

---

# PAC-Bayes Compression Bounds So Tight That They Can Explain Generalization

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 While there has been progress in developing non-vacuous generalization bounds  
2 for deep neural networks, these bounds tend to be uninformative about why deep  
3 learning works. In this paper, we develop a compression approach based on quan-  
4 tizing neural network parameters in a linear subspace, profoundly improving on  
5 previous results to provide state-of-the-art generalization bounds on a variety of  
6 tasks, including transfer learning. We use these tight bounds to better understand  
7 the role of model size, equivariance, and the implicit biases of optimization, for  
8 generalization in deep learning. Notably, we find large models can be compressed  
9 to a much greater extent than previously known, encapsulating Occam’s razor.

## 10 1 Introduction

11 There has been great effort to understand why deep learning models generalize so well on common  
12 datasets, despite the fact that these models have many more parameters than the number of training  
13 datapoints and can fit random labels [58]. These observations are not explainable through classical  
14 statistical learning theory such as VC dimension or Rademacher complexity which focus on uniform  
15 convergence over the hypothesis class. The PAC-Bayes framework, in contrast, gives a convenient  
16 way of constructing bounds where the generalization gap depends on the deep learning model (or  
17 models) that are actually found by training rather than the hypothesis set as a whole. Using this  
18 framework, many different potential explanations have been proposed drawing on properties of a  
19 deep learning model that are induced by the training dataset, such as low spectral norm [45], noise  
20 stability [2], flat minima [28], robustness and compression [2, 59].

21 In this work, we show that neural networks are substantially more compressible than previously  
22 believed when paired with structured training datasets. Constructing tighter generalization bounds  
23 than have been previously achieved, we show that this compression *alone* is sufficient to explain  
24 much of the generalization of neural networks.

## 25 Contributions of this Paper

- 26 1. We develop a new approach for training compressed neural networks that adapts the com-  
27 pressed size to the difficulty of the problem. We train in selected random subspaces of  
28 the parameters and perform learned quantization. Consequently, we achieve extremely low  
29 compressed sizes for neural networks at a given accuracy level, which is a requirement for  
30 our tight bounds.
- 31 2. Using a prior encoding Occam’s razor and our compression scheme, we construct the  
32 best generalization bounds to date on image datasets, both with data-dependent and data-  
33 independent priors. We also show how transfer learning provides improvements to com-  
34 pression and thus our generalization bounds, explaining the practical performance benefits  
35 of pre-training.

Table 1: Recent non-vacuous PAC-Bayes bounds obtained on popular image classification datasets in deep learning.  $\star$  indicates bounds obtained using data-dependent priors (Section 5.2).

Reference	Non-vacuous PAC-Bayes bounds (%)				
	Binary MNIST	MNIST	FMNIST	CIFAR-10	CIFAR-100
Dziugaite and Roy [17]	16.1	$\times$	$\times$	$\times$	$\times$
Rivasplata et al. [48]	2.2	$\times$	$\times$	$\times$	$\times$
Zhou et al. [59]	$\times$	< 46	$\times$	$\times$	$\times$
Dziugaite et al. [18]	$\times$	11 $\star$	38 $\star$	23 $\star$	$\times$
Pérez-Ortiz et al. [47]	$\times$	21.7/1.5 $\star$	$\times$	16.7 $\star$	$\times$
Our bounds	$\times$	11.6/1.4 $\star$	32.8/10.1 $\star$	58.2/16.6 $\star$	94.6/44.4 $\star$

36 3. Through the lens of compression and our bounds, we are able to explain why deep learning  
 37 models generalize on structured datasets like CIFAR-10, but not when structure is broken  
 38 such as by shuffling the pixels or shuffling the labels, the benefits of equivariant models,  
 39 and whether the implicit regularization of SGD is necessary for generalization.

## 40 2 Related Work

41 In this section, we delineate recent techniques developed for achieving non-vacuous PAC-Bayes  
 42 bounds for neural networks.

43 **Optimizing the PAC-Bayes bound** Dziugaite and Roy [17] obtain the first non-vacuous gen-  
 44 eralization bounds for deep stochastic neural networks on binary MNIST. To do so, the authors  
 45 construct a relaxation of the Langford and Seeger [36] bound and optimize it with the objective of  
 46 finding a posterior distribution that covers a large volume of low-loss solutions around a local mini-  
 47 mum obtained using SGD. Similarly, Rivasplata et al. [48] extend the optimization to other bounds  
 48 by developing novel relaxations of other PAC-Bayes bounds based on previous work from Blundell  
 49 et al. [5].

