feedback between layers.

This rapid growth in the strong component of the gradient is reflected in the plot of training coherence as measured by α_m/α_m^\perp ($m = 4$) over time since the $u^{(6)}$ and $w^{(6)}$ components (the strong, or the high coherence components) come to dominate the gradient for the deep model, and they have perfect coherence:

³⁵It is an abuse of terminology since the positive feedback is not absolute. There is a second “dampening” dependency through the residual which becomes more important closer to convergence.

³⁶ $r(u, w)$ is equal to $|r_i(u, w)|$, a quantity that is independent of i for (u, w) in the trajectory of GD due to the symmetries in the problem.

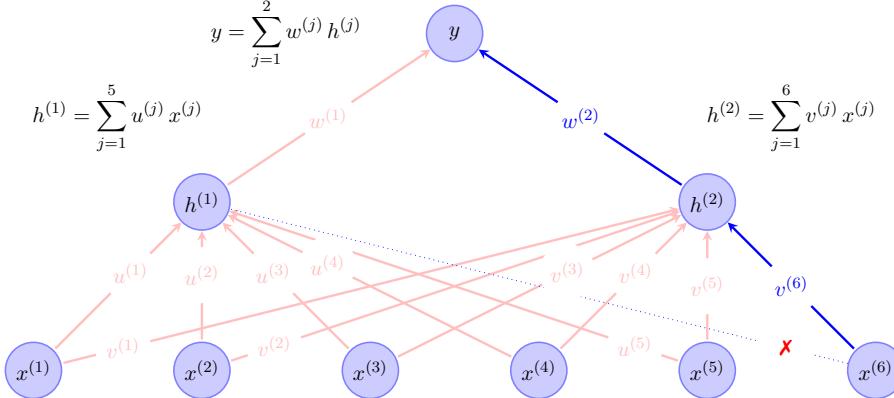
³⁷This is also reflected in the higher coherence of the $u^{(6)}$ and $w^{(6)}$ components compared to the rest as measured by (component-wise) α .



Thus although α_m/α_m^\perp for the deep model starts off at 2.5 (exactly the same as for linear regression), it rises to near maximum of 4 (while in the linear regression case it stays constant). \square

The previous example showed how feedback loops due to depth can amplify input signals or features that generalize better relative to noise signals. Feedback loops can also play a critical role in differentially amplifying *internal* signals in the network that generalize better than their peers. To see this, consider the following thought experiment that builds on Example 1 by adding a second layer to select between a “memorization” neuron and a “learning” neuron.

Example 8 (A Tale of Two Neurons). Consider the following linear neural network with one hidden layer with two neurons:



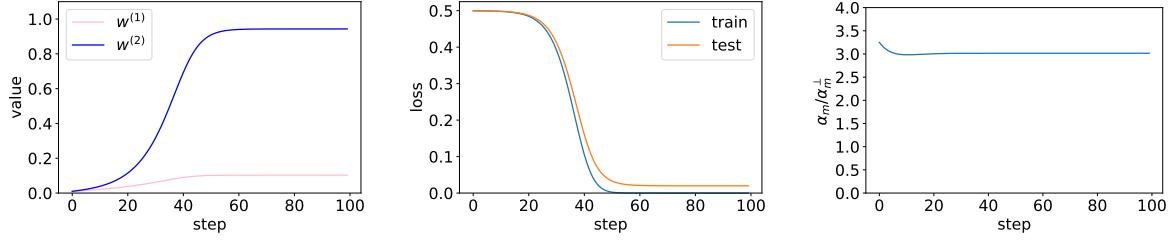
Now consider the task of fitting it to the following dataset (dataset **L** (“real”) from Example 1 that we have been using as a running example):

i	x_i	y_i
1	$\langle 1, 0, 0, 0, 0, 1 \rangle$	1
2	$\langle 0, -1, 0, 0, 0, -1 \rangle$	-1
3	$\langle 0, 0, -1, 0, 0, -1 \rangle$	-1
4	$\langle 0, 0, 0, 1, 0, 1 \rangle$	1
5	$\langle 0, 0, 0, 0, -1, -1 \rangle$	-1

Of the two neurons in the hidden layer, one, the “good” neuron, $h^{(2)}$, is connected to the common input feature $x^{(6)}$, whereas, the other, the “bad” neuron, $h^{(1)}$, is not. Therefore, the bad neuron can only help in memorizing the training data.³⁸

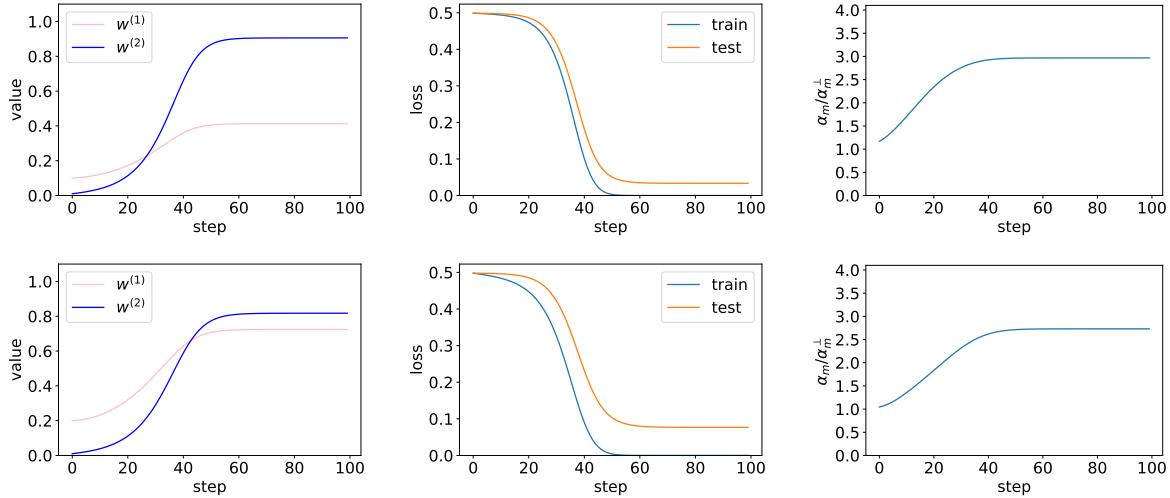
³⁸In terms of Example 1, the bad neuron “sees” the random dataset **M** instead of **L**.

When we perform gradient descent starting, as before, from 0.01 for all parameters, we find that the weight for the good neuron ($w^{(2)}$) grows much more rapidly than that for the bad neuron ($w^{(1)}$) and we have good generalization (as measured on (x_5, y_5)):



As in the previous example, the reason behind the rapid increase in $w^{(2)}$ relative to $w^{(1)}$ is due to the interaction between the layers. In brief, there is a positive feedback loop³⁹ between $w^{(1)}$ and each of $u^{(1)}, \dots, u^{(5)}$ and between $w^{(2)}$ and each of $v^{(1)}, \dots, v^{(5)}$. However, in addition, $w^{(2)}$ also has a feedback loop with $v^{(6)}$, a component with high coherence since it is the weight for the common input feature $x^{(6)}$. Consequently, $w^{(2)}$ grows exponentially faster than $w^{(1)}$.⁴⁰

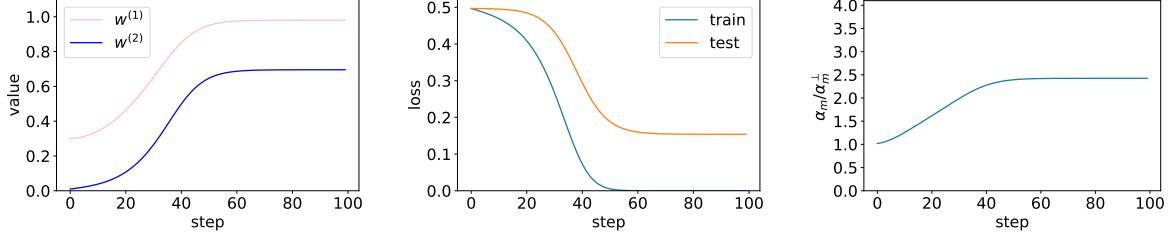
Adversarial initialization. The relative amplification effect is so strong that even if $w^{(1)}$ is very favorably initialized at 0.1 or 0.2 while keeping $w^{(2)}$ initialized at 0.01, in the course of training, $w^{(2)}$ eventually surpasses $w^{(1)}$ (we also show the test and train curves and α_m/α_m^\perp on the training set ($m = 4$) measured over the whole network for subsequent discussion):



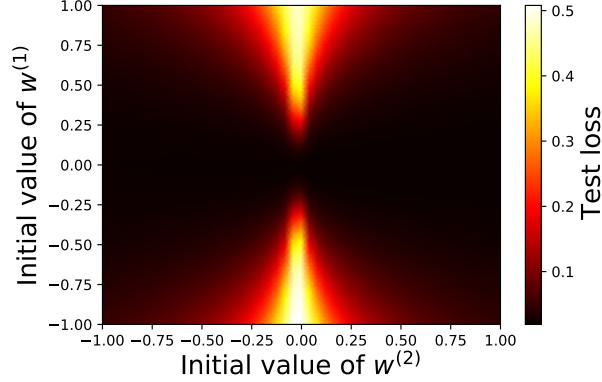
Of course, for a large enough adversarial initialization, $w^{(2)}$ can no longer overcome $w^{(1)}$ and we have poor generalization:

³⁹As in the previous example, there is a second dampening dependence through the residual that becomes important near convergence.

⁴⁰Indeed, due to the feedback loop, and the faster growth of $w^{(2)}$, even $v^{(1)}, \dots, v^{(4)}$ grow faster relative to $u^{(1)}, \dots, u^{(4)}$.



We can study this systematically by plotting a heatmap of the test loss at the end of training as a function of the initial values of $w^{(1)}$ and $w^{(2)}$ (all other parameters are initialized at 0.01):



The additional positive feedback loop between $w^{(2)}$ and $v^{(6)}$ ensures that from almost all initializations, the common feature has greater importance in the final output than any idiosyncratic feature. In summary,

Due to interaction between layers, the hidden neuron with access to the common feature gets a much higher weight compared to its peer that does not.

In this manner, gradient descent amplifies intermediate features that generalize better.⁴¹ □

Memorization and Layer-wise Coherence. In the previous example, it may be tempting to think that the difference in generalization capability between the good neuron and the bad neuron is reflected in the training coherence of their respective weights. That is not so. Both $w^{(1)}$ and $w^{(2)}$ have the same coherence over the training set during gradient descent for any of the above trajectories, that is, both have a component-wise α of 1. Thus,

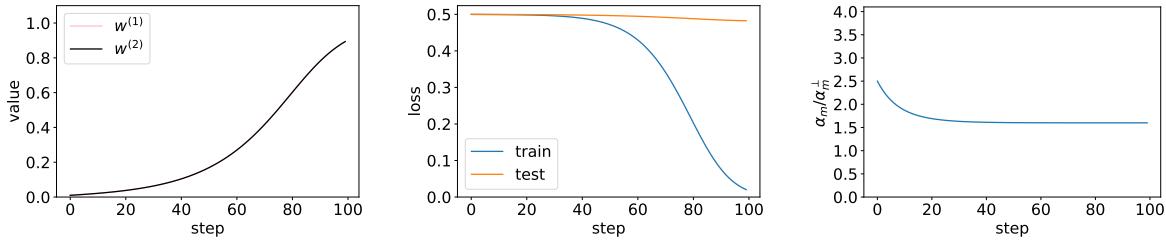
Coherence alone at a hidden layer cannot distinguish between a feature from the previous layer that does not generalize well and one that does.

This has two consequences for coherence measured on random data:

⁴¹Informally, a feature that generalizes better is one that is stable, that is, one that shows a smaller perturbation if the training set we slightly changed. A feature such as $h^{(1)}$ may be thought of as a “memorization” feature since it takes on very different values on a particular example depending on whether that example was in the training set or not. In this context, the method of counterfactual simulation presented in Chatterjee and Mishchenko [2020], a method to analyze the stability, and hence generalization of intermediate signals of a trained model, may be of interest.

- (1) The coherence of a given layer may be high if memorization happens in other layers. We believe that this is a likely reason for why high coherence is observed in AlexNet for the higher layers when it is trained on random data (Figure 3).
- (2) Coherence as measured over the whole network, that is, α computed using the entire gradient vector, can skew high since it is a weighted sum of the per-layer coherences, since many of the layers can have high coherence if memorization happens in only a few layers. This was seen in the case of AlexNet (Figure 2).

Example 9 (Tale of Two Neurons with Memorization). Both these consequences for random data can be illustrated by training the model used in Example 8 on the “random” dataset \mathbf{M} of Example 1 (instead of dataset \mathbf{L} as we have been doing so far):



It is easy to check, on the training set of \mathbf{M} ($m = 4$), that $\alpha_m/\alpha_m^\perp = 1$ for the first layer (where we may say memorization, or the case-by-case fitting, occurs) and $\alpha_m/\alpha_m^\perp = 4$ for the second layer (where all the examples are treated uniformly). Thus, the overall α_m/α_m^\perp is greater than 1 (the orthogonal limit), as seen above, even with memorization. \square

The observation above, that is, coherence of the entire network can be high on random data, is not specific to the α metric. It is indeed the case that the per-example gradients of different examples are similar in the components of the gradient corresponding to the higher (non-memorization) layers. Therefore, even other metrics for coherence such as dot product and stiffness that depend on the entire gradient vector would be skewed upwards.

To summarize, the examples in this section show how depth can help with feature selection through positive feedback loops between layers. Due to these feedback loops, features whose weights have high coherence get amplified beyond what coherence alone can accomplish.

11. WHAT SHOULD A THEORY OF GENERALIZATION LOOK LIKE?

The ultimate goal of a theory of generalization for deep learning would be to provide a tight guarantee on the generalization gap in problem setups of practical interest, but that appears out of reach at present. Even so, there are a number of different qualitative properties that such a theory should satisfy. For example, in the introduction, we already noted that the work of Zhang et al. [2017] puts a very strong constraint on the form of any theory of generalization:

[C1] Uniform explanation of memorization and generalization. It must simultaneously be able to explain why on some datasets, there is good generalization, when the exact same setup (architecture, learning rate, number of steps) is capable of memorizing a random dataset. Thus any non-trivial generalization bound must depend on the dataset in some non-trivial way (that is, beyond just the number of training examples). And ideally, in the case of label noise experiments, the bound should increase with increasing amount of label noise.

In addition to this constraint, practical experience with deep learning suggests several other constraints on a theory of generalization. Nagarajan and Kolter [2019b] propose the following criteria:

- [C2] **Incorporate architecture dependence.** It should provide an explanation of why the generalization error reduces with increasing width (or depth), as has been surprisingly observed in practice.
- [C3] **Apply to real networks (not post-processed networks, noisy variants, or infinite limits).** The theory should apply to the network learned by gradient descent without any modification or explicit regularization, or to idealizations.
- [C4] **Generalization should improve with more training data.** The generalization gap predicted by the theory should decrease with increase in the size of the training set (and not increase as is the case for many bounds, particularly those based on weight-norms!).

To these four criteria from the literature, we would like to add the following:

- [C5] **Uniform explanation of early stopping and training to convergence.** The theory should provide an explanation of why early stopping works in practice, and furthermore, such an explanation should ideally be uniform with respect to the two extremes of early stopping: perfect generalization at initialization (before any data is seen) and good generalization at end of training (that is, at a local minimum or a stationary point of the loss function), particularly since early stopping is not essential to good generalization (for example, see Figure 1).
- [C6] **Uniform handling of stochasticity of batch construction.** Although for many datasets and architectures of practical interest it is computationally prohibitive to do full-batch gradient descent, where such experiments have been performed (for example, Wu et al. [2018] and Chatterjee [2020]⁴²), or equivalent optimization outcome achieved in practice (for example, over-parameterized linear regression using full-batch gradient descent starting from the origin converges to the solution obtained by computing the pseudo-inverse using the singular value decomposition), we find that full-batch also leads to good generalization. Furthermore, careful experiments show that there is no clear connection between batch-size and generalization gap though it was suggested previously that there may be (see Shallue et al. [2019] for a comprehensive study). (Now, this is not to say that stochasticity does not help: for example, sampling noise from mini-batch construction could help get out of local minima, or perhaps even serve as a regularizer; but it appears that stochasticity is not essential for a non-trivial generalization gap.)
- [C7] **Right level of generality for the phenomenon.** Generalization in deep learning is an extremely broad phenomenon and has been observed in many situations that differ in specifics along dimensions such as type of activation function (for example, ReLU, sigmoid, sinusoidal, etc.); type of loss function; whether the task is a regression, classification (few classes or very many classes) or a generative task; discreteness of the problem (both in inputs, such as language models, and in weights and activations, such as extremely quantized models such as binarized neural networks [Hubara et al., 2016]); training schedules such as SGD with warm restarts [Loshchilov and Hutter, 2017]; etc.

Therefore, a satisfactory explanation of generalization in deep learning cannot be too closely tied to the specific properties of a system such as scale invariance of ReLU activations (the explanation would not work for other activation functions), or the learning problem being a classification problem with a few output classes such as a margin-based explanation (the explanation would not be natural for regression or generative models, or even

⁴²See data here: <https://openreview.net/forum?id=ryeFY0EFwS¬eId=rJg1FKTBoB>

classification with 100K output classes), or even the training steps monotonically reducing in size (would not work for warm restarts).

To this list, one might be tempted to add: *uniformity of the explanation with respect to the degree of over-parameterization*, that is, an uniform explanation for both under-parameterized models and over-parameterized models. However, given the observations of Belkin et al. [2019] around double descent, such a criterion may be difficult to satisfy. We believe that this dimension is an interesting one to explore in the future, but for now we focus only on the over-parameterized case.

We believe our approach satisfies all the criteria outlined in this section (even though some such as **[C2]** require further work). In the next section, we contrast our approach with the leading alternatives using these criteria as the framework for comparison.

12. COMPARISON WITH OTHER THEORIES AND EXPLANATIONS

There has been a lot of work to understand the nature of deep learning. At the highest level, we can categorize the efforts into two buckets: understanding optimization, and understanding generalization. This follows from decomposing the test loss as follows:

$$\text{test loss} = \text{training loss} + \text{generalization gap}$$

where the generalization gap is defined as the difference between the test and training losses.

The work on optimization focuses on trying to understand why given the non-convexity of the fitting problem in general, gradient descent can reach good global minima (and often zero training loss). On the other hand, the work on generalization tries to understand why the generalization gap is small. Although the work on optimization is fascinating in its own right, and there may be reasons to believe that there is an interesting interplay between optimization and generalization in deep learning (for example, our work suggests that when one can optimize quickly, that is also when one is most likely to generalize well), our focus here is on generalization. We refer the reader to Liu et al. [2020a] as a recent entry point into the optimization literature.

Approaches based on Uniform Convergence. Uniform convergence [Vapnik and Chervonenkis, 1971] is the earliest and most popular mathematical tool for proving generalization. This theory was initially developed to understand under what circumstances it is possible to identify a model or a hypothesis (from a possibly infinite family) based on only a finite training sample, and it provided a framework to bound the generalization gap based on the complexity of the family of hypotheses (quantified, for example, by notions such as Vapnik-Chervonenkis (VC) dimension) from which the model was selected relative to the size of the training set.

Early generalization bounds for neural networks were obtained by bounding their VC dimension by the parameter count [Baum and Haussler, 1989], but even then it was realized that far fewer training samples are needed for good generalization in practice than is implied by these bounds. Bartlett [1996, 1998] showed in the case of sigmoid networks, that as long as the learning algorithm finds a network with small weights, what matters for a small generalization gap is not the number of weights (i.e., the parameter count), but the size of weights (i.e., their norm).

The growing popularity of deep learning along with the surprising experimental results of Zhang et al. [2017] (showing that deep networks used in practice can easily fit random data) brought renewed attention to the generalization problem in deep learning. The results of Zhang et al. [2017] imply that the complexity of hypothesis class in the case of practical neural networks is already too large in comparison with the size of the training set to make these bounds non-vacuous (that is, they fail **[C1]**). This motivated a large number of new investigations that try to capture the implicit bias or regularization of gradient descent (for example, Bartlett et al. [2017], Golowich

et al. [2018], Arora et al. [2018], Neyshabur et al. [2018a], Dziugaite and Roy [2017], Zhou et al. [2019], Li and Liang [2018], Allen-Zhu et al. [2019a], Nagarajan and Kolter [2019a], Neyshabur et al. [2019]). While these bounds aim to satisfy **[C1]**, and often provide valuable insights into the nature of deep learning, Nagarajan and Kolter [2019b] observe that they generally fail to either have the correct architecture dependence (**[C2]**) or are not applicable to the network obtained from gradient descent, but to some modification thereof (**[C3]**). Furthermore, they point out that many of these approaches (particularly, the non-PAC-Bayesian ones), fail in a more fundamental way by failing to satisfy **[C4]**: the bounds on the generalization gap obtained from these approaches increases with more training examples. This happens because the bound depends on some norm of the weights (that is, distance from initialization) that typically *increases* with larger training sets.⁴³

Finally, Nagarajan and Kolter [2019b] construct examples of overparameterized problems involving linear models and neural networks where uniform convergence *provably* fails to explain generalization. They show that the set of models with low test error and actually explored by GD (that is, every model in the set corresponds to a GD solution for some training set) is too large to get anything other than an almost vacuous uniform convergence bound for its members.

For the expert in uniform convergence, we note that the fact that uniform convergence is necessary for learnability (in the setting of classification or regression) does not preclude the fact that uniform convergence may be unable to explain generalization in a particular algorithm for a particular distribution (which is the situation we are faced with in explaining the generalization mystery in deep learning). For a detailed discussion of this, please see Appendix I of Nagarajan and Kolter [2019b]. In fact, to their argument, we may add that what we are interested in is not just a particular algorithm for a particular distribution, but even a specific sample size. In other words, a result that is asymptotic in sample size may not be illuminating at all. For example, consider that there are only finitely many ResNet-50 networks when the parameters are quantized to 32 bits. Thus, for a large enough sample size, even the humble combination of Hoeffding’s inequality and the union bound would provide a non-trivial bound on the generalization gap, but that argument does not provide much insight on why ResNet-50 trained with gradient descent generalizes well on ImageNet given the tiny size of the actual ImageNet training set in comparison.

In comparison, our approach is not based on uniform convergence, but on a different tool from learning theory, namely algorithmic stability, and so does not suffer from the problems pointed out by Nagarajan and Kolter [2019b]. In particular, note that the bound in Theorem 1 has an explicit inverse dependence on the size of the training set, as expected.

Previous approaches based on Algorithmic Stability. In short order after the development of uniform convergence, an alternative framework was developed to explain generalization in models based on “local rules” (for example, nearest neighbors) where the complexity of the hypothesis class is not fixed but scales with the size of the training dataset [Rogers and Wagner, 1978, Devroye and Wagner, 1979]. Unlike uniform convergence where the focus is on the *size* of the space to be searched, the focus here is on *how* the algorithm searches the space. This approach is now known as Algorithmic Stability [Bousquet and Elisseeff, 2002, Shalev-Shwartz et al., 2010, Kearns and Ron, 1999], and the key idea, as we saw, is as follows: If the model learnt from the training set does not depend too much on any one training example (that is the presence or absence of any one example does not change the model very much) then we expect the model to generalize well.

Hardt et al. [2016] analyzed the stability of *stochastic* gradient descent in the setting of deep learning under the assumption that each batch consists of one example and the learning rate decays over time. Under these assumptions, they are able to bound the generalization gap by the

⁴³The PAC-Bayesian ones typically fail **[C3]** since they apply only to stochastic ensembles. We discuss them in greater detail in a little later.

number of training steps, a result they summarize by the slogan “train faster, generalize better.” Kuzborskij and Lampert [2018] tighten the Hardt et al. [2016] bound by taking into account the curvature of the loss function at initialization but at the cost of a more onerous assumption: no training example is looked at more than once.

This line of argument critically relies on the assumption that each training sample is seen rarely (or not at all) and so the presence or absence of a training sample cannot matter too much. But one should feel uneasy about such an explanation since it may throw the baby out with the bath water. For example, if we don’t see any example at all (that is, take 0 steps of training and stay at the randomly chosen initial point), then the generalization gap should be zero in expectation, that is, we should expect the training loss to be an unbiased estimate of the test loss. The previous slogan taken to its logical conclusion is: Train in no time, generalize perfectly!

So what has gone wrong? The problem is that the generalization mystery requires us to keep in mind simultaneously the real data case and the random data case, and what we need is a “differential” analysis of the two cases. In other words, in the random case, if you see each example rarely (or not at all), you will not be able to memorize it, but the fact that in practice neural networks memorize random labels means that they are run for long enough to see each example several times, and it is in *that* setting that we need to explain generalization on real data. In other words, the existing stability approaches fail to satisfy [C1] above. Finally, as noted in Section 7 of the extended version⁴⁴ of Hardt et al. [2016], this approach does not extend to full-batch case, that is, it fails [C6].

Our result may be seen as a natural extension of this approach that takes into account the dataset (more precisely, the interaction of the dataset and the architecture) which is necessary to explain the generalization mystery.

We note here a line of work that combines stability and PAC-Bayesian analysis to study noisy gradient descent (stochastic gradient Langevin dynamics) where isotropic Gaussian noise is added at every SGD update step (thus failing to meet [C3] and [C6]). Although a more detailed discussion is out of scope, please see Li et al. [2020] for a starting point into this literature.

Sharpness and Flatness of Minima. A popular explanation of why SGD generalizes well in large networks in practice is based on the empirical observation that the minima found by SGD are “flat,” that is, they are in regions where the overall loss (on the training set) does not change very much around the solution found by SGD [Hochreiter and Schmidhuber, 1997, Keskar et al., 2017]. While this is intuitively appealing (for example, see Figure 1 in [He et al., 2019]), and may be combined with a PAC-Bayesian analysis [McAllester, 1999] to provide non-vacuous bounds for stochastic (noisy) variants of real networks [Langford and Caruana, 2002, Dziugaite and Roy, 2017, Zhou et al., 2019], as a *fundamental* explanation of generalization it is unsatisfactory for several reasons:

- (a) **Does not work for linear regression.** Flatness is not useful in understanding why gradient descent generalizes well in the simple case of over-parameterized linear regression (such as Example 1) since the Hessian of the loss function is constant over the parameter space, and thus the same at both well and poorly generalizing minima [Zhang et al., 2021, Section 5].

⁴⁴<https://arxiv.org/abs/1509.01240>

- (b) **Not invariant to reparameterization in deep networks.** In deep neural networks, it is difficult to formalize the notion of flatness in a manner that is invariant to reparameterization. For example, Dinh et al. [2017] show that one can construct arbitrarily sharp minima through weight reparameterization without affecting the generalization performance.
- (c) **Not uniform with respect to early stopping.** That is, it fails [C5]. It provides different explanations for good generalization at start of training (“initialization being random is independent of the training data”), and for end of training (“gradient descent is biased towards producing solutions that lie in flat regions”).
- (d) **It is incomplete.** It does not provide an explanation for why gradient descent prefers flat minima in deep nets, only that it does so (but see Jastrzebski et al. [2020, 2021] for exciting recent progress in this direction based on choice of early learning rate). It is important to note that the PAC-Bayesian bounds on the generalization gap mentioned above only exploit this empirical observation but do not explain it (for example, see the last paragraph of Section 1.1 in Dziugaite and Roy [2017]).

From an algorithmic stability perspective, the relationship between the minima of different examples is more fundamental than the nature of those minima. For example, if the loss of each example has a sharp minimum (for whatever definition of sharpness one likes), but the minima of all the examples coincide, then one could expect better generalization than when each example has a shallow minimum, but the minima of different examples do not coincide. Of course, it is possible that in practice, the shallowness of the overall loss function is a consequence of better alignment of the minima of different examples, but the point is that the flatness of the overall loss is a secondary to the relationship between the minima of different examples.

Finally, we note that just as the flatness perspective provides inspiration to modify gradient descent to find flatter minima to improve generalization, as seen in Section 7, the Coherent Gradients perspective provides similar inspiration to design new algorithms to avoid weak directions.

On Minimum Norm Solutions. Another popular explanation of why GD generalizes well is that it finds solutions with small ℓ^2 norms. There has been renewed interest in understanding the over-parameterized linear regression case—where GD provably finds the minimum ℓ^2 norm solution⁴⁵—with the hope that these results may be extended to deep networks using observations from the infinite width scaling limit [Hastie et al., 2020, Bartlett et al., 2020]. However, an explanation based on the size of the norm obtained by gradient descent is also unsatisfactory as a *fundamental* explanation for much of the same reasons as before:

- (a) **It is incomplete or has wrong asymptotics.** Without an explanation for *why* gradient descent prefers small norm solutions in deep networks, any explanation for the generalization mystery based on smallness of ℓ^2 norm is incomplete. For example, Bartlett [1996] showed that certain deep networks generalize well *if* gradient descent finds small norm solutions. But this explanation is incomplete because the question for the generalization mystery then is: Why does gradient descent find a small norm solution in one case (real data) and not in the other (random)? That is, why does SGD not always find the large norm (memorization) solution? More recent work bounds the norm in terms of the size of the training set m , but as Nagarajan and Kolter [2019b, Section 1.1] observe, these bounds typically have the wrong asymptotics in m (they predict generalization gap increases with more training examples).

⁴⁵though only when initialized at $\mathbf{0}$

- (b) **Not invariant to reparameterization in deep networks.** For example, in a ReLU network, one can change the norm without changing the function, and hence, the generalization of the model.
- (c) **Not uniform with respect to initialization.** That is, it fails [C5]. In particular, there is good generalization at initialization regardless of the norm of the parameters (since the parameters are chosen independently of the data).

Furthermore, solutions with non-optimal ℓ^2 norm solutions can generalize better than those with minimum ℓ^2 norm (see Kawaguchi et al. [2017], Shah et al. [2020], or even the median solution in Example 1).

We believe that an explanation based on stability and alignment of gradients during training is a deeper and more general explanation than flatness or small ℓ^2 norm, especially since it provides an uniform explanation for generalization in over-parameterized linear regression and for non-minimal ℓ^2 norm solutions and early stopping.

Coherent Gradients may provide a deeper explanation of why SGD gravitates towards small norm or flat-minima solutions, and that would be an interesting direction for future work.

Kernel Methods and the Neural Tangent Kernel. Belkin et al. [2018] point out that traditional kernel methods show the same generalization mystery that deep neural networks do, and furthermore can be viewed as a special case of neural networks (a kernel machine is a two-layer neural network with the first layer frozen). Therefore, they argue, in order to understand deep learning, we need to understand “shallow” kernel methods first, and foreshadowing Nagarajan and Kolter [2019b], they point out that the traditional norm-based bounds for kernel methods have the wrong asymptotics. We believe that the notion of algorithmic stability (rather than uniform convergence) is the correct lens through which to understand generalization even in kernel machines, and that our theory for understanding generalization in deep learning can be extended to understand kernel learning as well.

Recently, there has been an increased appreciation that the behavior of certain classes of deep models become simpler in the limit that their width goes to infinity, and certain convergence results to zero training loss can be established [Du et al., 2019b, Li and Liang, 2018, Du et al., 2019a, Allen-Zhu et al., 2019b, Zou et al., 2020]. Jacot et al. [2018] showed that training behaves like kernel regression in this limit with a (constant) kernel that is derived from the linearization of the network around its initialization, called the *Neural Tangent Kernel* (NTK).

While this is a promising line of research, particularly with respect to understanding optimization, it does not directly shed light on the generalization mystery. In essence, it makes the observation that under some structural assumptions (for example, infinite width limit with linear output layer but no bottleneck layers [Liu et al., 2020b], or a suitably large initialization [Chizat et al., 2019]), a complex non-linear model can behave like a simple linear model (kernel regression). Thus the generalization mystery in those cases becomes: why does that simple kernel regression model generalize well in one case (real data) and not in another (random data)?

Our theory provides an answer for the latter question—as it does even for ordinary over-parameterized linear regression—but our explanation is general enough that it can be applied directly to the complex models, thus obviating the need for the structural assumptions that guarantee a constant kernel. From our list of criteria for a theory of generalization, this line of work falls short of [C3] and [C7] (and arguably [C1] to the extent it even aims to address that). We stress that our concern here is generalization alone: we do not address why an over-parameterized network can

fit the training data. The NTK line of research sheds valuable light on that fitting question, and points the way to more general theories of optimization for deep learning such as [Liu et al., 2020a].

It is unclear if real networks actually operate in the constant kernel regime. Chizat et al. [2019] argue that they do not, and show experimentally that the lazy training degrades performance.⁴⁶ The experiments of Arora et al. [2019] on convolutional nets also show a significant gap between the constant kernel, and finite convolutional networks (even with regularization turned off), and there is work to close this gap by identifying more sophisticated tangent constructions that, for example, depend not just on the input images but the output labels as well [Chen et al., 2020]. At any rate, Chizat et al. conclude that “the intriguing phenomenon that still defies theoretical understanding is [the non-lazy regime]: neural networks trained with gradient-based methods (and neurons that move) have the ability to perform high-dimensional feature selection through highly non-linear dynamics.” Our work tackles this head on (see Section 10), and we see that the non-linear dynamics help even in very simple (and very finite!) cases (such as Examples 7 and 8).

Finally, a “kernel” that varies during training (the “path kernel”) has been considered by Domingos [2020], but he does not investigate how the path taken by SGD is chosen, and thus does not answer the generalization mystery. His theory (implicitly) assumes that the paths for real and random data are different, and therefore does not explain why models along one path generalize better than those along the other. Our coherence and stability-based analysis answers that question.

Alignment of Per-Example Gradients. In addition to our own work outlining some preliminary results [Chatterjee, 2020, Zielinski et al., 2020, Chatterjee and Zielinski, 2020], there has been a small, but growing, number of other efforts to understand optimization and generalization in deep learning from the perspective of gradient alignment. Here we provide a brief overview of these efforts, and defer detailed discussions of theoretical connections and experimental results to Appendix B and Appendix D respectively.

Yin et al. [2018] identify a quantity called *gradient diversity* on the training set that they use to study the question of convergence speed saturation of SGD, that is, what is the point at which increasing batch size no longer helps reduce the number of passes needed through the data (leading to more aggregate computation). They also study the corresponding question for generalization, that is, when is the generalization with a larger batch size still similar to the batch-size of 1 case considered in Hardt et al. [2016]. They restrict their attention to the convex case and provide a bound on the batch size based on a quantity called *differential gradient diversity* (which like diversity is also computed on the training set).

Fort et al. [2020] defined two measures of per-example gradient alignment called *sign* and *cosine stiffness* and conducted experiments to study the connection between stiffness and generalization. Sankararaman et al. [2020] propose a measure for gradient alignment called confusion and study its effect on the convergence of SGD (as the depth and widths of the networks change). They do not study generalization. Recently, Mehta et al. [2021] present an adversarial initialization scheme for the sine activation [Liu et al., 2019], and study it using the lens of gradient alignment using a metric very similar to coherence. They show experimentally on small datasets such as CIFAR-10, SHVN, etc. that within a class, their metric correlates well with generalization, as the scale of initialization is varied.

We note here that these efforts, particularly those looking at generalization, study gradient alignment within (and sometimes across) classes which is unsatisfactory since class-based analyses do not extend naturally to large output spaces as in regression or generative models (see [C7]). Furthermore, in the context of generalization, their focus is primarily on the question of when GD generalizes well (that is, Question 1 in Section 1). In our work, we are also interested in

⁴⁶Liu et al. [2020b] however argue that constant kernel is not the same as lazy training, and the kernel may be constant even in cases where the parameters change significantly (depending on the Hessian).

understanding the causal mechanism that enables GD to find a well-generalizing solution when such a solution exists (that is, Question 2 in Section 1).

Although not directly based on per-example gradients, the notion of *local elasticity* discovered experimentally by He and Su [2020] is worth highlighting here. They find that an SGD update in a deep network for an input x does not change the predictions for a “dissimilar” input x' . Although they do not provide a general definition of dissimilarity, intuitively, images of a Persian cat and an Egyptian cat (to use an example from their paper) would be less dissimilar than images of a German shepherd and a trailer truck. They note that local elasticity is observed (a) only in deep models (and not, say, in linear models) and (b) is “pronounced only after several epochs of training” [Deng et al., 2021, footnote 1]. Thus, as an explanation of generalization of gradient descent it fails [C5] and [C7]. In later work, they have used the notion of local elasticity to suggest an improvement to the NTK [Chen et al., 2020] (this was discussed previously), as well as a new notion of algorithmic stability called *locally elastic stability* [Deng et al., 2021] that is an elegant refinement of the notion of uniform stability to incorporate data dependence. They show that SGD is locally elastic stable, but their proof is based on the same insight as Hardt et al. [2016] that each training example is seen rarely (due to stochasticity), and therefore, it does not naturally extend to the non-stochastic case (fails [C6]).

Finally, Liu et al. [2020c] is the closest work in the literature to ours, but their primary focus is also Question 1.⁴⁷ They define two quantities *Gradient Signal to Noise Ratio* (GSNR) (which is on a per-coordinate basis on the training set) and *One Step Generalization Ratio* (OSGR) (which links the one-step change in the training and test losses) and show how in *early* training, under the assumption that the average gradients in test and training are the same (their Assumption 2.3.1), in a full-batch setting, the two quantities are related, and high GSNR implies an OSGR close to 1. We stress that their result is a local result that holds only for one step of training, unlike Theorem 1 which holds globally. Furthermore, a reformulation in terms of {test, training} \times {overall, per-coordinate} quantities provides a simpler view of their result. See discussion of GSNR and OSGR in Appendix B for more details.

Furthermore, like us, they also observe that in deep networks, the alignment of per-example gradients can increase in the course of training (that is, there is coherence amplification). However, in their experiments, they only see it on real data and not with random labels in contrast to what we see in Section 6 with our metric on ResNet (where it is a small effect) and on AlexNet (where it is a large effect). They provide a heuristic argument for why the numerator of GSNR can increase due to feature learning, but stop short of their ultimate goal which is showing why GSNR itself can increase. Their heuristic explanation based on the dynamics of the training is similar to our discussion in Section 10 but see our discussion of the difference between coherence amplification and signal amplification. Of course, here it must be noted GSNR does not always increase (even in their data), and in our work, we show why that must be so as examples get fitted albeit in the context of our coherence metric α rather than GSNR (see Appendix F).

It should be noted that since they do not address Question 2 in Section 1, their theory is descriptive rather than causal, and does not directly suggest modifications such as those considered in Section 7 to reduce memorization and improve generalization.

13. DISCUSSION AND DIRECTIONS FOR FUTURE WORK

In this paper, we have presented a simple explanation for the generalization mystery in deep learning: When the gradients of different examples are well-aligned during training, that is, when

⁴⁷It is interesting to note that their paper appeared in parallel with our preliminary work [Chatterjee, 2020] which addressed the other half of the generalization mystery, that is, Question 2 in Section 1.

they are *coherent*, gradient descent is stable, and the resulting model is expected to generalize well. Otherwise, gradient descent may fail to generalize if there are too few examples or if it is run for too long.

This explanation has several attractive features. First, due to its simplicity, it is very general. It provides an uniform explanation for memorization and generalization, and scales from toy over-parameterized linear regression problems to deep models at ImageNet scale. Second, it separates optimization-related issues (which can be hard to reason about in the non-convex settings common to deep learning) from generalization-related ones, allowing us to make progress independently on these aspects. Third, it is a *causal* explanation for the implicit bias of gradient descent that provides perspective on existing regularization techniques, the relative hardness of examples, the role of depth, and motivates novel but simple modifications such as winsorization and M3 for improving the generalization of gradient descent. Last, but not least, it avoids critical (and possibly fatal) problems facing competing explanations that have been put forward to explain the generalization mystery.

There are many directions to explore from here, particularly, in the formal development of the theory, in establishing how widely it holds, in using it to improve our understanding of deep learning, and in developing new, more efficient, representation learning algorithms. We outline some specific ideas in greater detail below:

- (1) **Better metrics for coherence and a tighter bounds.** Our generalization bound depends on α as measured over entire per-example gradients. However, as we have seen, both on toy examples and some real architectures, some parts of the network may play a greater role in memorization than others, and a measurement over the entire per-example gradient averages over these differences leading to a necessarily weak bound. Can we obtain tighter bounds by accounting for the graphical structure of the network, say by developing some notion of “minimum coherence” cuts that control generalization?
- (2) **Incorporating the under-parameterized case.** Our goal here was to understand what happens in highly over-parameterized settings. While the inverse dependence on the number of examples m in our bound provides a limiting guarantee in the case of extreme under-parameterization (as $m \rightarrow \infty$), a more direct dependence on the dimension, or better yet, on the graphical structure of the model (as noted above), is desirable. The ideal theory would handle the over- and under-parameterized cases uniformly, thus providing a continuous explanation for the phenomenon of deep double descent [Belkin et al., 2019] and also apply to pathological learning problems that are simultaneously over- and under-parameterized (in different parts of the data distribution). See also Appendix H.
- (3) **Generalization bound based on training set only.** Is there an efficient way to measure stability degradation during training using *only* the training set in order to provide a non-vacuous bound on the generalization gap? Perhaps this could be done with a “stability accountant” similar to a privacy accountant in differential privacy.
- (4) **Confirmation on other datasets and architectures.** In this work we have only studied vision datasets. Does the explanation hold for other types of data (such as language and speech) and architectures (such as transformers)?
- (5) **Generalization and width.** Neyshabur et al. [2018b] found that wider networks generalize better. Can we now explain this? Intuitively, wider networks have more sub-networks at any given level, and so the sub-network with maximum coherence in a wider network may be more coherent than its counterpart in a thinner network, and hence generalize better. In other words, since—as discussed in Section 10—gradient descent is a feature selector that prioritizes well-generalizing (coherent) features, wider networks are likely to have better

features simply because they have more features. In this connection, see also the Lottery Ticket Hypothesis [Frankle and Carbin, 2018].

- (6) **Understanding texture and shape bias.** Can use coherence at random initialization and signal amplification to understand why vision networks have a bias towards features based on “bulk” properties such as texture as opposed to those based on more subtle properties such as shape (see, for example, Geirhos et al. [2019])?
- (7) **Practical use for new algorithms.** Can algorithms based on robust averaging (such as winsorized gradient descent, or M3) help improve the state-of-the-art in large-scale learning with noisy labels, learning in low data settings, and in private learning? Can they be extended to provide privacy guarantees?
- (8) **Hardware acceleration for private learning.** Modern hardware accelerators for deep learning typically have multiple nodes (for example, GPU or TPU cores) over which the examples in a mini-batch are distributed, and the gradient for a mini-batch is computed via a global all-reduce operation over the average gradient computed at each node [Xu et al., 2020, Kumar and Jouppi, 2020]. Platform support for robust averaging (by clipping outliers), perhaps by supporting multiple simultaneous operations as part of an all-reduce, would greatly improve the efficiency of algorithms such as M3 for the use cases outlined above.
- (9) **Using coherence to evaluate architectures.** Can we use coherence to get insights into the relative strengths and weaknesses of different network architectures, and to improve architectures to generalize better? For example, as we saw in Section 6, for random data, α of (entire) per-example gradients is much lower for ResNet-50 than for AlexNet. Does that mean that overfitting in ResNet-50 is harder than in AlexNet since it cannot be as “local” (confined to a few layers)?
- (10) **Using coherence to study optimizers.** In this work we have restricted our focus to vanilla gradient descent. However, the notion of coherence can be used to analyze the generalization behavior of other optimization algorithms such as ADAM, a topic of much interest (see, for example, Zhou et al. [2020] and references therein). For other optimizers, in addition to the coherence of the per-example gradients, one would have to consider the transformation applied by the optimizer to obtain the actual updates to the network parameters.
- (11) **Explore Implications of the Theory.** In order to gain further confidence in the theory, it would be good to check some of the following predictions which we believe are consequences of the theory:
 - (a) *Full-batch gradient descent and gradient flow should have non-trivial generalization on real data.* Based on our theory, finite learning rate, and mini-batch stochasticity are not necessary for generalization (though they may help). Note that we are only talking about generalization gap, not test loss: Finite learning rate and mini-batch stochasticity can help escape saddle points during training, and so in their absence we may get stuck with poor training loss, but any source of noise could be used to escape these saddle points. Thus, we would expect, for instance, that full-batch gradient descent on ImageNet with a tiny learning rate should show non-trivial generalization.
 - (b) *Assuming a small enough learning rate, as training progresses, the generalization gap cannot decrease.* This follows from the iterative stability analysis of training: with more steps, stability can only degrade. If this is violated in a practical setting, it would point to an interesting limitation of the theory.

- (c) *Continuity is not necessary for good generalization.* Although gradient descent involves real numbers (or close approximations), if the theory is right, it should be possible to engineer fully discrete representation learning systems using voting among training examples to make optimization decisions during fitting that generalize well. By avoiding floating point operations (or even high precision fixed point), such systems may be much more computationally efficient than gradient descent.

ACKNOWLEDGMENTS. We thank Alan Mishchenko and Shankar Krishnan for close collaborations on related projects, and the following for interesting discussions and feedback: Anand Babu, Sungmin Bae, Michele Covelle, Josh Dillon, Sergey Ioffe, Firdaus Janoos, Stanislaw Jastrzebski, Chandramouli Kashyap, Matt Streeter, Rahul Sukthankar, and Jay Yagnik.

A. MATHEMATICAL PROPERTIES OF α

Our metric for coherence α defined in Section 4 is not specific to gradients, but extends naturally to vectors in Euclidean spaces. It is convenient, therefore, to study its properties in that general setting. Let \mathcal{V} be a probability distribution on a collection of m vectors in an Euclidean space. In accordance with (3), we define $\alpha(\mathcal{V})$ to be

$$\alpha(\mathcal{V}) \equiv \frac{\mathbb{E}_{v \sim \mathcal{V}, v' \sim \mathcal{V}} [v \cdot v']}{\mathbb{E}_{v \sim \mathcal{V}} [v \cdot v]} \quad (6)$$

Note that $\mathbb{E}[v \cdot v] = 0$ implies $\mathbb{E}[v \cdot v'] = 0$. In what follows, we ignore the technicality of the denominator being 0 by always assuming that there is at least one non-zero vector in the support of \mathcal{V} .⁴⁸

The following example generalizes the coherence calculation of Example 1, and shows that greater the commonality between vectors, greater is α .

Example 10 (Commonality). For $1 \leq i \leq m$, suppose each v_i has a common component c and an idiosyncratic component u_i , that is, $v_i = c + u_i$ with $u_i \cdot u_j = 0$ for $1 \leq j \leq m$ and $j \neq i$; $u_i \cdot c = 0$; and say, $u_i \cdot u_i = \|u\|^2$ for some u . It is easy to see that α over the uniform distribution on v_i in this case is $\frac{1}{m} [1 + (m - 1) \cdot f]$ where $f = \|c\|^2 / (\|c\|^2 + \|u\|^2)$. \square

Theorem 2 (Boundedness and Scale Invariance). *We have,*

$$0 \leq \alpha(\mathcal{V}) \leq 1$$

where $\alpha(\mathcal{V}) = 0$ iff $\mathbb{E}_{v \sim \mathcal{V}}[v] = 0$ and $\alpha(\mathcal{V}) = 1$ iff all the vectors are equal.

Furthermore, for non-zero $k \in \mathbb{R}$, we have,

$$\alpha(k\mathcal{V}) = \alpha(\mathcal{V})$$

where $k\mathcal{V}$ denotes the distribution of the random variable kv where v is drawn from \mathcal{V} .

Proof. Boundedness. Since $v \cdot v \geq 0$ for any v , we have $\mathbb{E}_{v \sim \mathcal{V}}[v \cdot v] \geq 0$. Furthermore, it is easy to verify by expanding the expectations (in terms of the vectors and their corresponding probabilities) that

$$\mathbb{E}_{v \sim \mathcal{V}, v' \sim \mathcal{V}} [v \cdot v'] = \mathbb{E}_{v \sim \mathcal{V}}[v] \cdot \mathbb{E}_{v \sim \mathcal{V}}[v] \geq 0. \quad (7)$$

⁴⁸This is also generally true in our experiments, but if necessary, a small constant can be added to the denominator to ensure this quantity is always well-defined.

Therefore, $\alpha(\mathcal{V}) \geq 0$. Likewise, another direct computation shows that

$$0 \leq \mathbb{E}_{v' \sim \mathcal{V}} \left[(\mathbb{E}_{v \sim \mathcal{V}}[v] - v') \cdot (\mathbb{E}_{v \sim \mathcal{V}}[v] - v') \right] = \mathbb{E}_{v \sim \mathcal{V}}[v \cdot v] - \mathbb{E}_{v \sim \mathcal{V}}[v] \cdot \mathbb{E}_{v \sim \mathcal{V}}[v] \quad (8)$$

Since from Equation 7 we have $\mathbb{E}[v] \cdot \mathbb{E}[v] = \mathbb{E}[v \cdot v]$, it follows that $\alpha(\mathcal{V}) \leq 1$. Furthermore, since each term of the expectation on the left is non-negative, equality is attained only when all the vectors are equal.