50 **Compressing Models** It is well-established that neural networks are robust to small perturba-  
 51 tions [27, 28, 35, 34, 33, 45, 7]. The work of Arora et al. [2] leverages this connection to establish  
 52 a compression-based approach that uses noise stability to study generalization of neural networks.  
 53 Following this insight, Zhou et al. [59] develops a PAC-Bayes bound that uses the representation of a  
 54 compressed model in bits, and they also add noise stability through their use of Gaussian posteriors  
 55 and Gaussian mixture priors. Furthermore, they achieve the smallest model representations through  
 56 pruning and quantization following ideas of Han et al. [26] and Cheng et al. [9].

57 **Data-Dependent Priors** Dziugaite et al. [18] demonstrated that for linear PAC-Bayes bounds  
 58 such as Thiemann et al. [51], a tighter bound can be achieved by choosing the prior distribution to be  
 59 data-dependent. That is, the prior is trained to concentrate around low loss regions on held-out data.  
 60 More precisely, the authors show that the optimal data-dependent prior is the conditional expectation  
 61 of the posterior given a subset of the training data. They approximate this data-dependent prior by  
 62 solving a variational problem over Gaussian distributions. They evaluate the bounds for SGD-trained  
 63 networks on data-dependent priors obtaining tight bounds on MNIST, Fashion MNIST, and CIFAR-  
 64 10. In a similar vein, Pérez-Ortiz et al. [47] combine the data-dependent priors approach Dziugaite  
 65 et al. [18] with the PAC-Bayes with Backprop (PBB) approach from Rivasplata et al. [48] to obtain  
 66 *state-of-the-art* PAC-Bayes non-vacuous bounds for MNIST and CIFAR-10 with data-dependent  
 67 priors. Table 1 summarizes the PAC-Bayes bounds obtained by these works alongside our bounds.

68 Our upper bound on the KL term closely follows Zhou et al. [59], but with key improvements:  
 69 (i) we train in lower dimensional subspaces using intrinsic dimensionality and FiLM subspaces  
 70 which proves to be more effective and adaptable than pruning, (ii) we consider a more effective  
 71 quantization scheme with variable length codes and quantization aware training, and (iii) we take  
 72 advantage of the increased compression provided by transfer learning and data dependent priors. As  
 73 a result, we obtain state-of-the-art non-vacuous bounds for MNIST, Fashion MNIST, CIFAR-10, and  
 74 CIFAR-100 in both the from scratch and transfer learning setting as shown in Section 4. Motivated  
 75 by the tightness of these bounds, we seek to better understand the role of equivariance, stochastic  
 76 optimization, and data dependent bounds, for explaining generalization.

77 Ding et al. [15] is another closely related work which optimizes the PAC-Bayes bounds to predict the  
 78 transferability of various upstream models to a downstream task. In contrast, we aim to explain why  
 79 neural networks generalize through compression. This is different than predicting transferability  
 80 as there are several factors other than compression (which we demonstrate is sufficient, but not  
 81 necessary for generalization) that could make a model perform well on a downstream task.

### 82 3 PAC-Bayes Bounds

83 PAC-Bayes bounds are fundamentally an expression of Occam’s razor: simpler descriptions of the  
 84 data generalize better. As an illustration, consider the finite hypothesis generalization bound. Let  
 85  $\hat{R}(h) = \frac{1}{n} \sum_{i=1}^n \ell(h(x_i), y_i)$  be the empirical risk of a hypothesis  $h$  from a finite hypothesis set  
 86  $\mathcal{H}$  for the 0-1 loss  $\ell$ , and  $R(h) = \mathbb{E}[\hat{R}(h)]$  its population risk. With probability at least  $1 - \delta$ , the  
 87 population risk for a hypothesis  $h$  using  $n$  data samples satisfies

$$R(h) \leq \hat{R}(h) + \sqrt{\frac{\log |\mathcal{H}| + \log(1/\delta)}{2n}}, \quad (1)$$

88 bounded by the empirical risk and a complexity term  $\log |\mathcal{H}|$  counting the number of bits needed to  
 89 specify any  $h \in \mathcal{H}$ .

90 But what if we don’t believe that each hypothesis is equally likely? If we consider that a given prior  
 91 distribution  $P$  over  $\mathcal{H}$  concentrates around the optimal hypothesis, then we can construct a variable  
 92 length code that uses fewer bits for the more likely hypothesis. With the best possible code for  $P$ ,  
 93 a given hypothesis  $h$  will take  $\log_2 \frac{1}{P(h)}$  bits to represent. Using this prior may allow us to pay  
 94 for a smaller complexity term even if the hypothesis set is large, so long as the hypotheses that are  
 95 consistent with the data are also likely under the prior.