Scale invariance. We have,

$$\alpha(k\mathcal{V}) = \frac{\mathbb{E}_{v \sim k\mathcal{V}, v' \sim k\mathcal{V}} [v \cdot v']}{\mathbb{E}_{v \sim k\mathcal{V}} [v \cdot v]} = \frac{\mathbb{E}_{v \sim \mathcal{V}, v' \sim \mathcal{V}} [kv \cdot kv']}{\mathbb{E}_{v \sim \mathcal{V}} [kv \cdot kv]} = \frac{\mathbb{E}_{v \sim \mathcal{V}, v' \sim \mathcal{V}} [v \cdot v']}{\mathbb{E}_{v \sim \mathcal{V}} [v \cdot v]} = \alpha(\mathcal{V}) \quad (9)$$

□

Coherence of mini-batch gradients. We now prove a result that connects the α of (the uniform distribution on) a set of individual examples to that of the distribution of mini-batches constructed from those examples, by selecting k elements at random for each mini-batch with replacement. This is critical for measuring α in practical networks, and furthermore, allows us to estimate per-example α without computing per-example gradients which leads to a huge speedup in practice (see Appendix D).

Theorem 3 (Stylized mini-batching). *Let v_1, v_2, \dots, v_k be k i.i.d. variables drawn from \mathcal{V} . Let \mathcal{W} denote the distribution of the random variable $w = \frac{1}{k} \sum_{i=1}^k v_i$. We have,*

$$\alpha(\mathcal{W}) = \alpha(k\mathcal{W}) = \frac{k \cdot \alpha(\mathcal{V})}{1 + (k-1) \cdot \alpha(\mathcal{V})} \quad (10)$$

Furthermore, $\alpha(\mathcal{W}) \geq \alpha(\mathcal{V})$ with equality iff $\alpha(\mathcal{V}) = 0$ or $\alpha(\mathcal{V}) = 1$.

Proof. The first equality in (10) follows from Theorem 2. For the second equality, we have,

$$\alpha(k\mathcal{W}) = \frac{\mathbb{E}_{w \sim k\mathcal{W}, w' \sim k\mathcal{W}} [w \cdot w']}{\mathbb{E}_{w \sim k\mathcal{W}} [w \cdot w]} = \frac{\mathbb{E}_{v_1, \dots, v_k, v'_1, \dots, v'_k} [(\sum_i v_i) \cdot (\sum_i v'_i)]}{\mathbb{E}_{v_1, \dots, v_k} [(\sum_i v_i) \cdot (\sum_i v_i)]} = \frac{k^2 \mathbb{E}_{v \sim \mathcal{V}, v' \sim \mathcal{V}} [v \cdot v']}{k \mathbb{E}_{v \sim \mathcal{V}} [v \cdot v] + k \cdot (k-1) \mathbb{E}_{v \sim \mathcal{V}, v' \sim \mathcal{V}} [v \cdot v']}$$

By dividing the numerator and denominator of the last expression by $k \mathbb{E}_{v \sim \mathcal{V}} [v \cdot v]$ the required result follows.

Now, since $\alpha \leq 1$, we have $k \geq 1 + (k-1) \cdot \alpha$. Since $\alpha \geq 0$, multiplying both sides by $\frac{\alpha}{1+(k-1)\cdot\alpha}$ we have $\frac{k \cdot \alpha}{1+(k-1)\cdot\alpha} \geq \alpha$. Finally, it is easy to check that the two solutions of $\frac{k \cdot \alpha}{1+(k-1)\cdot\alpha} = \alpha$ are $\alpha = 0$ and $\alpha = 1$. □

Remark. This formulation provides a nice perspective on the type of results proved in Yin et al. [2018] and Jain et al. [2018]. When $\alpha \ll 1/k$ but non-zero (i.e., we have high gradient diversity), creating mini-batches of size k increases coherence almost k times. But, when $\alpha \approx 1$ (i.e., low diversity) there is not much point in creating mini-batches since there is little room for improvement. See also discussion of gradient diversity in Appendix B.

Coherence with some examples fitted. During training, as examples get fitted, their gradients become zero. The following lemma characterizes α when some fraction of the vectors are zero. It

shows that α decreases as larger fractions of the distribution become zero (all else remaining being equal), and therefore, as examples get fitted during training there is a natural tendency for α to decrease.

Lemma 4 (Effect of zero vectors). *If \mathcal{W} denotes the distribution where with probability $p > 0$ we pick a vector from \mathcal{V} and with probability $1 - p$ we pick the zero vector then $\alpha(\mathcal{W}) = p \cdot \alpha(\mathcal{V})$.*

Proof.

$$\alpha(\mathcal{W}) = \frac{\mathbb{E}_{w \sim \mathcal{W}, w' \sim \mathcal{W}} [w \cdot w']}{\mathbb{E}_{w \sim \mathcal{W}} [w \cdot w]} = \frac{p^2 \cdot \mathbb{E}_{v \sim \mathcal{V}, v' \sim \mathcal{V}} [v \cdot v']}{p \cdot \mathbb{E}_{v \sim \mathcal{V}} [v \cdot v]} = p \cdot \alpha(\mathcal{V}) \quad (11)$$

□

Example 11 (Coherence reduction with zero vectors). If we add k zero gradients to the collection of gradients constructed in Example 10, using Lemma 4, we get,

$$\alpha = \frac{m}{m+k} \cdot \frac{1}{m} [1 + (m-1) \cdot f] = \frac{1}{n} [1 + (n-k-1) \cdot f]$$

where $n = m+k$ is the size of this new sample. For a fixed n , as k increases, α decreases going down to $1/n$ (the orthogonal limit) when all but one vector in the sample is zero, i.e., $k = n-1$. □

Lemma 5 (Combining orthogonal distributions). *If \mathcal{W} denotes the distribution where with probability $p > 0$ we pick a vector from \mathcal{U} and with probability $1 - p$ we pick a vector from \mathcal{V} and all elements in the support of \mathcal{U} are orthogonal to those in the support of \mathcal{V} then we have*

$$\alpha(\mathcal{W}) \leq p \cdot \alpha(\mathcal{U}) + (1-p) \cdot \alpha(\mathcal{V}).$$

Proof.

$$\begin{aligned} \alpha(\mathcal{W}) &= \frac{\mathbb{E}_{w \sim \mathcal{W}, w' \sim \mathcal{W}} [w \cdot w']}{\mathbb{E}_{w \sim \mathcal{W}} [w \cdot w]} \\ &= \frac{p^2 \mathbb{E}_{u \sim \mathcal{U}, u' \sim \mathcal{U}} [u \cdot u'] + (1-p)^2 \mathbb{E}_{v \sim \mathcal{V}, v' \sim \mathcal{V}} [v \cdot v'] + 2 \cdot p \cdot (1-p) \mathbb{E}_{u \sim \mathcal{U}, v \sim \mathcal{V}} [u \cdot v]}{p \cdot \mathbb{E}_{u \sim \mathcal{U}} [u \cdot u] + (1-p) \cdot \mathbb{E}_{v \sim \mathcal{V}} [v \cdot v]} \\ &= \frac{p^2 \mathbb{E}_{u \sim \mathcal{U}, u' \sim \mathcal{U}} [u \cdot u'] + (1-p)^2 \mathbb{E}_{v \sim \mathcal{V}, v' \sim \mathcal{V}} [v \cdot v']}{p \cdot \mathbb{E}_{u \sim \mathcal{U}} [u \cdot u] + (1-p) \cdot \mathbb{E}_{v \sim \mathcal{V}} [v \cdot v]} \\ &= \frac{p^2 \mathbb{E}_{u \sim \mathcal{U}, u' \sim \mathcal{U}} [u \cdot u']}{p \cdot \mathbb{E}_{u \sim \mathcal{U}} [u \cdot u]} + \frac{(1-p)^2 \mathbb{E}_{v \sim \mathcal{V}, v' \sim \mathcal{V}} [v \cdot v']}{(1-p) \cdot \mathbb{E}_{v \sim \mathcal{V}} [v \cdot v]} \\ &\leq \frac{p^2 \mathbb{E}_{u \sim \mathcal{U}, u' \sim \mathcal{U}} [u \cdot u']}{p \cdot \mathbb{E}_{u \sim \mathcal{U}} [u \cdot u]} + \frac{(1-p)^2 \mathbb{E}_{v \sim \mathcal{V}, v' \sim \mathcal{V}} [v \cdot v']}{(1-p) \cdot \mathbb{E}_{v \sim \mathcal{V}} [v \cdot v]} \\ &= p \cdot \alpha(\mathcal{U}) + (1-p) \cdot \alpha(\mathcal{V}) \end{aligned}$$

□

Example 12 (Coherence reduction with “fitted” examples). Let U be a set with k vectors and V a set with $m - k$ vectors disjoint from U . The elements of V are pairwise orthogonal and also orthogonal to the elements of U . We may think of the elements of V as gradients of training examples that have already been “fitted” but have residual (“noisy”) gradients, and the elements of U as the gradients of training examples yet to be fitted.

Let \mathcal{U} , \mathcal{V} , and \mathcal{W} denote the uniform distributions on the sets U , V , and $U \cup V$, respectively. Applying Lemma 5, and observing that $\alpha(\mathcal{V}) = 1/(m - k)$, we get,

$$\alpha(\mathcal{W}) \leq \frac{k}{m} \alpha(\mathcal{U}) + \frac{m-k}{m} \frac{1}{m-k} = r \alpha(\mathcal{U}) + \frac{1}{m} \quad (12)$$

where $r \equiv k/m$ is the fraction of examples yet to be “fitted”. Thus as more examples get fitted, that is, as r goes to 0, the upper bound on \mathcal{W} goes towards the orthogonal limit $1/m$. \square

Lemma for stability. In order to analyze stability of gradient descent, it is useful to have a bound on the difference of two vectors in terms of α . This is used in the proof of Theorem 12 (which is used to establish Theorem 1).

Lemma 6 (Bound on Expected Difference). *We have,*

$$\mathbb{E}_{v \sim \mathcal{V}, v' \sim \mathcal{V}} [\|v - v'\|] \leq \sqrt{2 (1 - \alpha(\mathcal{V})) \mathbb{E}_{v \sim \mathcal{V}} [v \cdot v]} \quad (13)$$

Proof. Starting with Jensen’s inequality, we have,

$$\begin{aligned} \left(\mathbb{E}_{v \sim \mathcal{V}, v' \sim \mathcal{V}} [\|v - v'\|] \right)^2 &\leq \mathbb{E}_{v \sim \mathcal{V}, v' \sim \mathcal{V}} [\|v - v'\|^2] \\ &= \mathbb{E}_{v \sim \mathcal{V}, v' \sim \mathcal{V}} [(v - v') \cdot (v - v')] \\ &= 2 \left(\mathbb{E}_{v \sim \mathcal{V}} [v \cdot v] - \mathbb{E}_{v \sim \mathcal{V}, v' \sim \mathcal{V}} [v \cdot v'] \right) \\ &= 2 (1 - \alpha(\mathcal{V})) \mathbb{E}_{v \sim \mathcal{V}} [v \cdot v] \end{aligned}$$

from which the required result follows. \square

Note that in the case of perfect coherence, that is $\alpha(\mathcal{V}) = 1$, the expected norm of the difference between any two vectors is zero.

Decomposition. Finally, the following theorem shows that the overall α of a network (that is, α as measured over entire per-example gradients) is a convex combination of that of its sub-components (α as measured, say, over different layers or even individual parameters). This is useful in understanding the coherence of a network in terms of the coherences of its parts (for example, see the discussion of AlexNet in Section 6).

Theorem 7 (Decomposition). *Let V be the support of the distribution \mathcal{V} . Furthermore, let*

$$V = V_1 \oplus V_2 \oplus \cdots \oplus V_k$$

where the subspaces V_i ($1 \leq i \leq k$) are orthogonal to each other, that is, V is the orthogonal direct sum of the V_i . \mathcal{V} induces a distribution \mathcal{V}_i on each V_i . Suppose each $\alpha(\mathcal{V}_i)$ exists. Then,

$$\alpha(\mathcal{V}) = f_1 \cdot \alpha(\mathcal{V}_1) + f_2 \cdot \alpha(\mathcal{V}_2) + \cdots + f_k \cdot \alpha(\mathcal{V}_k)$$

where $f_i \equiv \mathbb{E}_{v_i \sim \mathcal{V}_i} [v_i \cdot v_i] / (\sum_{i=1}^k \mathbb{E}_{v_i \sim \mathcal{V}_i} [v_i \cdot v_i])$ and $0 \leq f_i \leq 1$ and $\sum_{i=0}^k f_i = 1$.

Proof. We have,

$$\alpha(\mathcal{V}) = \frac{\mathbb{E}_{v \sim \mathcal{V}, v' \sim \mathcal{V}} [v \cdot v']}{\mathbb{E}_{v \sim \mathcal{V}} [v \cdot v]} = \frac{\sum_{i=1}^k \mathbb{E}_{v_i \sim \mathcal{V}_i, v'_i \sim \mathcal{V}_i} [v_i \cdot v'_i]}{\sum_{i=1}^k \mathbb{E}_{v_i \sim \mathcal{V}_i} [v_i \cdot v_i]} = \frac{\sum_{i=1}^k \left(\alpha(\mathcal{V}_i) \cdot \mathbb{E}_{v_i \sim \mathcal{V}_i} [v_i \cdot v_i] \right)}{\sum_{i=1}^k \mathbb{E}_{v_i \sim \mathcal{V}_i} [v_i \cdot v_i]} = \sum_{i=1}^k f_i \cdot \alpha(\mathcal{V}_i).$$

□

B. COMPARISON OF α WITH OTHER METRICS

The coherence metric presented in Section 4 was first presented in Chatterjee and Zielinski [2020]. Here, we review some of the other metrics in the literature for studying gradient alignment, most of which were proposed concurrently. See also Section 12 for broader discussions of the papers in which these metrics were proposed.

Stiffness. Fort et al. [2020] study two variants of the average pairwise dot product that they call *sign stiffness* and *cosine stiffness*. In our notation these are

$$S_{\text{sign}} \equiv \mathbb{E}_{\substack{z \sim \mathcal{D}, z' \sim \mathcal{D} \\ z \neq z'}} [\text{sign}(g_z \cdot g_{z'})] \quad \text{and} \quad S_{\text{cos}} \equiv \mathbb{E}_{\substack{z \sim \mathcal{D}, z' \sim \mathcal{D} \\ z \neq z'}} \left[\frac{g_z}{\|g_z\|} \cdot \frac{g_{z'}}{\|g_{z'}\|} \right].$$

These are meant to capture how a small gradient step based on one input example affects the loss on a *different* input example. Although Fort et al. do not describe why they choose to transform the gradients in these specific ways, we expect it is to normalize the dot product so that it can be tracked in the course of training. In their experience, they found sign stiffness to be more useful to analyze stiffness between classes whereas cosine stiffness was more useful within a class.

Gradient Confusion. Sankararaman et al. [2020] introduce the notion of a gradient confusion bound. The *gradient confusion bound* is $\zeta \geq 0$ if for all $z, z' \in Z$ and $z \neq z'$, we have, $g_z \cdot g_{z'} \geq -\zeta$. They use this concept to study theoretically the convergence rate of gradient descent, but in their experimental results they measure the minimum cosine similarity between gradients, i.e.,

$$\min_{\substack{z \in Z, z' \in Z \\ z \neq z'}} \left[\frac{g_z}{\|g_z\|} \cdot \frac{g_{z'}}{\|g_{z'}\|} \right].$$

There are several advantages to using α and, by extension, α_m/α_m^\perp over these metrics:

- **Computational Efficiency.** For a sample of size m , due to (3), α can be computed exactly in $O(m)$ time in contrast to $O(m^2)$ time required for stiffness and cosine dot products. Furthermore, it can be computed in a streaming fashion by keeping a running sum for the numerator and a separate one for the denominator, thereby removing the need to store per-example gradients. Thus, in our experiments we are able to use sample sizes a couple of orders of magnitude higher than those in Fort et al. [2020] and Sankararaman et al. [2020]. Furthermore, Theorem 3 leads to a huge speedup in estimating per-example coherence by avoiding the need to compute per-example gradients.
- **Mathematical Simplicity.** We believe our definition is cleaner mathematically. This allows us to reason about the metric more easily. For example,
 - (1) We can show that α of minibatch gradients is greater than that of individual examples (Theorem 3). Therefore, care must be taken if minibatch gradients are used in lieu of example gradients in computing coherence (for example, as in Sankararaman et al. [2020]).

- (2) Explicitly ruling out $z \neq z'$ as in done in stiffness and cosine similarity to eliminate self-correlation is unnatural and can get tricky in practice due to near-duplicates or multiple examples leading to same or very similar gradients. We obtain meaningful values without imposing those conditions, but if one insists on removing self-correlations, then subtracting $1/m$ from α_m or 1 from α_m/α_m^\perp is a more principled way to do it.
- (3) The non-linearities in stiffness and cosine similarity amplify small per-example gradients potentially overstating their importance, and lead to a discontinuity (or undefined behavior) with zero gradients. However, we can cleanly account for the effect of negligible gradients in our observations (e.g., see Lemma 4).
- **Interpretability.** Finally, as discussed in detail above, they are normalized and yet easily interpretable due to the natural connection with loss. In contrast, the non-linearities (and to a lesser extent the $z \neq z'$ restriction) make it hard to tie stiffness or minimum cosine similarity to what happens during training; specifically, to the change in the loss function as a result of a gradient step which is the expectation over *all* per-example gradients.

GSNR and OSGR. Liu et al. [2020c] propose two metrics to study gradient alignment. The first one is *gradient signal-to-noise ratio* (GSNR) and is defined for the i th gradient coordinate g_z^i as follows

$$r^i \equiv \frac{[\mathbb{E}_{z \sim \mathcal{D}} [g_z^i]]^2}{\text{Var}_{z \sim \mathcal{D}} [g_z^i]}.$$

They measure GSNR on the training sample.

To study generalization, they define another quantity called *one-step-generalization ratio* (OSGR) denoted by $\mathbf{R}(n)$ which measures the ratio of the expected change in the test loss (on a test sample of size n) due to one step of gradient descent relative to the expected change in the training loss (on a training sample of size n). They show that under the assumption that the expected gradient over the training sample and test sample have the same distribution (the *early training assumption*), GSNR is related to OSGR (their equation 22), and confirm this with experiments on small datasets such as MNIST and CIFAR10.

In comparison to two different quantities (one over the training sample and on a per-coordinate basis, and another over both training and test and on a aggregate basis across all coordinates), we define only a single quantity that can be measured on either the test sample or the training sample, and on a per-coordinate basis, or on an aggregated basis. This leads to a simpler theory, and in our formulation the equivalent result to their equation 22 factors into using their early training assumption to argue that training and test coherence (either aggregate or per-coordinate) is equal (we also confirm this with experiments in Section 6), and then relating the aggregate quantities to the per-coordinate quantities using Theorem 7.

Furthermore, in order to understand the generalization mystery, we are interested in a bound not just for one step (as they have), or even for early training, but for *all* of training. We note here that Theorem 1 does not depend on an early training assumption such as theirs, and indeed as their experiments and ours show, that assumption does not hold later on in training.

Gradient Diversity. The reciprocal of α appears in the theory literature as *gradient diversity*. This was used by Yin et al. [2018] in theoretical bounds to understand the effect of mini-batching on convergence of SGD. (A similar result appears for least squares regression in Jain et al. [2018].) They show that the greater is the gradient diversity, the more effective are large mini-batches in speeding up SGD (see also Theorem 3 and the remark following it). Although they support their theoretical analysis with experiments on CIFAR-10 (where they replicate $1/r$ of the dataset r times and show that greater the value of r less the effectiveness of mini-batching to speed up) they never

actually measure the gradient diversity in their experiments, or further study its properties. Note that they compute gradient diversity only on the training set since they use it for reasoning about optimization only.

For our purposes, α is a better choice than $1/\alpha$ —not just because coherence rather than incoherence is what leads to generalization—but also since the latter can diverge: g can be 0 without all g_z being zero (e.g., at the end of training in an under-parameterized setting). Furthermore, a better measure of the lack of coherence is $1 - \alpha$ which is the quantity that appears in our bound in Theorem 1.

Yin et al. [2018] also define a related quantity called *differential gradient diversity* $\Delta(w, w')$ which they use to reason about the relative stability of SGD with mini-batch size b compared to SGD with mini-batch size of 1 (the case analyzed in Hardt et al. [2016]). If $g_i(w)$ denotes the gradient of the i th training example at (parameter) w , the differential gradient diversity between w and w' (for $w \neq w'$)⁴⁹ is defined as follows:

$$\Delta(w, w') \equiv \frac{\sum_{i=1}^m \|g_i(w) - g_i(w')\|^2}{\|\sum_{i=1}^m (g_i(w) - g_i(w'))\|^2}$$

where m is the size of the training set.

Their main result is that in the convex and strongly convex cases, if with high probability b is smaller than $m\Delta(w, w')$ for all $w \neq w'$, then SGD with mini-batch size b and SGD with mini-batch size 1 can be both stable in roughly the same range of step-sizes, and their generalization errors are roughly equal. (They do not consider the non-convex case.)

In contrast, our focus is on understanding generalization in the non-convex case, and for all batch sizes, including the full batch case.

Finally, when studying a form of adversarial initialization, Mehta et al. [2021] (independently) introduce a notion very similar to α except (like stiffness and cosine similarity) they explicitly require $z \neq z'$ (but see discussion on this point above under mathematical simplicity). They show experimentally on small datasets such as CIFAR-10, SHVN, etc. that within a class, their metric correlates well with generalization ability of the net when scale of initialization is varied.

C. PROOF OF THE GENERALIZATION THEOREM

We present a proof of Theorem 1 which bounds the expected generalization gap of gradient descent (GD) in terms of the coherence as measured by α during training. The proof goes through the notion of algorithmic stability [Bousquet and Elisseeff, 2002, Kearns and Ron, 1999, Devroye and Wagner, 1979], using the iterative framework first established by Hardt et al. [2016] (using uniform stability) and later extended by Kuzborskij and Lampert [2018] (using a more general notion of stability that we use as well).

Both Hardt et al. [2016] and Kuzborskij and Lampert [2018] prove stability of stochastic gradient descent in the non-convex case (that is, for neural network training) by relying solely on *early stopping*. Their arguments require that no example is seen “too often.” The guarantee provided by Hardt et al. [2016] is independent of the training dataset, and therefore, as noted before, cannot explain the generalization mystery. Indeed, their result says that stochastic gradient descent (SGD) is stable if run for a few steps with a step size that decays as $1/t$ (they summarize it under the slogan “train faster, generalize better”), but such a short run is insufficient to memorize random labels. And, if SGD is run long enough to memorize random labels, then clearly it is not stable (since if a training example with a random label is missing in the dataset, the loss on it is expected to be much higher than otherwise). The bound in Kuzborskij and Lampert [2018], while improving on that

⁴⁹Though that does not guarantee that the denominator is non-zero.

of Hardt et al. [2016], by taking into account the curvature of the loss landscape at initialization (and thus incorporating some data-dependence), comes at the price of requiring that no example be seen more than once.

Our approach towards arguing the data-dependent stability of gradient descent is fundamentally different since it is based on the coherence seen during training, and does not rely on solely on early stopping. Through coherence it is able to differentiate between the real data case and the random data case, providing different bounds in the two cases. An interesting consequence of this line of argument is that our proof applies equally well to non-stochastic gradient descent (that is, full-batch gradient descent) thus resolving an open question in Hardt et al. [2016].

Preliminaries. Let Z be a domain of examples, and let $\mathcal{D}(z)$ be a distribution on Z from which the training and test examples are sampled. Let $\ell(w, z)$ denote the loss of the model with parameters w on an example z . The goal is to find model with small *population risk*, defined as:

$$R(w) \equiv \mathbb{E}_{z \sim \mathcal{D}} \ell(w, z).$$

Let $S \equiv (z_1, \dots, z_m)$ be a sample of m examples drawn i.i.d. from \mathcal{D} . The *empirical risk* of a model w on the sample S , denoted by $\hat{R}(w, S)$ is defined as follows:

$$\hat{R}(w, S) \equiv \frac{1}{m} \sum_{i \in [m]} \ell(w, z_i).$$

where $[n]$ denotes the set $\{1, 2, \dots, n\}$.

Let A be a (possibly) randomized training algorithm. We are interested in the *expected generalization gap* of A when training with m examples from a distribution \mathcal{D} which is given by,

$$\text{gap}(\mathcal{D}, m) \equiv \mathbb{E}_{S \sim \mathcal{D}^m} \mathbb{E}_{\theta} \left[R(\mathbf{A}_\theta(S)) - \hat{R}(\mathbf{A}_\theta(S), S) \right]$$

where $\mathbf{A}_\theta(S)$ denotes the model obtained by running A on S and θ denotes the sequence of coin tosses used by A in a given run. (We omit the distribution of θ in the expression for expectation to reduce clutter.) The expected generalization gap is similar to the notion on on-average generalization in Shalev-Shwartz et al. [2010].⁵⁰

As before, let $S \equiv (z_1, \dots, z_m)$ be a sample of m examples drawn i.i.d. from \mathcal{D} . Furthermore, let $S' \equiv (z'_1, \dots, z'_m)$ be a second sample of m examples drawn i.i.d. from \mathcal{D} . The *expected stability* of A is given by,

$$\text{stab}(\mathcal{D}, m) \equiv \mathbb{E}_{S \sim \mathcal{D}^m} \mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{\theta} \left[\frac{1}{m} \sum_{i \in [m]} \left[\ell(\mathbf{A}_\theta(S^{(i)}), z_i) - \ell(\mathbf{A}_\theta(S), z_i) \right] \right]$$

where $S^{(i)} \equiv (z_1, \dots, z_{i-1}, z'_i, z_{i+1}, \dots, z_m)$, that is, $S^{(i)}$ is S with z_i replaced by z'_i .

Our notion of stability, expected stability, most closely resembles the notion of *average-RO (replace-one) stability* introduced by Shalev-Shwartz et al. [2010] but extended to randomized algorithms.⁵¹ Like their Lemma 11 (and also the first equality of Lemma 7 in Bousquet and Elisseeff [2002]), we have the following connection between expected stability and the expected generalization gap:

⁵⁰Also see Lemma 14 in Shalev-Shwartz et al. [2010] for how on-average generalization relates to a stronger notion of generalization under certain asymptotic assumptions on empirical risk minimization.

⁵¹Kuzborskij and Lampert [2018] also extend on-average stability to randomized algorithms, but they compute the supremum over i instead of the mean.

Theorem 8 (Stability equals generalization). *When training with m samples from a distribution \mathcal{D} , we have,*

$$\text{gap}(\mathcal{D}, m) = \text{stab}(\mathcal{D}, m).$$

Proof. This proof is similar to that of Lemma 11 in Shalev-Shwartz et al. [2010] and the first equality of Lemma 7 in Bousquet and Elisseeff [2002], and could be also derived from those arguments by considering \mathbf{A}_θ as a deterministic algorithm, and then taking expectations over θ .

From linearity of expectation, $\text{stab}(\mathcal{D}, m)$ is equal to

$$\mathbb{E}_{S \sim \mathcal{D}^m} \mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_\theta \left[\frac{1}{m} \sum_{i \in [m]} \ell(\mathbf{A}_\theta(S^{(i)}), z_i) \right] - \mathbb{E}_{S \sim \mathcal{D}^m} \mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_\theta \left[\frac{1}{m} \sum_{i \in [m]} \ell(\mathbf{A}_\theta(S), z_i) \right]$$

Now, from the definition of \hat{R} , we can rewrite the second term of the difference as follows:

$$\mathbb{E}_{S \sim \mathcal{D}^m} \mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_\theta \left[\frac{1}{m} \sum_{i \in [m]} \ell(\mathbf{A}_\theta(S), z_i) \right] = \mathbb{E}_{S \sim \mathcal{D}^m} \mathbb{E}_\theta [\hat{R}(\mathbf{A}_\theta(S), S)],$$

From linearity and exchangeability, we can rewrite the first term of the difference as follows:

$$\begin{aligned} \mathbb{E}_{S \sim \mathcal{D}^m} \mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_\theta \left[\frac{1}{m} \sum_{i \in [m]} \ell(\mathbf{A}_\theta(S^{(i)}), z_i) \right] &= \frac{1}{m} \sum_{i \in [m]} \mathbb{E}_{S \sim \mathcal{D}^m} \mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_\theta [\ell(\mathbf{A}_\theta(S^{(i)}), z_i)] \\ &= \frac{1}{m} \sum_{i \in [m]} \mathbb{E}_{S \sim \mathcal{D}^m} \mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_\theta [\ell(\mathbf{A}_\theta(S), z'_i)] \quad (\text{from exchangeability}) \\ &= \frac{1}{m} \sum_{i \in [m]} \mathbb{E}_{S \sim \mathcal{D}^m} \mathbb{E}_{z'_i \sim \mathcal{D}} \mathbb{E}_\theta [\ell(\mathbf{A}_\theta(S), z'_i)] \\ &= \frac{1}{m} \sum_{i \in [m]} \mathbb{E}_{S \sim \mathcal{D}^m} \mathbb{E}_\theta \mathbb{E}_{z'_i \sim \mathcal{D}} [\ell(\mathbf{A}_\theta(S), z'_i)] \\ &= \frac{1}{m} \sum_{i \in [m]} \mathbb{E}_{S \sim \mathcal{D}^m} \mathbb{E}_\theta [R(\mathbf{A}_\theta(S))] \\ &= \mathbb{E}_{S \sim \mathcal{D}^m} \mathbb{E}_\theta [R(\mathbf{A}_\theta(S))]. \end{aligned}$$

Therefore,

$$\begin{aligned} \text{stab}(\mathcal{D}, m) &= \mathbb{E}_{S \sim \mathcal{D}^m} \mathbb{E}_\theta [R(\mathbf{A}_\theta(S))] - \mathbb{E}_{S \sim \mathcal{D}^m} \mathbb{E}_\theta [\hat{R}(\mathbf{A}_\theta(S), S)] \\ &= \mathbb{E}_{S \sim \mathcal{D}^m} \mathbb{E}_\theta [R(\mathbf{A}_\theta(S)) - \hat{R}(\mathbf{A}_\theta(S), S)] \\ &= \text{gap}(\mathcal{D}, m) \end{aligned}$$

□

Smoothness Assumptions. We make two standard assumptions on the smoothness of the loss function. These assumptions are also made by Hardt et al. [2016] and Kuzborskij and Lampert [2018], and they are (recalling that $g(w, z)$ denotes $\nabla_w \ell(w, z)$):

(1) We assume that $\ell(\cdot, z)$ is L -Lipschitz and differentiable for every $z \in Z$, that is,

$$|\ell(w, z) - \ell(w', z)| \leq L \|w - w'\| \quad (14)$$

and

$$\|g(w, z)\| \leq L. \quad (15)$$

(2) We also assume that $\ell(\cdot, z)$ is β -smooth for every $z \in Z$. That is, we assume,

$$\|g(w, z) - g(w', z)\| \leq \beta \|w - w'\|. \quad (16)$$

Stability of (Stochastic) Gradient Descent. In order to analyze the stability of gradient descent, we compare what happens when we train with S and when we train with $S^{(i)}$ for T steps. Let w_0, w_1, \dots, w_T and w'_0, w'_1, \dots, w'_T denote the sequence of weights when training with S and $S^{(i)}$ respectively. Since both runs begin with the same random initialization, we have $w_0 = w'_0$. Note that w_t for $t \in [T]$ corresponds to the weight *after* the descent step at time t has taken place, and therefore, the gradients used for step t are evaluated at w_{t-1} . (As usual $[T]$ denotes the set $\{1, 2, \dots, T\}$.)

We treat both stochastic and full-batch gradient descent in the same manner, with the latter being a special case of the former, where the batch size b is equal to the number of training samples m . Each mini-batch is constructed (independently of other mini-batches) by selecting b training examples at random *without replacement* from the m training examples.⁵² Let θ denote the sequence of coin tosses made by gradient descent to select the examples in each mini-batch (these tosses are, of course, independent of the training sample). Let $I_t(\theta)$ be an indicator variable that is 1 if the i th training example is selected in the mini-batch used at time step $t \in [T]$, and 0 otherwise. Therefore,

$$\mathbb{E}_{\theta} [I_t(\theta)] = b/m \quad (17)$$

Let z_{tj} and z'_{tj} be the training examples selected in the two runs as the j th mini-batch entry at time t where $j \in [b]$ and $t \in [T]$. Let the step size at step t be η_t . From the update rule for stochastic gradient descent, we have,

$$w_t = w_{t-1} - \eta_t \frac{1}{b} \sum_{j \in [b]} g(w_{t-1}, z_{tj})$$

and

$$w'_t = w'_{t-1} - \eta_t \frac{1}{b} \sum_{j \in [b]} g(w'_{t-1}, z'_{tj})$$

Therefore,

$$(w_t - w'_t) = (w_{t-1} - w'_{t-1}) - \eta_t \frac{1}{b} \sum_{j \in [b]} [g(w_{t-1}, z_{tj}) - g(w'_{t-1}, z'_{tj})]$$

Let $\delta_t = \|w_t - w'_t\|$, and let $\Delta''_t(b) = \frac{1}{b} \sum_{j \in [b]} [g(w_{t-1}, z_{tj}) - g(w'_{t-1}, z'_{tj})]$. From the above equation, taking norms on both sides, and the triangle inequality, we have,

$$\delta_t \leq \delta_{t-1} + \eta_t \|\Delta''_t(b)\| \quad (18)$$

⁵²We choose this specific scheme to analyze since it highlights the essential similarity between the stochastic and the non-stochastic case and is simple to analyze, but this choice is not essential. Other schemes can also be analyzed similarly.

We now express this expansion of the difference at step t between the two runs *purely* in terms of the i th training example, that is, the example that is different in S and $S^{(i)}$. Intuitively, the following lemma shows that the expansion has two components: (1) one that is always present, and due to different starting points that gets magnified with the step, and (2) an occasional extra expansion if the mini-batch happens to include the i th example.

Lemma 9 (1-Step Expansion). *For $t \in [T]$, we have,*

$$\delta_t \leq (1 + \eta_t \beta) \cdot \delta_{t-1} + \frac{\eta_t I_t(\theta)}{b} \cdot \|g(w_{t-1}, z_i) - g(w_{t-1}, z'_i)\| \quad (19)$$

Proof. Define,

$$\Delta'_t(s) = \frac{1}{s} \sum_{j \in [s]} [g(w_{t-1}, z_{tj}) - g(w'_{t-1}, z_{tj})]$$

Now, if $I_t(\theta) = 0$, then $z_{tj} = z'_{tj}$ for all $j \in [b]$, and we have,

$$\Delta''_t(b) = \frac{1}{b} \sum_{j \in [b]} [g(w_{t-1}, z_{tj}) - g(w'_{t-1}, z_{tj})] = \Delta'_t(b)$$

On the other hand, if $I_t(\theta) = 1$, then there is a $k \in [b]$ s.t. $z_{tk} = z_i$ and $z'_{tk} = z'_i$, and $z_{tj} = z'_{tj}$ for all $j \in [b] \setminus \{k\}$. Without loss of generality assume that $k = b$. Therefore, in this case,

$$\Delta''_t(b) = \frac{1}{b} (g(w_{t-1}, z_i) - g(w'_{t-1}, z'_i)) + \frac{b-1}{b} \Delta'_t(b-1)$$

Putting the $I_t(\theta) = 0$ and the $I_t(\theta) = 1$ cases together, we get,

$$\Delta''_t(b) = I_t(\theta) \left(\frac{1}{b} (g(w_{t-1}, z_i) - g(w'_{t-1}, z'_i)) + \frac{b-1}{b} \Delta'_t(b-1) \right) + (1 - I_t(\theta)) \Delta'_t(b)$$

Taking norms, and applying the triangle inequality, we get,

$$\|\Delta''_t(b)\| \leq I_t(\theta) \left(\frac{1}{b} \|g(w_{t-1}, z_i) - g(w'_{t-1}, z'_i)\| + \frac{b-1}{b} \|\Delta'_t(b-1)\| \right) + (1 - I_t(\theta)) \|\Delta'_t(b)\|$$

Now, from the smoothness condition on the loss functions (16), we have $\|\Delta'_t(s)\| \leq \beta \delta_t$ for any $s \in [b]$, and this simplifies to

$$\|\Delta''_t(b)\| \leq I_t(\theta) \left(\frac{1}{b} \|g(w_{t-1}, z_i) - g(w'_{t-1}, z'_i)\| + \frac{b-1}{b} \beta \delta_t \right) + (1 - I_t(\theta)) \beta \delta_t.$$

Now, since

$$\begin{aligned} \|g(w_{t-1}, z_i) - g(w'_{t-1}, z'_i)\| &= \|g(w_{t-1}, z_i) - g(w_{t-1}, z'_i) + g(w_{t-1}, z'_i) - g(w'_{t-1}, z'_i)\| \\ &\leq \|g(w_{t-1}, z_i) - g(w_{t-1}, z'_i)\| + \|g(w_{t-1}, z'_i) - g(w'_{t-1}, z'_i)\| \\ &\leq \|g(w_{t-1}, z_i) - g(w_{t-1}, z'_i)\| + \beta \delta_t \end{aligned}$$

we get,

$$\begin{aligned} \|\Delta''_t(b)\| &\leq I_t(\theta) \left(\frac{1}{b} (\|g(w_{t-1}, z_i) - g(w_{t-1}, z'_i)\| + \beta \delta_t) + \frac{b-1}{b} \beta \delta_t \right) + (1 - I_t(\theta)) \beta \delta_t \\ &= I_t(\theta) \left(\frac{1}{b} (\|g(w_{t-1}, z_i) - g(w_{t-1}, z'_i)\|) + \beta \delta_t \right) + (1 - I_t(\theta)) \beta \delta_t \\ &= I_t(\theta) \left(\frac{1}{b} (\|g(w_{t-1}, z_i) - g(w_{t-1}, z'_i)\|) \right) + \beta \delta_t \end{aligned}$$

Recall from (18) that

$$\delta_{t+1} \leq \delta_t + \eta_t \|\Delta''_t(b)\|$$

Therefore, we have,

$$\delta_t \leq \delta_{t-1} + \eta_t \left(I_t(\theta) \left(\frac{1}{b} (\|g(w_{t-1}, z_i) - g(w_{t-1}, z'_i)\|) \right) + \beta \delta_{t-1} \right)$$

simplifying which the required result follows. \square

Next, we unroll the recursion to bound the difference in weights at the end of training δ_T in terms of the differences in gradients in S and $S^{(i)}$ at each step. (In what follows we use the standard convention for \prod that an empty product is equal to 1.)

Lemma 10 (Unrolling recursion). *We have,*

$$\delta_T \leq \sum_{t \in [T]} \left(\frac{\eta_t I_t(\theta)}{b} \cdot \|g(w_{t-1}, z_i) - g(w_{t-1}, z'_i)\| \cdot \prod_{k=t+1}^T (1 + \eta_k \beta) \right) \quad (20)$$

Proof. By unrolling (19) and using the fact that $\delta_0 = 0$. \square

To reduce clutter, it is convenient to define a couple of abbreviations:

$$[\eta_k \beta]_{k=t_0}^{t_1} \equiv \prod_{k=t_0}^{t_1} (1 + \eta_k \beta) \quad \text{and} \quad \bar{g}(w_t) \equiv \sqrt{\mathbb{E}_{z \sim \mathcal{D}} [\|g(w_t, z)\|^2]}.$$

Intuitively, $[\eta_k \beta]_{k=t_0}^{t_1}$ denotes the expansion incurred from time step t_0 to t_1 , and $\bar{g}(w_t)$ is the average norm of per example gradients at w_t (roughly speaking). Note that by the L -Lipschitz condition on $\ell(\cdot, z)$, that is, by (15), we have,

$$\bar{g}(w_t) \leq L. \quad (21)$$

Next, we bound the difference between examples by the coherence as measured by α .

Lemma 11. *We have,*

$$\mathbb{E}_{S \sim \mathcal{D}^m} \mathbb{E}_{z'_i \sim \mathcal{D}} \mathbb{E}_\theta \left[\left| \ell(\mathbf{A}_\theta(S^{(i)}), z_i) - \ell(\mathbf{A}_\theta(S), z_i) \right| \right] \leq \frac{L^2}{m} \sum_{t \in [T]} [\eta_k \beta]_{k=t+1}^T \cdot \eta_t \cdot \sqrt{2(1 - \alpha(w_{t-1}))}$$

Proof. Since, by definition, $w_T = \mathbf{A}_\theta(S)$, $w'_T = \mathbf{A}_\theta(S^{(i)})$, and $\delta_t = \|w_T - w'_T\|$, applying the Lipschitz condition on ℓ (14), and then applying Lemma 10 we get,

$$\begin{aligned} \left| \ell(\mathbf{A}_\theta(S^{(i)}), z_i) - \ell(\mathbf{A}_\theta(S), z_i) \right| &\leq L \delta_T \\ &= L \sum_{t \in [T]} \left(\frac{\eta_t I_t(\theta)}{b} \cdot [\eta_k \beta]_{k=t+1}^T \cdot \|g(w_{t-1}, z_i) - g(w_{t-1}, z'_i)\| \right) \end{aligned}$$

Now, taking expectations on both sides, and expanding into the sum, we get the following for each term of the sum:

$$\begin{aligned}
& \mathbb{E}_{S \sim \mathcal{D}^m} \mathbb{E}_{z'_i \sim \mathcal{D}} \mathbb{E}_{\theta} \left[[\eta_k \beta]_{k=t+1}^T \cdot \frac{\eta_t I_t(\theta)}{b} \cdot \|g(w_{t-1}, z_i) - g(w_{t-1}, z'_i)\| \right] \\
&= [\eta_k \beta]_{k=t+1}^T \cdot \left(\frac{\eta_t}{b} \mathbb{E}_{\theta} [I_t(\theta)] \right) \cdot \left(\mathbb{E}_{S \sim \mathcal{D}^m} \mathbb{E}_{z'_i \sim \mathcal{D}} \|g(w_{t-1}, z_i) - g(w_{t-1}, z'_i)\| \right) \\
&\leq [\eta_k \beta]_{k=t+1}^T \cdot \left(\frac{\eta_t}{m} \right) \cdot \sqrt{2 (1 - \alpha(w_{t-1}))} \cdot \bar{g}(w_{t-1}) \\
&\leq [\eta_k \beta]_{k=t+1}^T \cdot \left(\frac{\eta_t}{m} \right) \cdot \sqrt{2 (1 - \alpha(w_{t-1}))} \cdot L
\end{aligned}$$

The second step follows from the first since θ is independent of S and z'_i . The third step follows from $\mathbb{E}_{\theta} [I_t(\theta)] = \frac{b}{m}$ (see (17)), and from Lemma 6 applied to bound the difference in the gradients for z_i and z'_i by the coherence of the per-example gradients at w_t . The last step follows from (21).

Finally, by summing up the bounds for each term and rearranging, we get the required result. \square

A straightforward consequence of the above lemma is a bound on the stability of gradient descent.

Theorem 12 (Stability Theorem).

$$|\text{stab}(\mathcal{D}, m)| \leq \frac{L^2}{m} \sum_{t \in [T]} [\eta_k \beta]_{k=t+1}^T \cdot \eta_t \cdot \sqrt{2 (1 - \alpha(w_{t-1}))} \quad (22)$$

Proof. This is a straightforward consequence of Jensen's inequality and Lemma 11.

$$\begin{aligned}
|\text{stab}(\mathcal{D}, m)| &= \left| \mathbb{E}_{S \sim \mathcal{D}^m} \mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{\theta} \left[\frac{1}{m} \sum_{i \in [m]} [\ell(\mathbf{A}_{\theta}(S^{(i)}), z_i) - \ell(\mathbf{A}_{\theta}(S), z_i)] \right] \right| \\
&\leq \frac{1}{m} \sum_{i \in [m]} \mathbb{E}_{S \sim \mathcal{D}^m} \mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{\theta} \left[|\ell(\mathbf{A}_{\theta}(S^{(i)}), z_i) - \ell(\mathbf{A}_{\theta}(S), z_i)| \right] \\
&= \frac{1}{m} \sum_{i \in [m]} \mathbb{E}_{S \sim \mathcal{D}^m} \mathbb{E}_{z'_i \sim \mathcal{D}} \mathbb{E}_{\theta} \left[|\ell(\mathbf{A}_{\theta}(S^{(i)}), z_i) - \ell(\mathbf{A}_{\theta}(S), z_i)| \right] \\
&\leq \frac{L^2}{m} \sum_{t \in [T]} [\eta_k \beta]_{k=t+1}^T \cdot \eta_t \cdot \sqrt{2 (1 - \alpha(w_{t-1}))}
\end{aligned}$$

\square

Remark. Note that in the penultimate step we bound the following quantity (using Lemma 11):

$$\frac{1}{m} \sum_{i \in [m]} \mathbb{E}_{S \sim \mathcal{D}^m} \mathbb{E}_{z'_i \sim \mathcal{D}} \mathbb{E}_{\theta} \left[|\ell(\mathbf{A}_{\theta}(S^{(i)}), z_i) - \ell(\mathbf{A}_{\theta}(S), z_i)| \right]$$

which corresponds to a stronger notion of stability than the notion of expected stability that we use to link to generalization. Indeed, it is closer to the notion of *point-wise hypothesis stability* [Bousquet and Elisseeff, 2002, Kearns and Ron, 1999] with the main differences being that it is formulated in terms of replace-one (instead of leave-one-out), for a randomized setting (instead of a deterministic one), and does not assume symmetry in the training examples. Therefore, it may

be possible to extend our results to provide a concentration guarantee along the lines of Theorem 11 in Bousquet and Elisseeff [2002].

We are now ready to prove our main theorem.

Theorem 1 (Generalization Theorem).

$$|\text{gap}(\mathcal{D}, m)| \leq \frac{L^2}{m} \sum_{t \in [T]} [\eta_k \beta]_{k=t+1}^T \cdot \eta_t \cdot \sqrt{2(1 - \alpha(w_{t-1}))} \quad (23)$$

Proof. This is a straightforward consequence of Theorem 8 and Theorem 12. \square

We now consider special cases for two common learning rate schedules. These allow us to replace the expansion term in the bound with an exponential.

Corollary 13. *If the step sizes are fixed, that is, $\eta_t = \eta$ then*

$$|\text{gap}(\mathcal{D}, m)| \leq \frac{L^2 \eta}{m} \sum_{t \in [T]} \exp((T-t) \eta \beta) \cdot \sqrt{2(1 - \alpha(w_{t-1}))} \quad (24)$$

Proof. For any $t \geq 0$,

$$[\eta_k \beta]_{k=t+1}^T \leq \exp \left(\sum_{k=t+1}^T \eta_k \beta \right) \leq \exp \left(\eta \beta \sum_{k=t+1}^T 1 \right) \leq \exp(\eta \beta(T-t))$$

where the first inequality follows from $(1+x) \leq \exp(x)$ for any x . The required result follows substituting this bound on the expansion term in Theorem 1. \square

Remark. Note that the bound in Corollary 13 may be simplified further to

$$|\text{gap}(\mathcal{D}, m)| \leq \frac{L^2 \eta \exp(T \eta \beta)}{m} \sum_{t \in [T]} \sqrt{2(1 - \alpha(w_{t-1}))}. \quad (25)$$

Corollary 14. *If we assume as in Hardt et al. [2016] that step sizes decay linearly, that is, for some $\eta > 0$ we have $\eta_t \leq \eta/t$ then*

$$|\text{gap}(\mathcal{D}, m)| \leq \frac{L^2 \eta T^{\eta \beta}}{m} \sum_{t \in [T]} \sqrt{2(1 - \alpha(w_{t-1}))} \quad (26)$$

Proof. For any $t \geq 0$,

$$[\eta_k \beta]_{k=t+1}^T \leq \exp \left(\sum_{k=t+1}^T \eta_k \beta \right) \leq \exp \left(\eta \beta \sum_{k=t+1}^T 1/t \right) \leq \exp(\eta \beta \log T) = T^{\eta \beta}$$

where, as in the previous proof, the first inequality follows from $(1+x) \leq \exp(x)$ for any x . The required result follows substituting this bound in Theorem 1 and from noticing that $\eta_t \leq \eta$ for all t . \square

D. METHODS TO MEASURE α

The Direct and Imputed Methods. In theory, α at any point in training can be directly measured on a test or training sample by computing per-example gradients, and computing it from the definition (3):

$$\alpha \equiv \frac{\mathbb{E}_{z \sim \mathcal{D}} [g_z] \cdot \mathbb{E}_{z \sim \mathcal{D}} [g_z]}{\mathbb{E}_{z \sim \mathcal{D}} [g_z \cdot g_z]}$$

where \mathcal{D} is the empirical distribution of the sample. α can be computed without storing per-example gradients by keeping separate running sums for $\mathbb{E}_{z \sim \mathcal{D}} [g_z]$ and for $\mathbb{E}_{z \sim \mathcal{D}} [g_z \cdot g_z]$. We call this method the *direct* method.

In practice, however, for most modern networks, the presence of batch normalization layers [Ioffe and Szegedy, 2015] makes it challenging to compute (or even define!) per-example gradients, since the gradient of an example depends on all the other examples in the mini-batch, and one cannot calculate a gradient of a single example in isolation from the others. One way to get around this is to estimate per-example gradients in evaluation (or inference) mode by using the moving averages of the population for normalization which disentangles the gradient of an example from the others in the batch. Although we used that technique in our preliminary work [Chatterjee and Zielinski, 2020], we found two problems with that approach:

- (1) **Lag in the moving averages.** During training, particularly, when the loss changes rapidly, the moving averages fall behind resulting in highly inflated coherence measurements.
- (2) **Incorrect gradients.** With batch normalization, an example contributes not just to its own loss, but also to the losses of the other examples through its contribution to the batch statistics. By using the population statistics instead of batch statistics, we ignore this component of the gradient of the example which leads to a disconnect between the gradients used to compute coherence and the ones actually used to update weights which is unsatisfactory.

Even though the first problem can be addressed (at an added compute cost) by updating the moving averages with the population statistics before computing the per-example gradients, there does not appear to be an easy solution for the second.⁵³

Therefore, in our experiments, we handle batch normalization in a different manner by appealing to Theorem 3 which relates the coherence of mini-batch gradients $\alpha(\mathcal{W})$ to the coherence of per-example gradients $\alpha(\mathcal{V})$. We invert equation (10) to get:

$$\alpha(\mathcal{V}) = \frac{\alpha(\mathcal{W})}{k - (k - 1) \cdot \alpha(\mathcal{W})} \tag{27}$$

where k is the batch size.

In this method, which we call the *imputed* method, to compute per-example coherence for m examples, we compute the gradient of each batch in the usual manner for training, and then compute the coherence between the m/k *batch gradients* from the definition. Finally, we *impute* a per-example coherence using (27). When computing the corresponding α_m/α_m^\perp for the imputed per-example gradients we use m for the sample size (rather than m/k).

Note that Theorem 3 is proven under the assumption that each mini-batch is constructed by picking examples at random *with* replacement, although in practice, for efficiency, mini-batches are constructed at random *without* replacement. However, in the case of networks that do not use batch normalization, when we compare the per-example coherences computed using the direct and

⁵³Note: this is indeed the main research challenge in replacing batch normalization with a technique that would work even for batches of size 1.

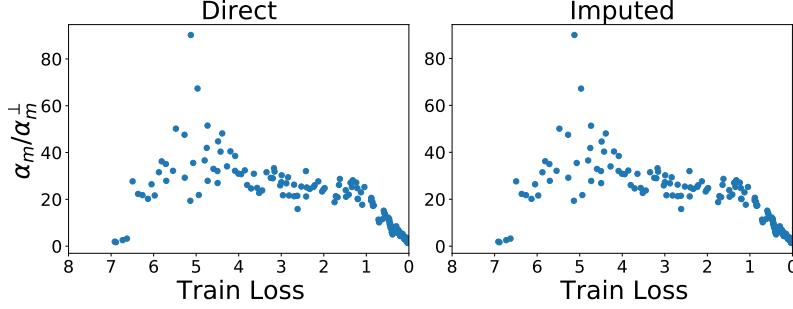


FIGURE 10. α computed using direct and imputed methods on AlexNet (without batch normalization) for ImageNet. The methods produce virtually identical results.

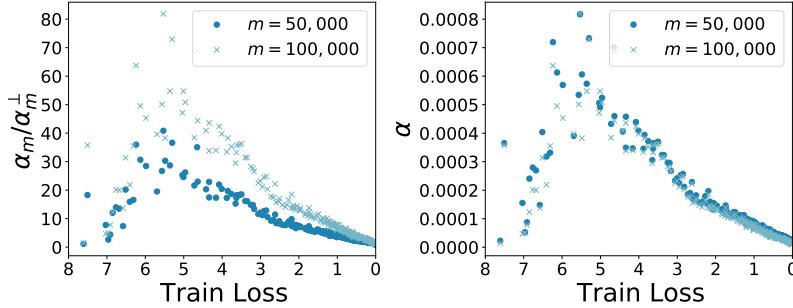


FIGURE 11. α_m and α_m/α_m^\perp on training samples of sizes $m = 50,000$ and $m = 100,000$ for a ResNet-50 being trained on ImageNet. The exact size of the sample used to measure α does not matter much if the sample is large enough. However, as discussed in Section 4, the smaller sample produces a (slightly) greater estimate of α than the larger sample. To a first order, α_m/α_m^\perp simply scales with the sample size (since it is $m\alpha$), but note that at the end of training each sample reaches its (respective) orthogonal limit ($\alpha_m/\alpha_m^\perp = 1$).

the imputed methods, we find them to be in close agreement (for example, see Figure 10). In our experiments we use the direct method wherever possible (that is, for networks without batch normalization) and resort to the imputed method where needed (for batch normalization). Note though that the imputed method is much more efficient than the direct method since it does not require the computation of per-example gradients.

Effect of sample size used for estimation. Please see Figure 11.

E. MEASURING α ON ADDITIONAL DATASETS AND ARCHITECTURES

We measured coherence during training on the test set and a training sample (of the same size as test) for the following networks and datasets:

- (1) Simple feed forward network with a single hidden layer of 2048 nodes on MNIST [LeCun and Cortes, 2010] with random labels as a source of noise (Figure 12).
- (2) A simple convolutional network (AlexNet [Krizhevsky et al., 2012]) on CIFAR-10 [Krizhevsky et al., 2009] with random labels as a source of noise (Figure 13).

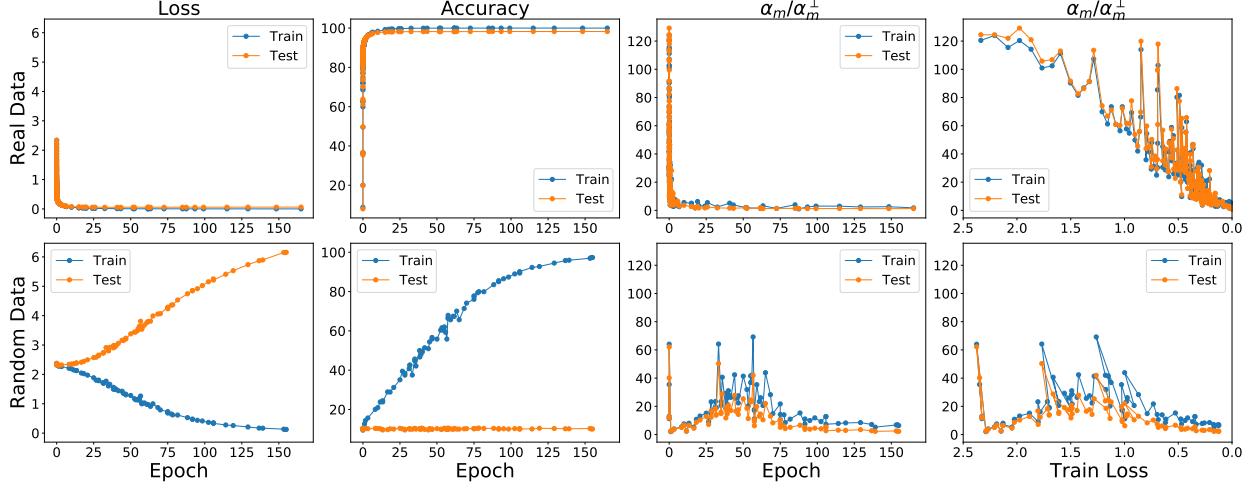


FIGURE 12. The evolution of alignment of per-example gradients during training as measured with α_m/α_m^\perp on samples of size $m = 10,000$ on MNIST dataset. The model is a simple feed-forward network with a single hidden layer containing 2048 neurons. We use SGD with a constant learning rate of 0.1. No regularizers were used.

- (3) AlexNet on ImageNet [Deng et al., 2009] with random pixels (Figure 2) and with random labels (Figure 15).
- (4) A deep residual convolutional network (ResNet-50 [He et al., 2016]) on ImageNet with random pixels (Figure 1) and with random labels (Figure 14).

Each network was trained on the original dataset (0% noise), as well as the corresponding dataset with either the labels or pixels randomized (100% noise). In all cases, we used vanilla stochastic gradient descent with a fixed learning rate. The batch size, learning rate, and the total number of steps was chosen so as to reach almost perfect training accuracy for random data in a reasonable amount of training time. Since our goal is to study the implicit regularization of gradient descent, no explicit regularizers (such as dropout, weight-decay, or input augmentation) were used.

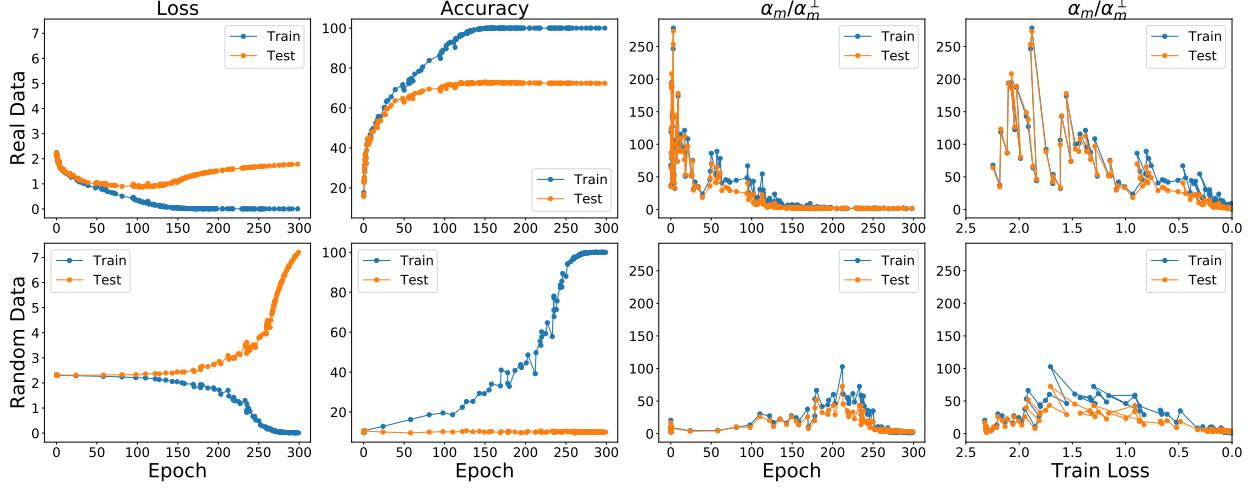


FIGURE 13. The evolution of alignment of per-example gradients during training as measured with α_m/α_m^\perp on samples of size $m = 10,000$ on CIFAR10 dataset. The model is a AlexNet architecture without dropout layer. We train for 300 epochs with SGD using constant learning rate of 0.005. No regularizers were used. Input images were rescaled to 224x224. Additional runs can be found in Figure 23.

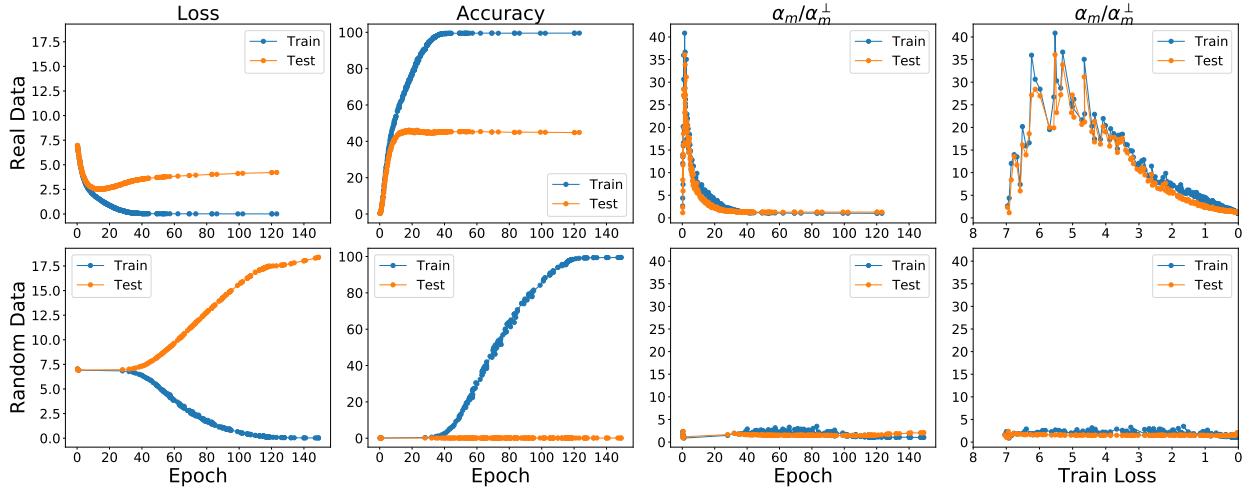


FIGURE 14. The evolution of alignment of per-example gradients during training as measured with α_m/α_m^\perp on samples of size $m = 50,000$ on ImageNet dataset. Noise was added through labels randomization. The model is a Resnet-50. Additional runs can be found in Figure 24.

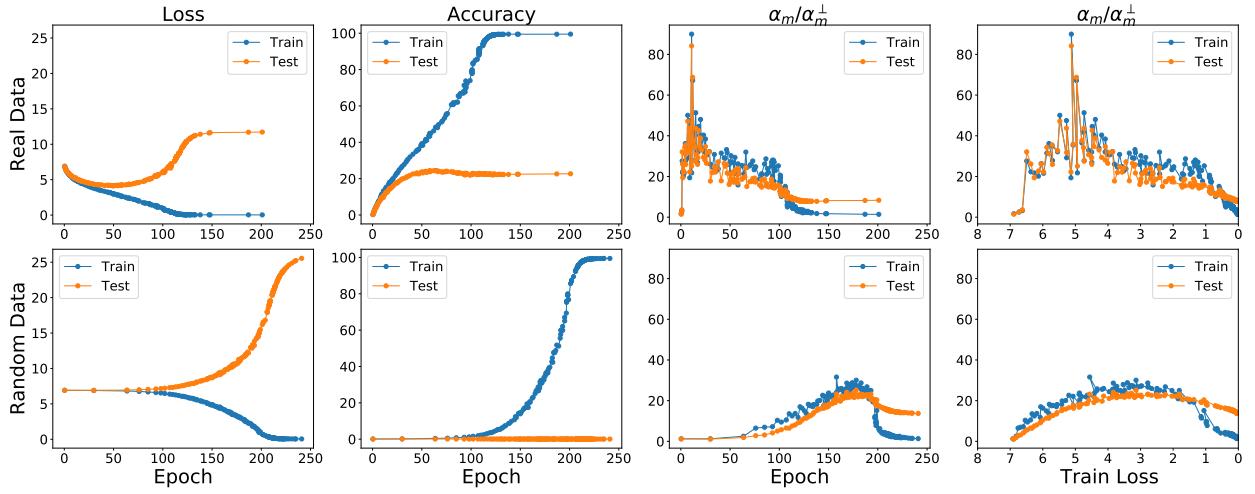


FIGURE 15. The evolution of alignment of per-example gradients during training as measured with α_m/α_m^\perp on samples of size $m = 50,000$ on ImageNet dataset. Noise was added through labels randomization. The model is an AlexNet architecture without dropout layer. Additional runs can be found in Figure 25.

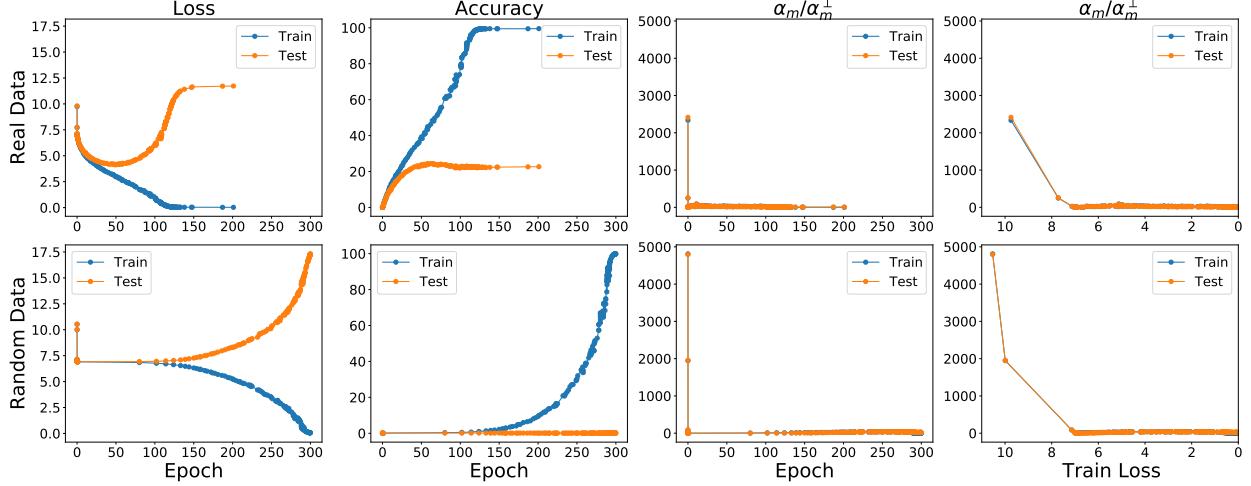


FIGURE 16. In the first few steps of training α_m/α_m^\perp can be very high due to imperfect initialization and in a sample of 50,000 training examples, each example helps thousands of others, even in case of 100% noise. In our usual plots of α_m/α_m^\perp , we omit this initial transient. This plot shows the transient when training AlexNet on ImageNet (Figure 2).

F. THE EVOLUTION OF COHERENCE

As mentioned in Section 6, experiments across several architectures and datasets show a common pattern in how coherence as measured by α_m/α_m^\perp (or equivalently, α) changes during training. Ignoring the initial transient in the first few steps of training (which we discuss a little later), coherence follows a roughly parabolic pattern: It starts off at a low value, rises to a peak, and then comes back down to the orthogonal limit.⁵⁴ This happens regardless of whether the dataset is random or real, indicating that this is an optimization (as opposed to a generalization) effect.

This parabolic evolution of coherence is noisy with many minor fluctuations and deviations from the broad pattern. This suggests that the dynamics of coherence can be understood as being governed by two competing forces during training: one that causes coherence to increase (“creation”) and one that causes coherence to decrease (“consumption”).

Of the two, coherence consumption is easier to explain. If we assume that as examples get fitted during training, their gradients either become negligible, or orthogonal to the remaining (unfitted) examples, then it can be shown that as the fraction of “fitted examples” increases, α must decrease. For more details see Lemma 4 and Lemma 5 and associated examples in Appendix A.

Coherence creation, on the other hand, is harder to explain. One way to reason about it is through the decomposition of the overall coherence in terms of the coherence of the constituent parts. As shown in Theorem 7, α as measured over the entire network (that is, over the entire gradient vector) is a convex combination of the α_i of its constituent parts (projections of the gradient vector). In particular, we have that

$$\alpha = \sum_i f_i \alpha_i$$

where α_i is the coherence of the i th parameter of the network, and $0 \leq f_i \leq 1$, and $\sum_i f_i = 1$. Given this decomposition, there are two (non mutually exclusive) ways in which the overall coherence α

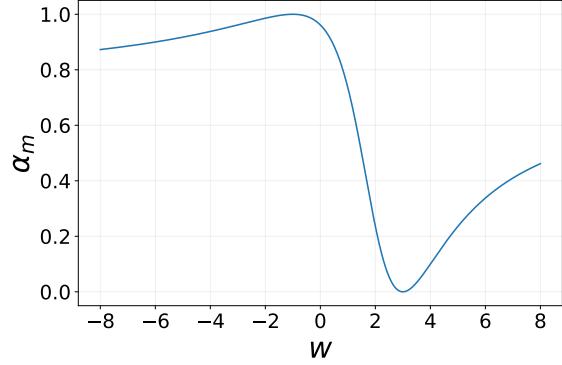
⁵⁴In rare cases, such as a fully connected network on MNIST where the signal is strong and easy to find, coherence starts off high.

can increase during training: (1) there is a shift in weight from a parameter i with low α_i to one with higher α_i , or (2) α_i for some parameter i increases.

Example 7 provides a good example for first. Observe that the α for each $u^{(i)}$ and $w^{(i)}$ is constant during training, but the overall α increases. This happens due to shift in weight from the other idiosyncratic components to $u^{(6)}$ and $w^{(6)}$ due to amplification. The following example demonstrates the second, and in fact demonstrates how α can follow a parabolic trajectory.

Example 13 (Single parameter trajectory). Consider the task of fitting a linear model $y = w \cdot x$ with a single parameter $w \in \mathbb{R}$ to the $m = 2$ points (x_1, y_1) and (x_2, y_2) under the usual mean squared loss. It is easy to check that the point of minimum α_m is when the per-example gradients are equal in magnitude and opposite in direction, that is, at $w^* = (y_1 x_1 + y_2 x_2) / (x_1^2 + x_2^2)$. (Note that this is also the w that optimizes the average loss.) The point of maximum α_m is when the per-example gradients are equal in both magnitude and direction, that is, at $w^\dagger = (y_1 x_1 - y_2 x_2) / (x_1^2 - x_2^2)$.

Now let $(x_1, y_1) = (2, 3)$ and $(x_2, y_2) = (1, 9)$. Therefore, $w^* = 3$ and $w^\dagger = -1$. If we plot α_m as a function of w for this training set, we see the following:



If gradient descent started is initialized at a point to the left of w^\dagger , along its trajectory, we would observe an initial increase in α_m as the gradients of the two examples equalized in magnitude. Once they are equal and α_m is 1, coherence begins to decrease until the two gradients are equal and opposite at the final solution w^* and α_m is 0. Note that the overall trajectory of α_m is roughly parabolic. \square

Coherence in first few steps of training. As mentioned in Section 6, very early in training, α_m/α_m^\perp can be very high even for random data due to imperfect initialization. All the training examples are coordinated in moving the network to a more reasonable point in parameter space. As may be expected from our theory, this movement generalizes well: the test loss decreases in concert with training loss in this period. Rapid changes to the network early in training is well documented (see for example, the need for learning rate warmup in He et al. [2016] and Goyal et al. [2017]).

Figure 16 shows an example of high coherence in the first few steps of training for AlexNet (in the same experiment as Figure 2). We observe that transient in the first few steps when the loss is above what might be expected from an uniform distribution, that is, when the loss is above $\log(1000) \approx 6.9$ (recall that ImageNet has 1000 labels).

G. EXPERIMENTAL DETAILS OF EASY AND HARD EXAMPLES

Hardness (difficulty) of examples. We estimate the difficulty of each of the 1,281,167 examples from the ImageNet training dataset by training a ResNet-50 to 50% training accuracy, and recording for each example, whether it was correctly classified or not. We repeat this 8 times from different random initializations, and assign a *hardness score* to each example by counting how many times it failed to be correctly classified. A score of 0 indicates that the example was correctly classified on all 8 runs (a “super easy” example) whereas a score of 8 indicates it was never correctly classified (“super hard”). The distribution of the hardness score is as follows:

hardness score	count
8 (super hard)	248,551
7	121,290
6	101,469
5	95,780
4	97,069
3	103,182
2	118,701
1	142,939
0 (super easy)	252,186

We construct the Easy ImageNet dataset described in Section 8 by randomly picking 550K examples that have a score of 3 or lower, and randomly allocating them into train (500K examples) and test (50K examples) subsets. Hard ImageNet is constructed similarly by randomly picking 550K examples that have a score of 5 or greater. (Examples with a score of 4 are discarded.)

In situ measurements. We can track coherence on examples of different levels of difficulty during regular ImageNet training.⁵⁵ Figure 17 shows the coherence as measured by α_m/α_m^\perp for 50K training samples of super easy and super hard examples during regular ImageNet training. As expected from our theory, the peak coherence of super easy examples is much higher than those of super hard examples.

Whether an ImageNet example is learned early in training or not does not appear to be architecture specific. Figure 18 tracks the loss, accuracy, and α_m/α_m^\perp for (samples of) the super easy and super hard examples (from ResNet-50) during *AlexNet* training on ImageNet. We see that the super hard examples are learned later than the super easy examples by AlexNet as well, and the coherence of super easy examples is higher than that of the super hard examples. But the difference in α_m/α_m^\perp is less stark for AlexNet than for ResNet-50. This mirrors the situation with real and random data analyzed in Section 6, and a more granular per-layer plot of α_m/α_m^\perp may reveal more significant per-layer differences.

Do intrinsic easy patterns exist? Since AlexNet seems to find the same examples to be easy (or hard) as ResNet-50, one may wonder if there is something intrinsic to an example that controls its difficulty.

To test this we train AlexNet on a training set containing either a single super easy example or a single super hard example, and count how many steps it takes for the training loss to fall below

⁵⁵One issue with tracking coherence on two different subsets of the training set is that the batch normalization statistics in each subset set may be different than those in the overall training set since the distribution of examples is different, which may impact the coherence measurement.

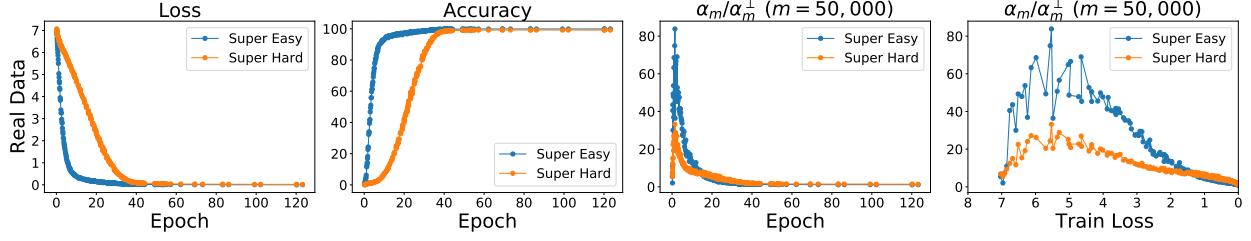


FIGURE 17. Loss, accuracy, and α_m/α_m^\perp on 50K samples of super easy and super hard examples measured during (normal) ImageNet training. As expected, the super easy examples have significantly higher coherence than super hard examples. The network used is a ResNet-50.

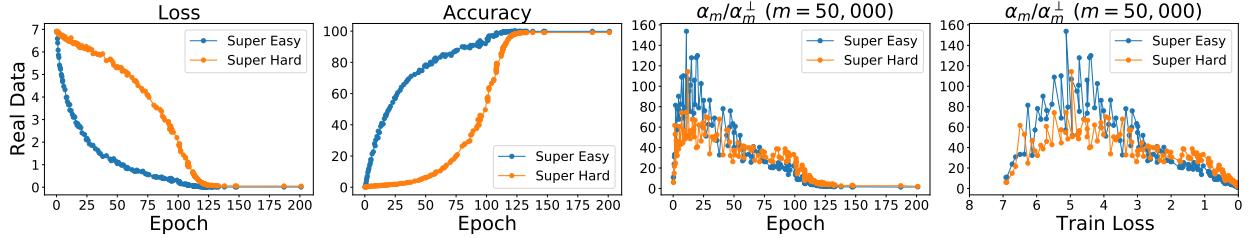


FIGURE 18. Loss, accuracy, and α_m/α_m^\perp on 50K samples of super easy and super hard examples measured during (normal) ImageNet training. This time the network being measured is an AlexNet although the hardness of the examples was measured on ResNet-50. The examples that are easy for ResNet-50 are also easy for AlexNet, suggesting that hardness is not very closely tied to the specific architecture used.

0.05.⁵⁶ We repeat this for 500 super easy and 500 super hard examples picked at random. Learning a single super hard example took on average 40.76 training steps⁵⁷ which was less than the average for a super easy example which was 43.12 steps! However, the difference between the two classes was not statistically significant.

This appears to suggest that:

Hardness is not an intrinsic property of examples, but only emerges via interaction with other examples in the training set.

This is a stronger statement than the one starting Section 8 that “easy examples are those examples that have a lot in common with other examples where commonality is measured by the alignment of the gradients of the examples” since the latter does not exclude the additional possibility of

⁵⁶Since AlexNet does not use batch normalization, we can actually train with a single example.

⁵⁷We use a very low learning rate of 0.00001 for this experiment.

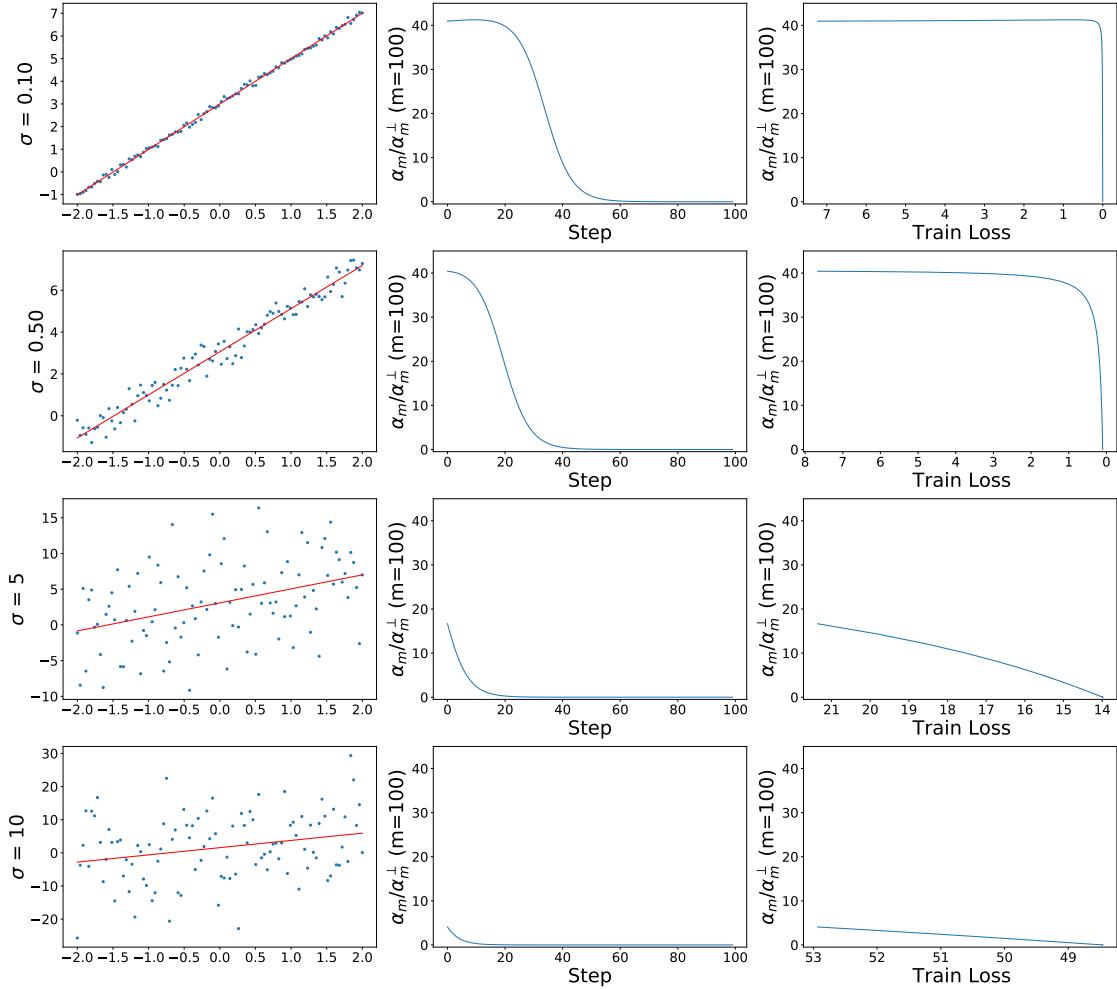
something intrinsic to examples that controls their hardness. The above experiment seems to rule out this possibility.

H. THE UNDER-PARAMETERIZED CASE: A PRELIMINARY LOOK

Coherence. Although our theory has been motivated by the over-parameterized case, one may wonder what coherence looks like in an under-parameterized setting.⁵⁸ To that end, in this section, we take a look at coherence in the case of a 2D linear regression example. Consider a sample of $m = 100$ points (x, y) where x is chosen uniformly at random from $[-2.0, 2.0]$, ϵ from $\mathcal{N}(0, \sigma^2)$, and y is given by

$$y = 2x + 3 + \epsilon.$$

We fit the linear model $y = \beta_1 x + \beta_0$ to this sample using full-batch gradient descent and measure the alignment of per-example gradients α_m/α_m^\perp on the sample:



Each row above corresponds to a different value of σ . The first panel shows a plot of the training sample and the line learned by gradient descent at the end of training. The second and third panel show α_m/α_m^\perp as a function of the number of training steps taken and the training loss respectively.

⁵⁸Note that under-parameterized systems may have good generalization, that is small gap between test and training loss, even if coherence is low, and this is qualitatively reflected in the $1/m$ dependence of the bound in Theorem 1.

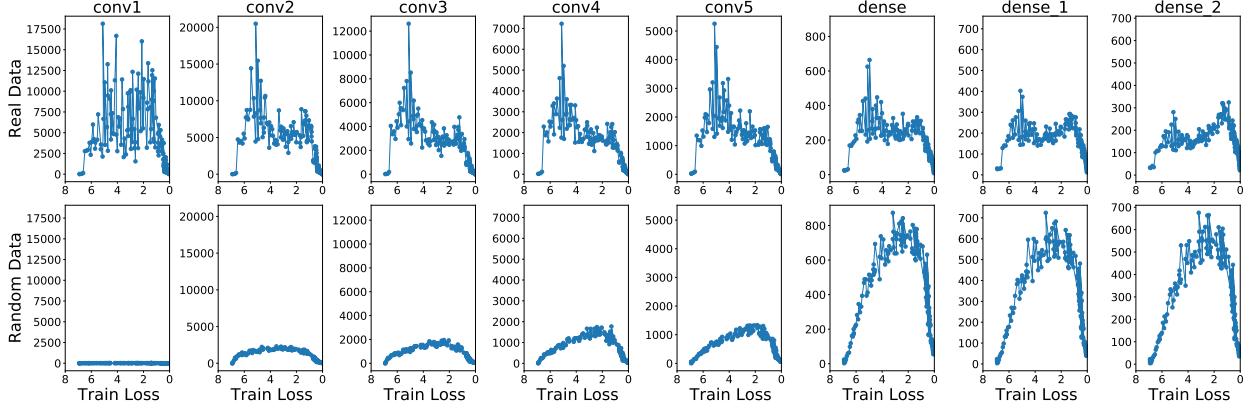
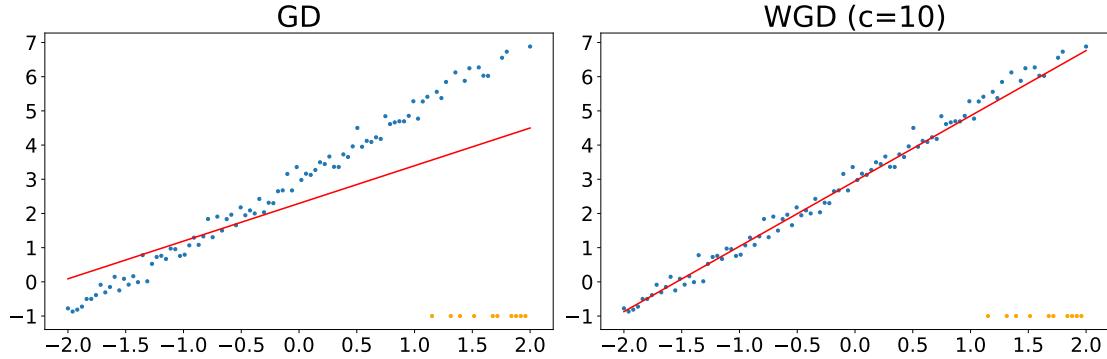


FIGURE 19. Figure 3 but with α_m/α_m^\perp measured on the *full* ImageNet training set ($m = 1,281,167$) using the imputed method.

Our main observation is that coherence depends on the amount of noise, and for low noise levels can be surprisingly high. For example, for $\sigma = 0.10$, each training example in our sample of 100 helps 40 other examples whereas for $\sigma = 10$, each example helps at most 5 examples.

Suppressing Weak Directions. The techniques described in Section 7 can be useful in under-parameterized situations as well. To illustrate this, we introduce outliers to our dataset by choosing 10 points at random from the points that have $x \geq 1.0$ and setting their corresponding y coordinates to -1.0. We train two models on this new dataset: one with gradient descent (GD) and another one with winsorized gradient descent (WGD) with $c = 10$. The following plots show the dataset and the regression line at the end of GD and WGD (the outliers are shown in orange):



We see that the model trained with WGD effectively ignores the outliers, producing a solution better fitting the rest of the data.

I. ADDITIONAL DATA

This section collects additional plots for variations or multiple runs of experiments previously discussed.

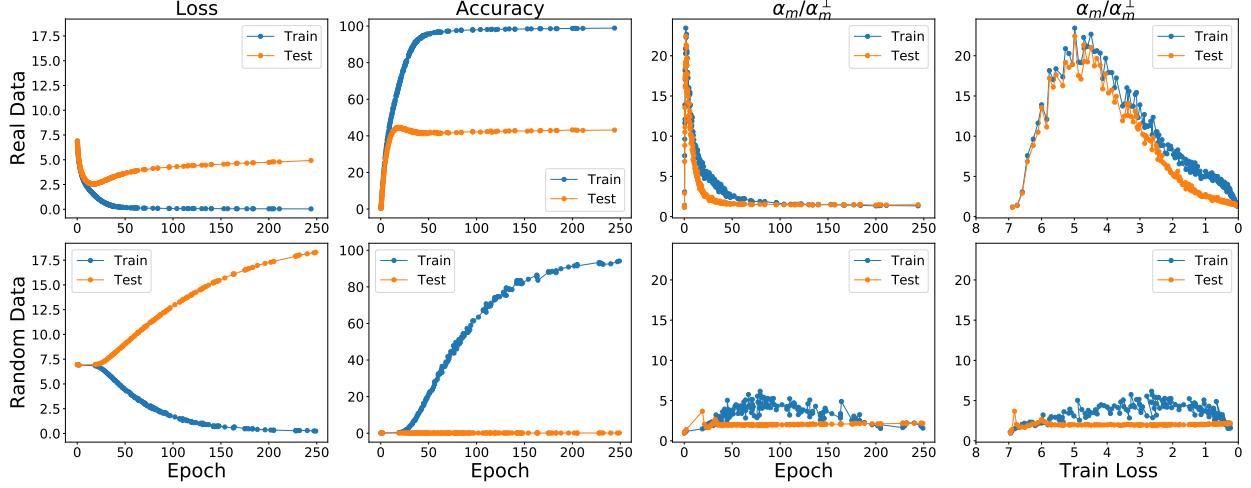


FIGURE 20. Training setup identical with Figure 14, but minibatch gradients were normalized during training and (imputed) α_m/α_m^\perp calculation to have the length of 1.0. Observe that α_m/α_m^\perp on real data is still much higher than on random.

REFERENCES

- Zeyuan Allen-Zhu, Yuanzhi Li, and Yingyu Liang. Learning and generalization in overparameterized neural networks, going beyond two layers. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019a. [12](#)
- Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over-parameterization. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 242–252. PMLR, 09–15 Jun 2019b. URL <https://proceedings.mlr.press/v97/allen-zhu19a.html>. [12](#)
- Sanjeev Arora, Rong Ge, Behnam Neyshabur, and Yi Zhang. Stronger generalization bounds for deep nets via a compression approach. In Jennifer G. Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pages 254–263. PMLR, 2018. URL <http://proceedings.mlr.press/v80/arora18b.html>. [12](#)
- Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Russ R Salakhutdinov, and Ruosong Wang. On exact computation with an infinitely wide neural net. In *Advances in Neural Information Processing Systems*, volume 32, 2019. URL <https://proceedings.neurips.cc/paper/2019/file/dbc4d84bfccfe2284ba11beffb853a8c4-Paper.pdf>. [12](#)
- Devansh Arpit, Stanislaw K. Jastrzebski, Nicolas Ballas, David Krueger, Emmanuel Bengio, Maxinder S. Kanwal, Tegan Maharaj, Asja Fischer, Aaron C. Courville, Yoshua Bengio, and Simon Lacoste-Julien. A closer look at memorization in deep networks. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, pages 233–242, 2017. URL <http://proceedings.mlr.press/v70/arpit17a.html>. [8, 8, 32](#)
- Randall Balestrieri, Jerome Pesenti, and Yann LeCun. Learning in high dimension always amounts to extrapolation, 2021. [8](#)
- Peter L. Bartlett. For valid generalization, the size of the weights is more important than the size of the network. In *Proceedings of the 9th International Conference on Neural Information Processing Systems*, NIPS’96, page 134–140, Cambridge, MA, USA, 1996. MIT Press. [1, 12, \(a\)](#)

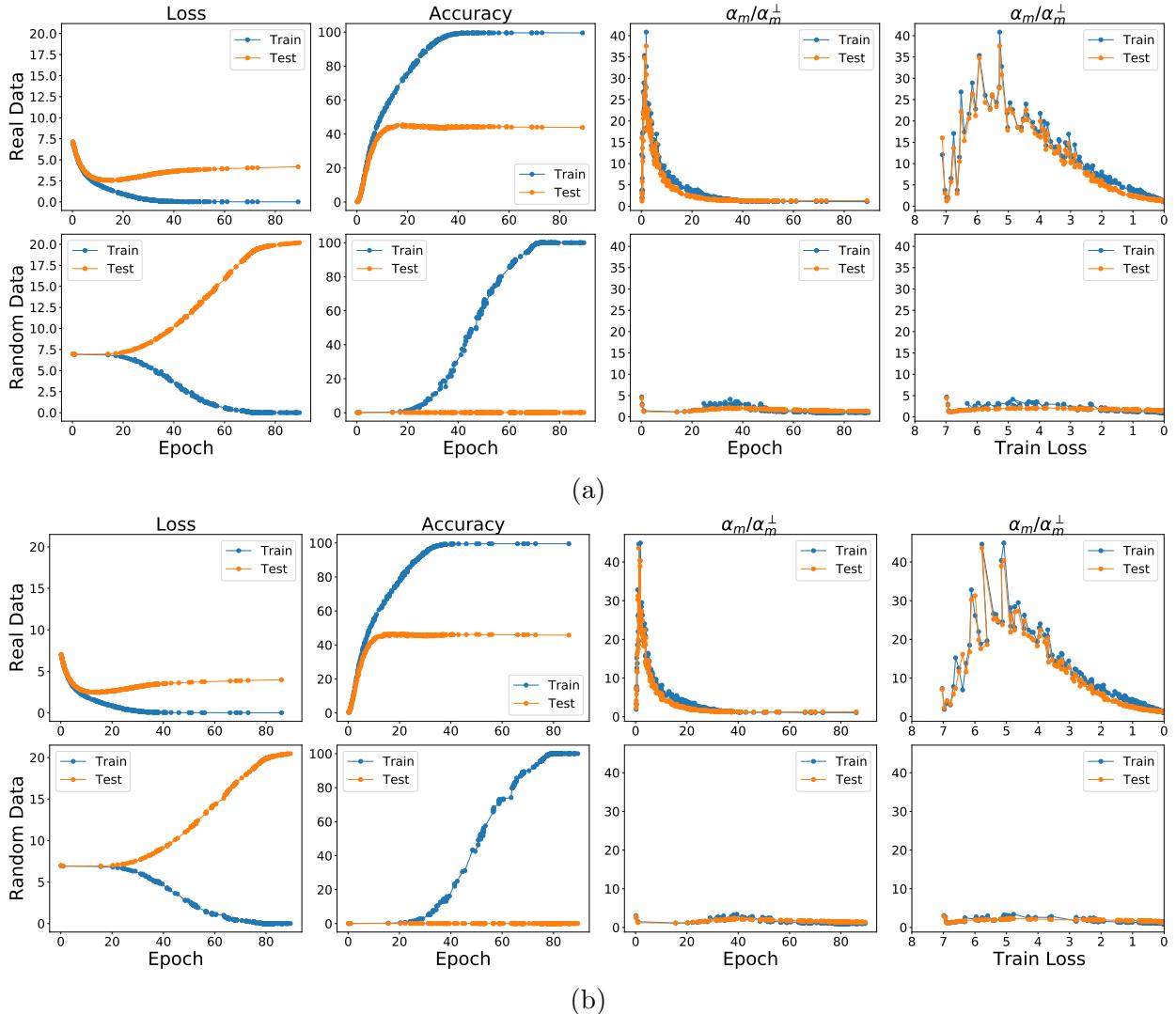


FIGURE 21. Two additional runs of the experiment in Figure 1.

Peter L Bartlett, Dylan J Foster, and Matus J Telgarsky. Spectrally-normalized margin bounds for neural networks. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 6240–6249. Curran Associates, Inc., 2017. [12](#)

Peter L. Bartlett, Philip M. Long, Gábor Lugosi, and Alexander Tsigler. Benign overfitting in linear regression. *Proceedings of the National Academy of Sciences*, 117(48):30063–30070, 2020. ISSN 0027-8424. doi: 10.1073/pnas.1907378117. URL <https://www.pnas.org/content/117/48/30063>. [12](#)

P.L. Bartlett. The sample complexity of pattern classification with neural networks: the size of the weights is more important than the size of the network. *IEEE Transactions on Information Theory*, 44(2):525–536, 1998. doi: 10.1109/18.661502. [12](#)

Eric Baum and David Haussler. What size net gives valid generalization? In D. Touretzky, editor, *Advances in Neural Information Processing Systems*, volume 1. Morgan-Kaufmann, 1989. URL <https://proceedings.neurips.cc/paper/1988/file>

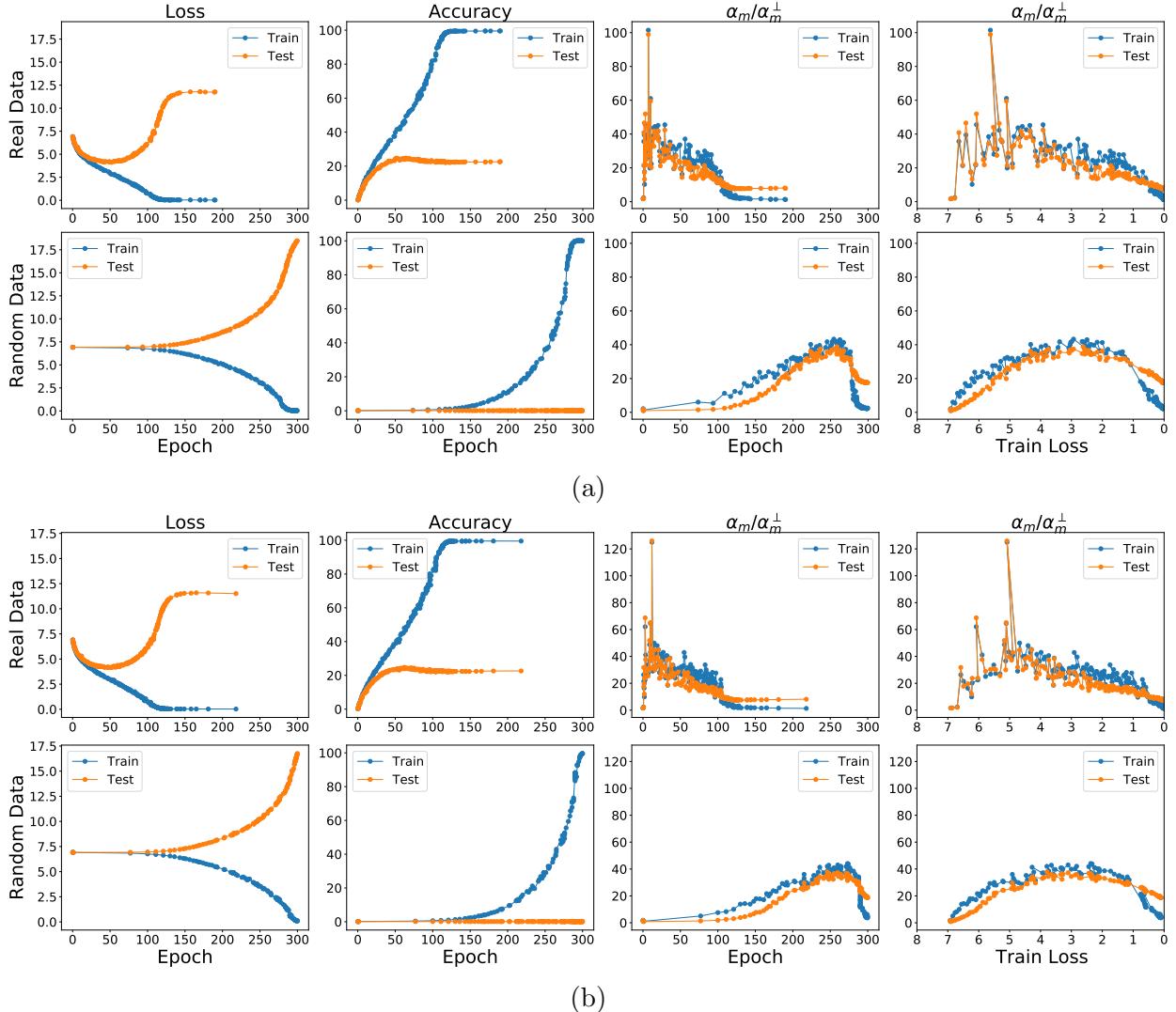


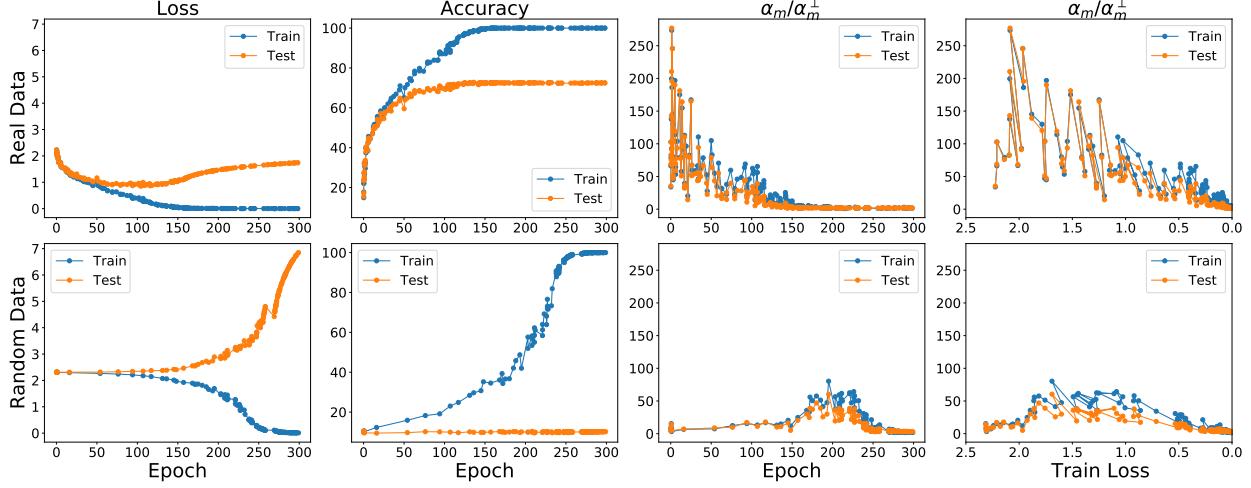
FIGURE 22. Two additional runs of the experiment in Figure 2.

[1d7f7abc18fc...-Paper.pdf](https://arxiv.org/pdf/1d7f7abc18fc...-Paper.pdf). 12

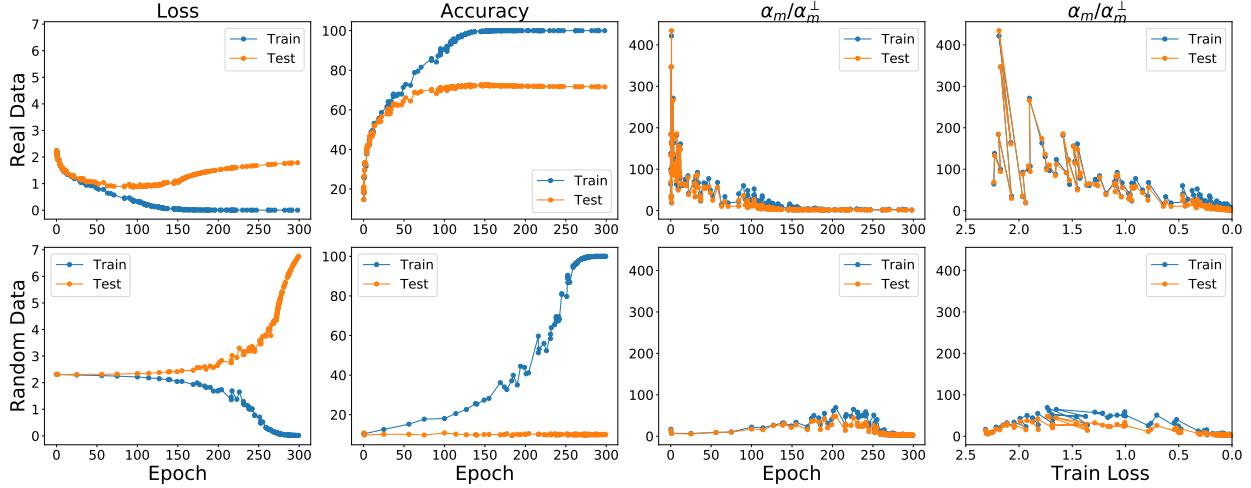
Mikhail Belkin, Siyuan Ma, and Soumik Mandal. To understand deep learning we need to understand kernel learning. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 541–549. PMLR, 10–15 Jul 2018. URL <http://proceedings.mlr.press/v80/belkin18a.html>. 12

Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. Reconciling modern machine-learning practice and the classical bias–variance trade-off. *Proceedings of the National Academy of Sciences*, 116(32):15849–15854, 2019. ISSN 0027-8424. doi: 10.1073/pnas.1903070116. URL <https://www.pnas.org/content/116/32/15849>. 11, 2

Olivier Bousquet and André Elisseeff. Stability and generalization. *J. Mach. Learn. Res.*, 2:499–526, March 2002. ISSN 1532-4435. doi: 10.1162/153244302760200704. URL <https://doi.org/10.1162/153244302760200704>. 1, 2, 12, C, C, C



(a)



(b)

FIGURE 23. Two additional runs of the experiment in Figure 13.

Satrajit Chatterjee. Coherent gradients: An approach to understanding generalization in gradient descent-based optimization. In *Proceedings of the International Conference on Learning Representations ICLR*, 2020. URL <https://openreview.net/forum?id=ryeFY0EFwS>. 3, 32, [C6], 12, 47

Satrajit Chatterjee and Alan Mishchenko. Circuit-based intrinsic methods to detect overfitting. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, Vienna, Austria, PMLR 119*, 2020. URL <https://arxiv.org/abs/1907.01991>. 2, 41

Satrajit Chatterjee and Piotr Zieliński. Making coherence out of nothing at all: Measuring the evolution of gradient alignment. *CoRR*, abs/2008.01217, 2020. URL <https://arxiv.org/abs/2008.01217>. 3, 12, B, D

Shuxiao Chen, Hangfeng He, and Weijie Su. Label-aware neural tangent kernel: Toward better generalization and local elasticity. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 15847–15858. Curran Associates, Inc., 2020. URL <https://proceedings.neurips.cc/paper/2020/>

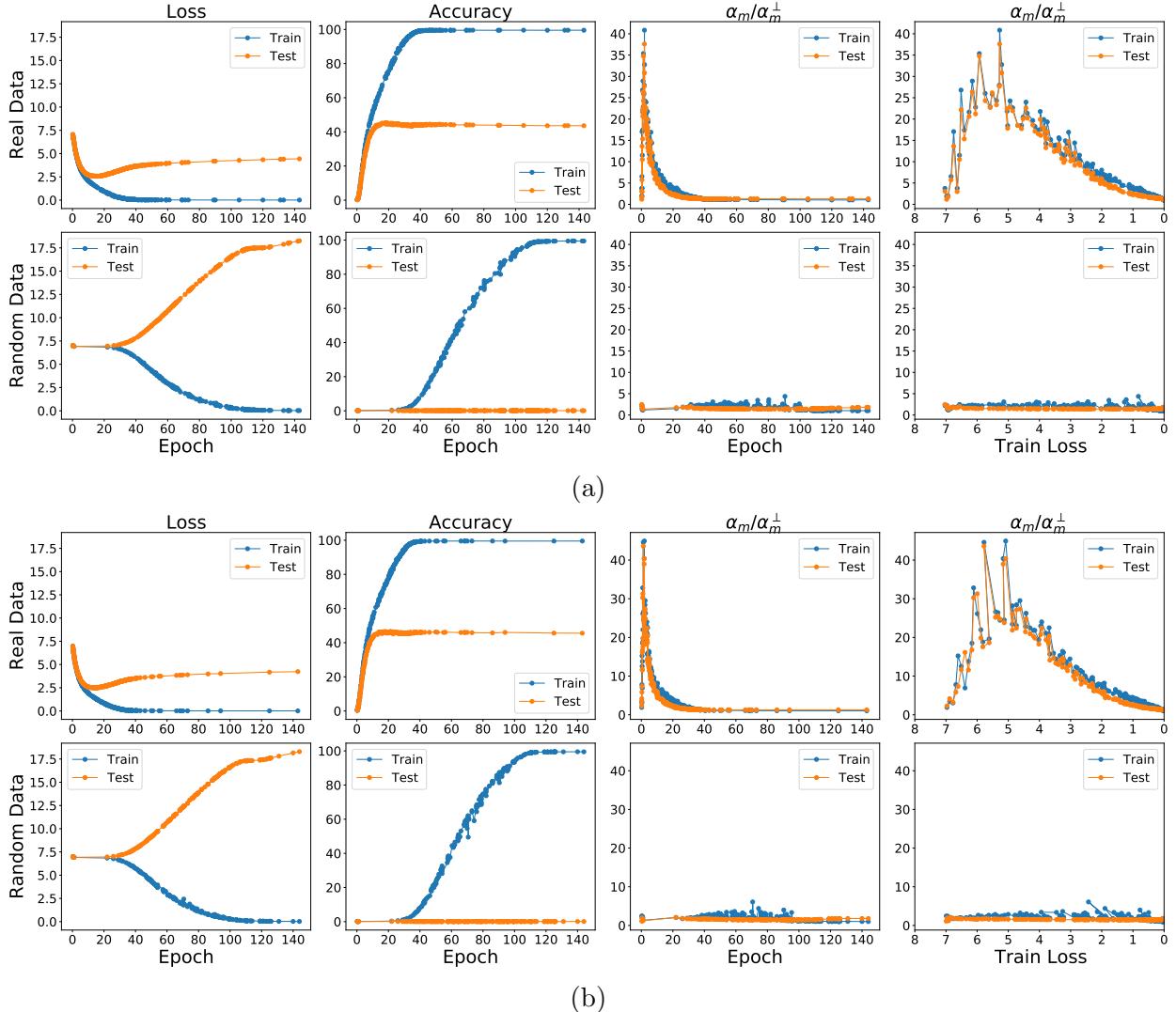


FIGURE 24. Two additional runs of the experiment in Figure 14.

[file/b6b90237b3ebd1e462a5d11dbc5c4dae-Paper.pdf](https://arxiv.org/pdf/1912.00356.pdf). 12
Lénaïc Chizat, Edouard Oyallon, and Francis Bach. On lazy training in differentiable programming. In *Advances in Neural Information Processing Systems*, volume 32, 2019. URL <https://proceedings.neurips.cc/paper/2019/file/ae614c557843b1df326cb29c57225459-Paper.pdf>. 12

Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 3

Zhun Deng, Hangfeng He, and Weijie Su. Toward better generalization bounds with locally elastic stability. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 2590–2600. PMLR, 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/deng21b.html>. 12

L. Devroye and T. Wagner. Distribution-free performance bounds for potential function rules. *IEEE Transactions on Information Theory*, 25(5):601–604, September 1979. ISSN 0018-9448.

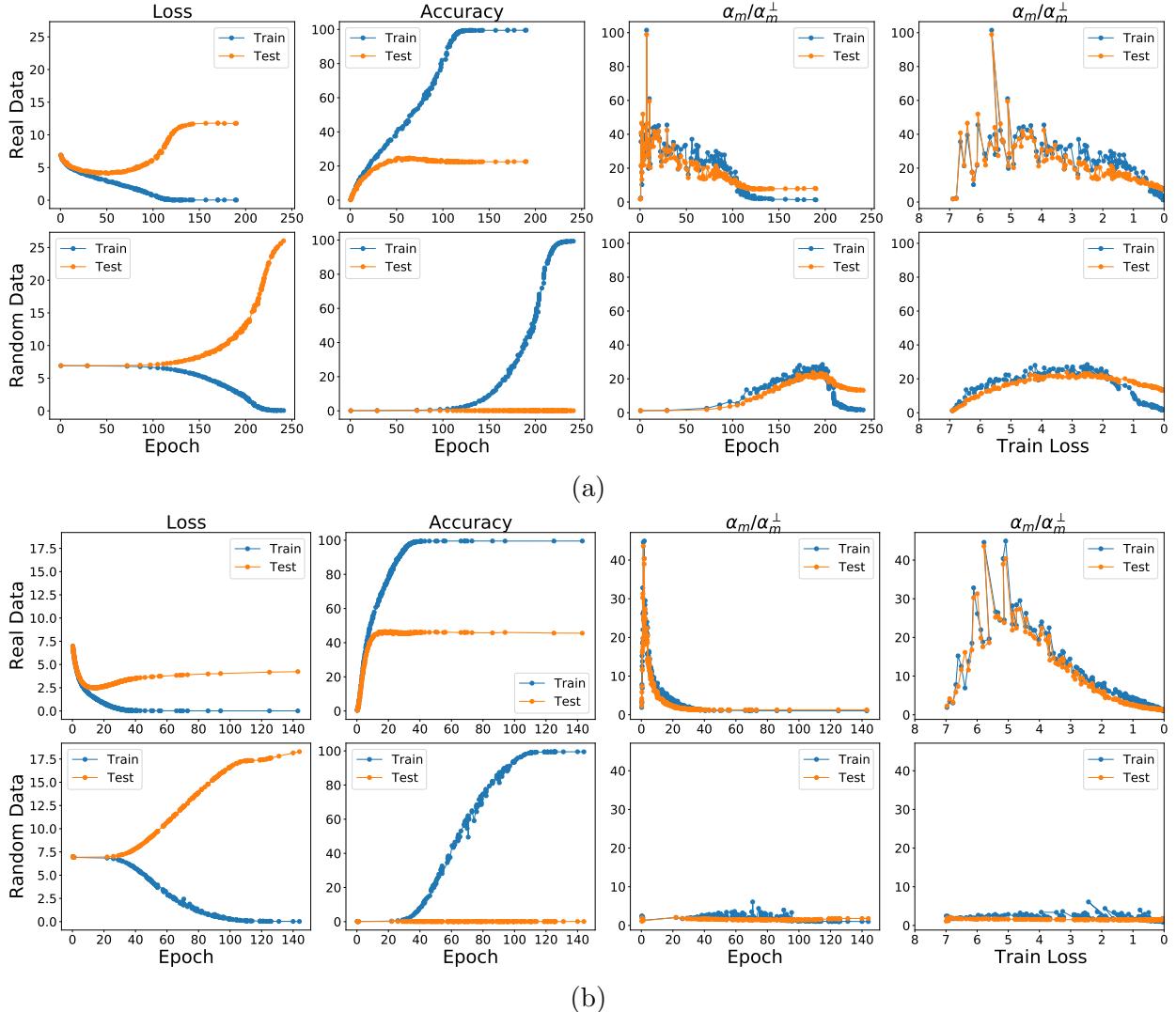


FIGURE 25. Two additional runs of the experiment in Figure 15.

doi: 10.1109/TIT.1979.1056087. 1, 2, 12, C

Laurent Dinh, Razvan Pascanu, Samy Bengio, and Yoshua Bengio. Sharp minima can generalize for deep nets. In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 1019–1028. PMLR, 06–11 Aug 2017. URL <http://proceedings.mlr.press/v70/dinh17b.html>. (b)

Pedro Domingos. Every model learned by gradient descent is approximately a kernel machine, 2020. 12

Simon Du, Jason Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds global minima of deep neural networks. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 1675–1685. PMLR, 09–15 Jun 2019a. URL <https://proceedings.mlr.press/v97/du19c.html>. 12

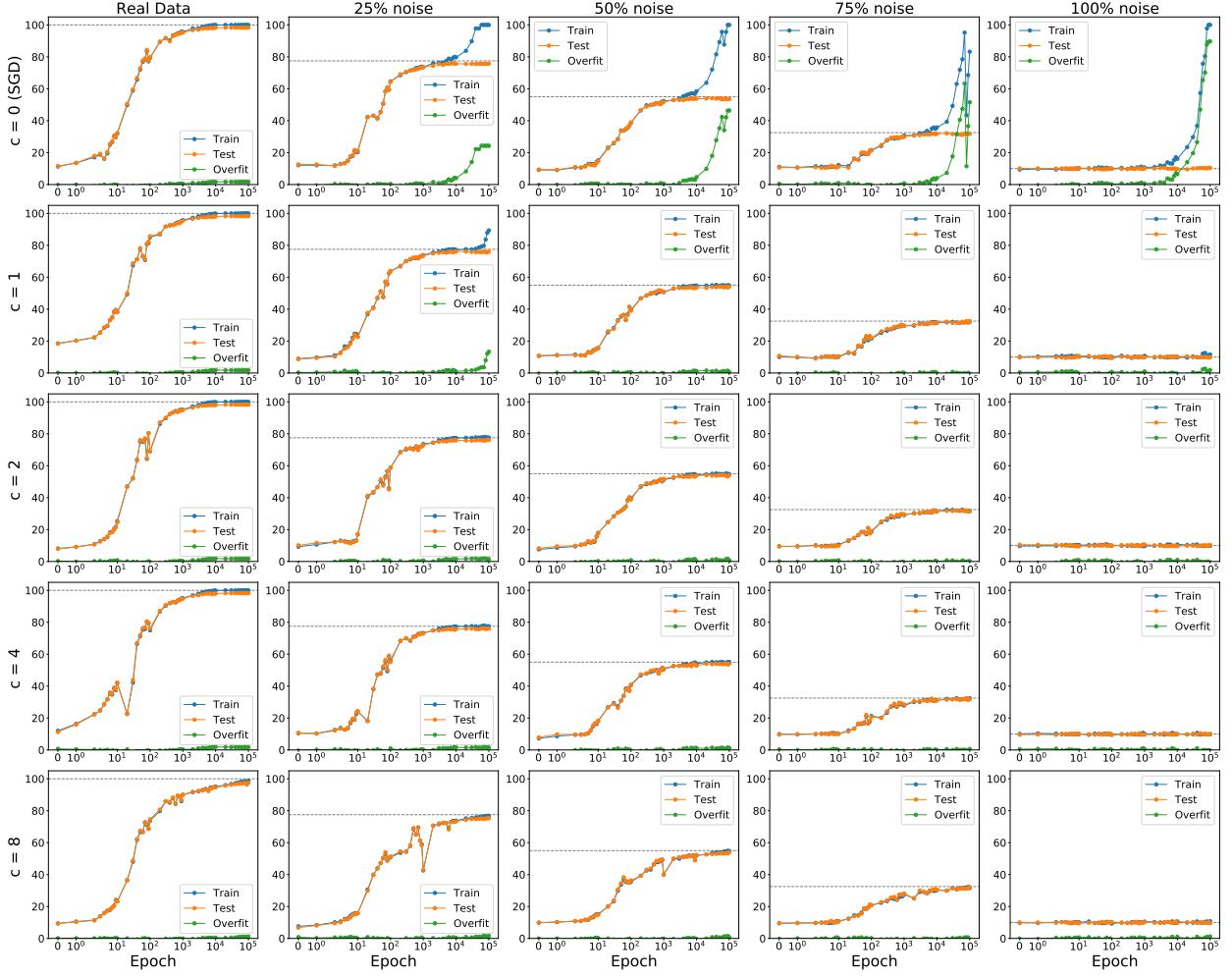


FIGURE 26. Winsorization on MNIST with random pixels. Each column represents a dataset with different noise level, e.g. the third column shows dataset with half of the examples replaced with Gaussian noise. See Figure 4 for experiments with random labels.

Simon S. Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh. Gradient descent provably optimizes over-parameterized neural networks. In *International Conference on Learning Representations*, 2019b. URL <https://openreview.net/forum?id=S1eK3i09YQ>. 12

Gintare Karolina Dziugaite and Daniel M. Roy. Computing nonvacuous generalization bounds for deep (stochastic) neural networks with many more parameters than training data. In Gal Elidan, Kristian Kersting, and Alexander T. Ihler, editors, *Proceedings of the Thirty-Third Conference on Uncertainty in Artificial Intelligence, UAI 2017, Sydney, Australia, August 11-15, 2017*. AUAI Press, 2017. URL <http://auai.org/uai2017/proceedings/papers/173.pdf>. 12, (d)

Stanislav Fort, Paweł Krzysztof Nowak, Stanisław Jastrzebski, and Srini Narayanan. Stiffness: A new perspective on generalization in neural networks. *CoRR*, abs/1901.09491v3, 2020. URL <http://arxiv.org/abs/1901.09491v3>. 1, 4, 12, B

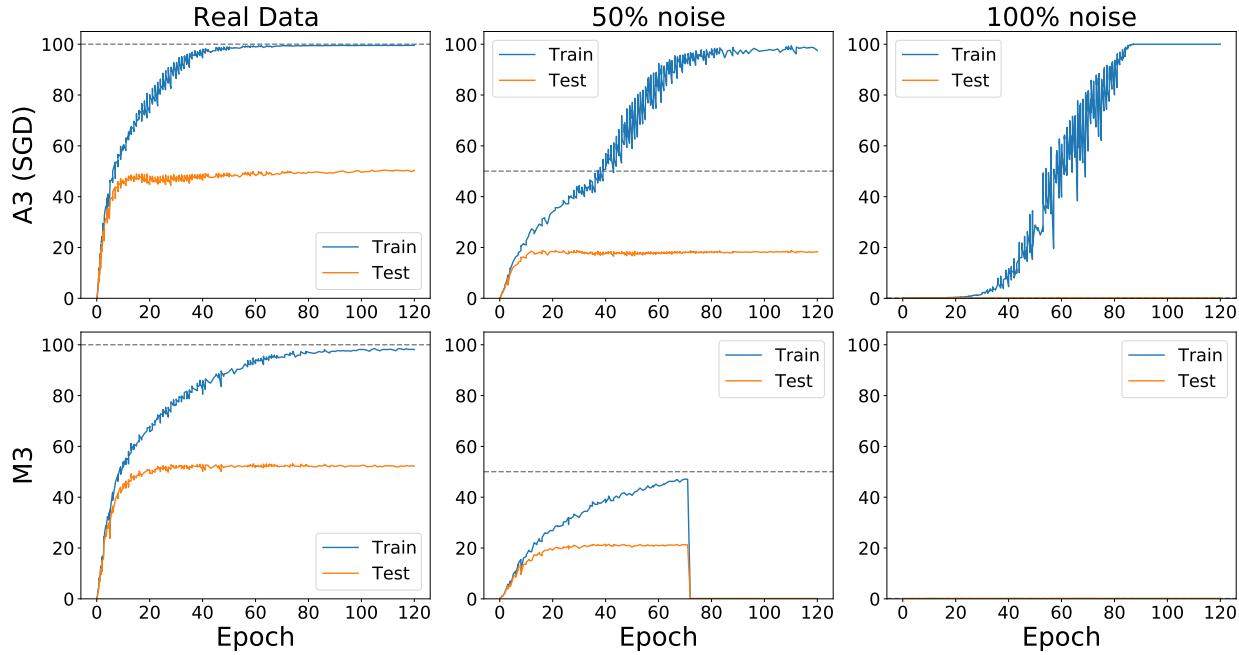


FIGURE 27. Unstable M3 in 50% noise case, learning rate of 0.1. See Figure 5 for a stable run with a different learning rate.

Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Training pruned neural networks. *CoRR*, abs/1803.03635, 2018. URL <http://arxiv.org/abs/1803.03635>. 5

Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A. Wichmann, and Wieland Brendel. Imagenet-trained CNNs are biased towards texture; increasing shape bias improves accuracy and robustness. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=Bygh9j09KX>. 6

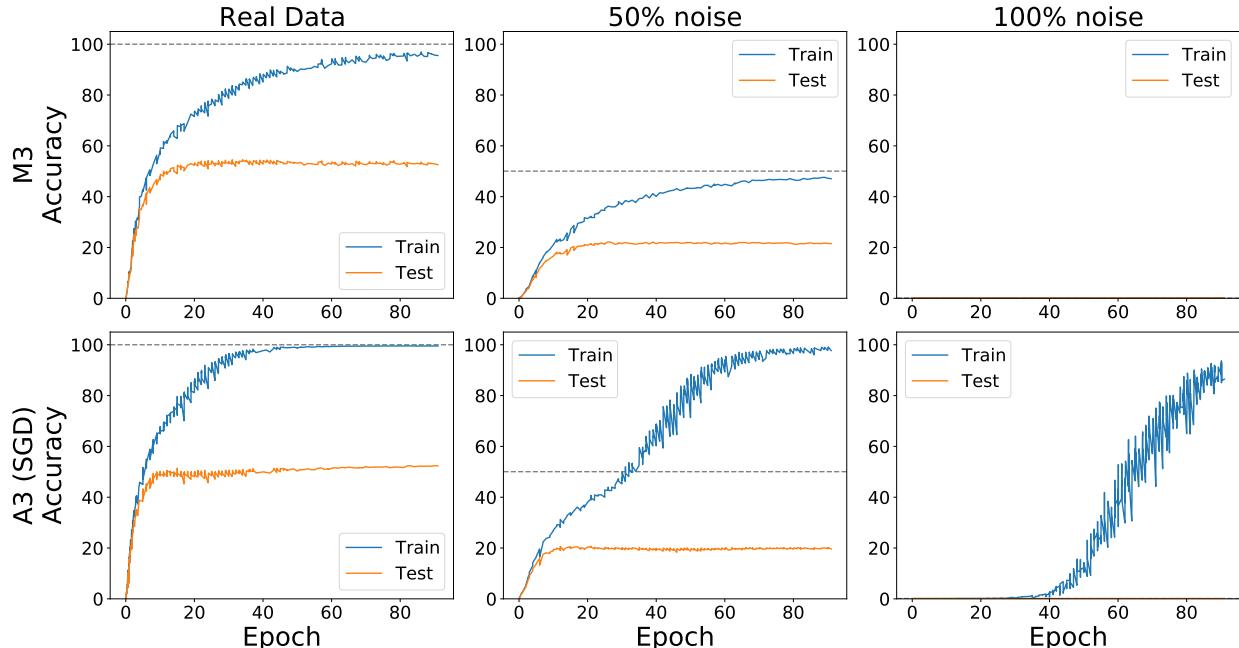
Noah Golowich, Alexander Rakhlin, and Ohad Shamir. Size-independent sample complexity of neural networks. In Sébastien Bubeck, Vianney Perchet, and Philippe Rigollet, editors, *Proceedings of the 31st Conference On Learning Theory*, volume 75 of *Proceedings of Machine Learning Research*, pages 297–299. PMLR, 06–09 Jul 2018. URL <http://proceedings.mlr.press/v75/golowich18a.html>. 12

Priya Goyal, Piotr Dollár, Ross B. Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, and Kaiming He. Accurate, large minibatch SGD: training imagenet in 1 hour. *CoRR*, abs/1706.02677, 2017. URL <http://arxiv.org/abs/1706.02677>. 19, F

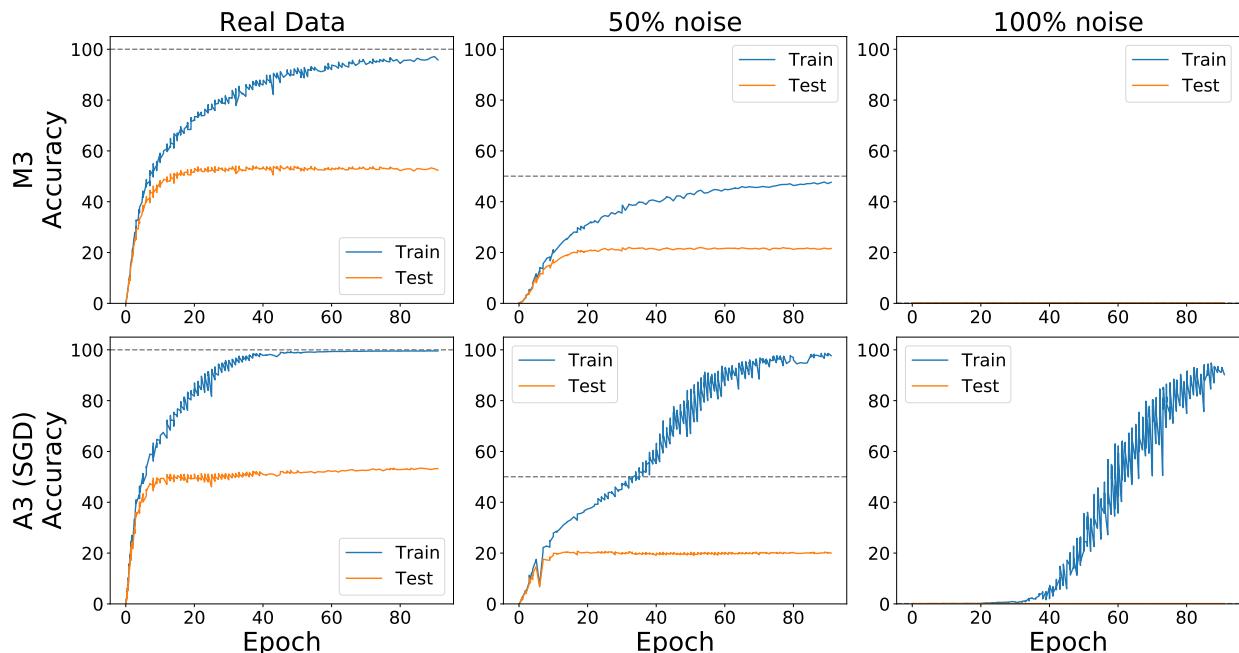
Moritz Hardt, Benjamin Recht, and Yoram Singer. Train faster, generalize better: Stability of stochastic gradient descent. In *Proceedings of the 33rd International Conference on International Conference on Machine Learning - Volume 48*, ICML’16, pages 1225–1234. JMLR.org, 2016. URL <https://arxiv.org/abs/1509.01240>. 1, 5, 5, 12, 12, B, C, C, 14

Trevor Hastie, Andrea Montanari, Saharon Rosset, and Ryan J. Tibshirani. Surprises in high-dimensional ridgeless least squares interpolation, 2020. 12

Hangfeng He and Weijie Su. The local elasticity of neural networks. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=HJxMYAntPH>. 1,



(a)



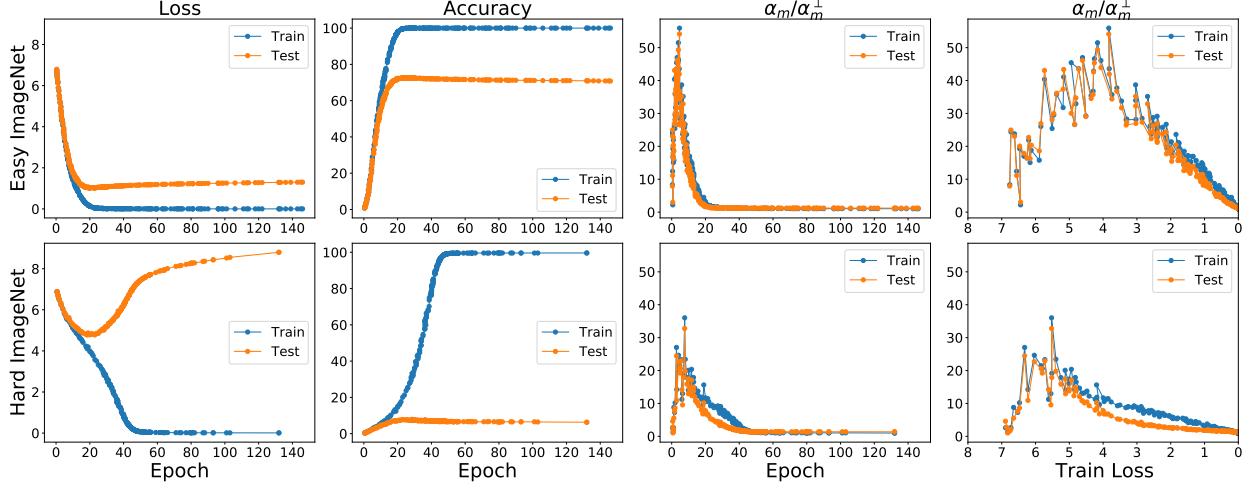
(b)

FIGURE 28. Two additional runs of the experiment in Figure 5.

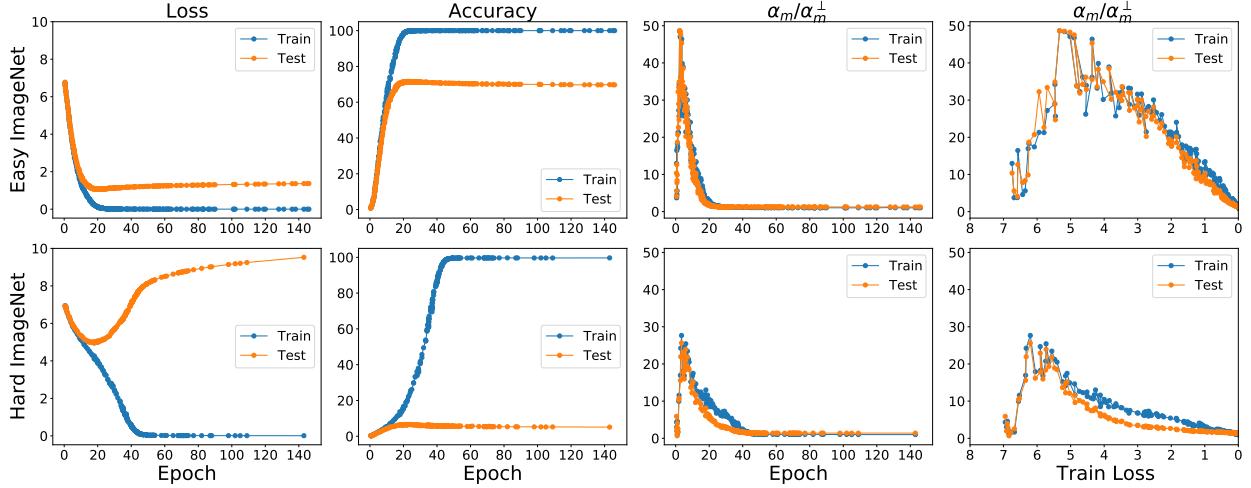
12

Haowei He, Gao Huang, and Yang Yuan. Asymmetric valleys: Beyond sharp and flat local minima. In *Advances in Neural Information Processing Systems*, volume 32, 2019. 12

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR*



(a)



(b)

FIGURE 29. Two additional runs of the experiment in Figure 7.

2016, Las Vegas, NV, USA, June 27-30, 2016, pages 770–778, 2016. doi: 10.1109/CVPR.2016.90. URL <https://doi.org/10.1109/CVPR.2016.90>. 19, 4, F

Sepp Hochreiter and Jürgen Schmidhuber. Flat minima. *Neural Comput.*, 9(1):1–42, January 1997. ISSN 0899-7667. doi: 10.1162/neco.1997.9.1.1. URL <http://dx.doi.org/10.1162/neco.1997.9.1.1>. 9.1.1. 12

Itay Hubara, Matthieu Courbariaux, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. Binarized neural networks. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016. URL <https://proceedings.neurips.cc/paper/2016/file/d8330f857a17c53d217014ee776bfd50-Paper.pdf>. [C7]

Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In Francis Bach and David Blei, editors, *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 448–456, Lille, France, 07–09 Jul 2015. PMLR. URL <http://proceedings.mlr.press/v37/>

mlr.press/v37/ioffe15.html. 6, D

Arthur Jacot, Franck Gabriel, and Clement Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018. URL <https://proceedings.neurips.cc/paper/2018/file/5a4be1fa34e62bb8a6ec6b91d2462f5a-Paper.pdf>. 12

Prateek Jain, Sham M. Kakade, Rahul Kidambi, Praneeth Netrapalli, and Aaron Sidford. Parallelizing stochastic gradient descent for least squares regression: Mini-batching, averaging, and model misspecification. *Journal of Machine Learning Research*, 18(223):1–42, 2018. URL <http://jmlr.org/papers/v18/16-595.html>. A, B

Stanislaw Jastrzebski, Maciej Szymczak, Stanislav Fort, Devansh Arpit, Jacek Tabor, Kyunghyun Cho*, and Krzysztof Geras*. The break-even point on optimization trajectories of deep neural networks. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=r1g87C4KwB>. (d)

Stanislaw Jastrzebski, Devansh Arpit, Oliver Astrand, Giancarlo B Kerg, Huan Wang, Caiming Xiong, Richard Socher, Kyunghyun Cho, and Krzysztof J Geras. Catastrophic fisher explosion: Early phase fisher matrix impacts generalization. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 4772–4784. PMLR, 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/jastrzebski21a.html>. (d)

K. Kawaguchi, L. Pack Kaelbling, and Y. Bengio. Generalization in Deep Learning. *ArXiv e-prints*, December 2017. URL <https://arxiv.org/abs/1710.05468v2>. 12

Michael Kearns and Dana Ron. Algorithmic stability and sanity-check bounds for leave-one-out cross-validation. *Neural Computation*, 11(6):1427–1453, 1999. doi: 10.1162/089976699300016304. 12, C, C

Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Peter Tang. On large-batch training for deep learning: Generalization gap and sharp minima. In *International Conference on Learning Representations*, 2017. 12

Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Cifar-10 (canadian institute for advanced research), 2009. URL <http://www.cs.toronto.edu/~kriz/cifar.html>. 2

Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 25. Curran Associates, Inc., 2012. URL <https://proceedings.neurips.cc/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf>. 2

Sameer Kumar and Norm Jouppi. Highly available data parallel ML training on mesh networks. *CoRR*, abs/2011.03605, 2020. URL <https://arxiv.org/abs/2011.03605>. 8

Ilja Kuzborskij and Christoph Lampert. Data-dependent stability of stochastic gradient descent. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 2815–2824, Stockholm Sweden, 10–15 Jul 2018. PMLR. 5, 12, C, 51, C

John Langford and Rich Caruana. (not) bounding the true error. In T. Dietterich, S. Becker, and Z. Ghahramani, editors, *Advances in Neural Information Processing Systems*, volume 14. MIT Press, 2002. URL <https://proceedings.neurips.cc/paper/2001/file/98c7242894844ecd6ec94af67ac8247d-Paper.pdf>. 12

Yann LeCun and Corinna Cortes. MNIST handwritten digit database, 2010. URL <http://yann.lecun.com/exdb/mnist/>. 1

Jian Li, Xuanyuan Luo, and Mingda Qiao. On generalization error bounds of noisy gradient methods for non-convex learning. In *8th International Conference on Learning Representations*,

- ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020.* OpenReview.net, 2020. URL <https://openreview.net/forum?id=SkxxtgHKPS>. 12
- Yuanzhi Li and Yingyu Liang. Learning overparameterized neural networks via stochastic gradient descent on structured data. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018. 12, 12
- Qianli Liao, Brando Miranda, Andrzej Banburski, Jack Hidary, and Tomaso A. Poggio. A surprising linear relationship predicts test performance in deep networks. *CoRR*, abs/1807.09659, 2018. URL <http://arxiv.org/abs/1807.09659>. 8
- Chaoyue Liu, Libin Zhu, and Mikhail Belkin. Toward a theory of optimization for over-parameterized systems of non-linear equations: the lessons of deep learning. *CoRR*, abs/2003.00307, 03 2020a. URL <https://arxiv.org/abs/2003.00307>. 12, 12
- Chaoyue Liu, Libin Zhu, and Misha Belkin. On the linearity of large non-linear models: when and why the tangent kernel is constant. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 15954–15964. Curran Associates, Inc., 12 2020b. URL <https://proceedings.neurips.cc/paper/2020/file/b7ae8fecf15b8b6c3c69eceae636d203-Paper.pdf>. 12, 46
- Jinlong Liu, Yunzhi Bai, Guoqing Jiang, Ting Chen, and Huayan Wang. Understanding why neural networks generalize well through gsnr of parameters. In *International Conference on Learning Representations*, 2020c. URL <https://openreview.net/forum?id=HyevIJStwH>. 1, 4, 12, B
- Shengchao Liu, Dimitris S. Papailiopoulos, and Dimitris Achlioptas. Bad global minima exist and SGD can reach them. *CoRR*, abs/1906.02613, 2019. URL <http://arxiv.org/abs/1906.02613>. 8, 12
- Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. In *International Conference on Learning Representations*, 2017. [C7]
- Jürgen Mayer, Khaled Khairy, and Jonathon Howard. Drawing an elephant with four complex parameters. *American Journal of Physics*, 78(6):648–649, 2010. doi: 10.1119/1.3254017. 1
- David A. McAllester. Pac-bayesian model averaging. In *Proceedings of the Twelfth Annual Conference on Computational Learning Theory*, COLT ’99, page 164–170, New York, NY, USA, 1999. Association for Computing Machinery. ISBN 1581131674. doi: 10.1145/307400.307435. URL <https://doi.org/10.1145/307400.307435>. 12
- Harsh Mehta, Ashok Cutkosky, and Behnam Neyshabur. Extreme memorization via scale of initialization. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=Z4R1vxLbRLO>. 1, 12, B
- Stanislav Minsker. Geometric median and robust estimation in banach spaces. *Bernoulli*, 21, 08 2013. doi: 10.3150/14-BEJ645. 7
- E. A. Nadaraya. On estimating regression. *Theory of Probability and its Applications*, 9:141–142, 1964. 8
- Vaishnavh Nagarajan and Zico Kolter. Deterministic PAC-bayesian generalization bounds for deep networks via generalizing noise-resilience. In *International Conference on Learning Representations*, 2019a. 12
- Vaishnavh Nagarajan and Zico Kolter. Uniform convergence may be unable to explain generalization in deep learning. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, 8-14 December 2019, Vancouver, BC, Canada*, pages 11611–11622, 2019b. 1, 15, 11, 12, (a), 12
- Preetum Nakkiran, Gal Kaplun, Dimitris Kalimeris, Tristan Yang, Benjamin L. Edelman, Fred Zhang, and Boaz Barak. SGD on neural networks learns functions of increasing complexity. *CoRR*, abs/1905.11604, 2019. URL <http://arxiv.org/abs/1905.11604>. 8

- B. Neyshabur, R. Tomioka, and N. Srebro. In Search of the Real Inductive Bias: On the Role of Implicit Regularization in Deep Learning. *arXiv e-prints*, December 2014. [1](#)
- Behnam Neyshabur, Srinadh Bhojanapalli, David McAllester, and Nathan Srebro. A PAC-Bayesian approach to spectrally-normalized margin bounds for neural networks. In *Proceedings of the International Conference on Learning Representations ICLR*, 2018a. [12](#)
- Behnam Neyshabur, Zhiyuan Li, Srinadh Bhojanapalli, Yann LeCun, and Nathan Srebro. Towards understanding the role of over-parametrization in generalization of neural networks. *CoRR*, abs/1805.12076, 2018b. URL <http://arxiv.org/abs/1805.12076>. [5](#)
- Behnam Neyshabur, Zhiyuan Li, Srinadh Bhojanapalli, Yann LeCun, and Nathan Srebro. The role of over-parametrization in generalization of neural networks. In *International Conference on Learning Representations*, 2019. [12](#)
- Nasim Rahaman, Aristide Baratin, Devansh Arpit, Felix Draxler, Min Lin, Fred Hamprecht, Yoshua Bengio, and Aaron Courville. On the spectral bias of neural networks. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 5301–5310, Long Beach, California, USA, 09–15 Jun 2019. PMLR. URL <http://proceedings.mlr.press/v97/rahaman19a.html>. [8](#)
- W. H. Rogers and T. J. Wagner. A Finite Sample Distribution-Free Performance Bound for Local Discrimination Rules. *The Annals of Statistics*, 6(3):506 – 514, 1978. doi: 10.1214/aos/1176344196. URL <https://doi.org/10.1214/aos/1176344196>. [12](#)
- David Rolnick, Andreas Veit, Serge J. Belongie, and Nir Shavit. Deep learning is robust to massive label noise. *CoRR*, abs/1705.10694, 2017. URL <http://arxiv.org/abs/1705.10694>. [9](#)
- Karthik Abinav Sankararaman, Soham De, Zheng Xu, W. Ronny Huang, and Tom Goldstein. The impact of neural network overparameterization on gradient confusion and stochastic gradient descent. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 8469–8479. PMLR, 13–18 Jul 2020. URL <http://proceedings.mlr.press/v119/sankararaman20a.html>. [1, 12, B, 1](#)
- Vatsal Shah, Soumya Basu, Anastasios Kyrillidis, and Sujay Sanghavi. On generalization of adaptive methods for over-parameterized linear regression, 2020. [12](#)
- Shai Shalev-Shwartz, Ohad Shamir, Nathan Srebro, and Karthik Sridharan. Learnability, stability and uniform convergence. *J. Mach. Learn. Res.*, 11:2635–2670, December 2010. ISSN 1532-4435. URL <http://dl.acm.org/citation.cfm?id=1756006.1953019>. [2, 12, C, 50, C](#)
- Christopher J. Shallue, Jaehoon Lee, Joseph Antognini, Jascha Sohl-Dickstein, Roy Frostig, and George E. Dahl. Measuring the effects of data parallelism on neural network training. *Journal of Machine Learning Research*, 20(112):1–49, 2019. URL <http://jmlr.org/papers/v20/18-789.html>. [\[C6\]](#)
- Hwanjun Song, Minseok Kim, Dongmin Park, and Jae-Gil Lee. Learning from noisy labels with deep neural networks: A survey. *CoRR*, abs/2007.08199, 2020. URL <https://arxiv.org/abs/2007.08199>. [9](#)
- Guillermo Valle-Perez, Chico Q. Camargo, and Ard A. Louis. Deep learning generalizes because the parameter-function map is biased towards simple functions. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=rye4g3AqFm>. [8](#)
- V. N. Vapnik and A. Ya. Chervonenkis. On the uniform convergence of relative frequencies of events to their probabilities. *Theory of Probability and its Applications*, 16(2):264–280, 1971. [1, 12](#)
- Geoffrey S. Watson. Smooth regression analysis. *Sankhyā Ser.*, 26:359–372, 1964. [8](#)

- Lei Wu, Chao Ma, and Weinan E. How sgd selects the global minima in over-parameterized learning: A dynamical stability perspective. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems 31*, pages 8279–8288. Curran Associates, Inc., 2018. **[C6]**
- Yuanzhong Xu, HyoukJoong Lee, Dehao Chen, Hongjun Choi, Blake A. Hechtman, and Shibo Wang. Automatic cross-replica sharding of weight update in data-parallel training. *CoRR*, abs/2004.13336, 2020. URL <https://arxiv.org/abs/2004.13336>. **8**
- Dong Yin, Ashwin Pananjady, Max Lam, Dimitris Papailiopoulos, Kannan Ramchandran, and Peter Bartlett. Gradient diversity: a key ingredient for scalable distributed learning. In Amos Storkey and Fernando Perez-Cruz, editors, *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*, volume 84 of *Proceedings of Machine Learning Research*, pages 1998–2007, Playa Blanca, Lanzarote, Canary Islands, 09–11 Apr 2018. PMLR. **1, 12, A, B**
- Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. In *Proceedings of the International Conference on Learning Representations ICLR*, 2017. **1, 4, 14, 1, 6, 11, 12**
- Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning (still) requires rethinking generalization. *Commun. ACM*, 64(3):107–115, February 2021. ISSN 0001-0782. doi: 10.1145/3446776. URL <https://doi.org/10.1145/3446776>. **1, (a)**
- Pan Zhou, Jiashi Feng, Chao Ma, Caiming Xiong, Steven Hoi, and Weinan E. Towards theoretically understanding why sgd generalizes better than adam in deep learning. In *Neural Information Processing Systems*, 2020. **10**
- Wenda Zhou, Victor Veitch, Morgane Austern, Ryan P. Adams, and Peter Orbanz. Non-vacuous generalization bounds at the imagenet scale: A pac-bayesian compression approach. In *Proceedings of the International Conference on Learning Representations ICLR*, 2019. **12**
- Piotr Zielinski, Shankar Krishnan, and Satrajit Chatterjee. Weak and strong gradient directions: Explaining memorization, generalization, and hardness of examples at scale. *ArXiv*, abs/2003.07422, 2020. URL <https://arxiv.org/abs/2003.07422>. **3, 32, 12**
- Difan Zou, Yuan Cao, Dongruo Zhou, and Quanquan Gu. Gradient descent optimizes over-parameterized deep relu networks. *Mach. Learn.*, 109(3):467–492, 2020. doi: 10.1007/s10994-019-05839-6. URL <https://doi.org/10.1007/s10994-019-05839-6>. **12**