96 The number of bits paid can further be reduced from  $\log_2 \frac{1}{P(h)}$  to  $\mathbb{KL}(Q \parallel P)$  by considering a  
 97 distribution of “good” solutions  $Q$ . If we don’t care which element of  $Q$  we arrive at, we can gain  
 98 some bits *back* from this insensitivity (which could be used to code a separate message). The average  
 99 number of bits to code a sample from  $Q$  using the prior  $P$  is  $\mathbb{H}(Q, P)$ , and one gets  $H(Q)$  bits back  
 100 from being agnostic about which sample  $h \sim Q$  to use, yielding the KL-divergence between  $Q$  and  
 101  $P$ :  $\mathbb{H}(Q, P) - \mathbb{H}(Q) = \mathbb{KL}(Q \parallel P)$ .

102 With these improvements of the finite hypothesis bound: replacing  $\log |\mathcal{H}|$  with  $\mathbb{KL}(Q \parallel P)$ , and  
 103 sampling a hypothesis  $h \in \mathcal{H}$ , one arrives (with minor bookkeeping) at the PAC-Bayes bound  
 104 introduced in McAllester [43]. With probability at least  $1 - \delta$ ,

$$\mathbb{E}_{h \sim Q} [R(h)] \leq \mathbb{E}_{h \sim Q} [\hat{R}(h)] + \sqrt{\frac{\mathbb{KL}(Q \parallel P) + \log(n/\delta) + 2}{2n - 1}}. \quad (2)$$

105 Many refinements of Eq. (2) have been developed [36, 42, 6, 51] but retain the same character. The  
 106 lower the ratio of the  $\mathbb{KL}$  to the number of data points  $n$ , the tighter the gap between empirical  
 107 and expected risk. In this work, we use the tighter Catoni [6] variant of the PAC-Bayes bound (see  
 108 Appendix I for details).

109 **Universal Prior.** Connecting back to Occam’s razor, one can define a prior that explicitly penalizes  
 110 the minimum compressed length of the hypothesis, also known as the universal prior [50]:  $P(h) =$   
 111  $2^{-K(h)}/Z$ , where  $K$  is the *prefix* Kolmogorov complexity [30] of  $h$  (the length of the shortest  
 112 program that produces  $h$  and also delimits itself), and  $Z \leq 1$ .<sup>1</sup> Using a point mass posterior on a  
 113 single hypothesis  $h^*$ ,  $\mathbb{KL}(Q \parallel P)$  from the bound reduces to

$$\mathbb{KL}(\mathbf{1}_{[h=h^*]} \parallel P(h)) = \log \frac{1}{P(h^*)} \leq K(h) \log 2 \leq l(h) \log 2 + 2 \log l(h),$$

114 where  $l(h)$  is the length of a given program that reproduces  $h$  not including the delimiter, which we  
 115 can use when upper bounding the KL. For convenience, we can condition on using the same method  
 116 for compression and decompression for all elements of the prior. If we can reduce the compressed

<sup>1</sup>The universal prior similar to the discrete hypothesis prior from Zhou et al. [59] but setting  $m(h) = 2^{-2 \log_2 l(h)}$  rather than the flat  $m(h) = 2^{-72}$ .

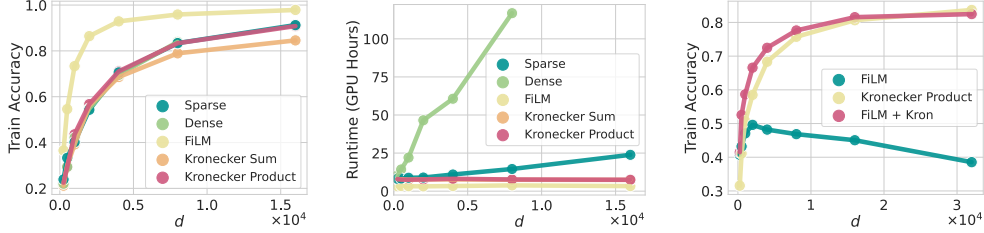


Figure 1: **(Left)** Different projection operators  $P$  (Section 4.1) used for transfer learning from Imagenet to CIFAR-10 on a ResNet-34 across different subspace dimensions  $d$ . Kronecker product, Sparse, and Dense perform almost identically **(Center)** Kronecker product runs with substantially reduced the runtime cost compared to the Sparse or Fastfood matrices used by Li et al. [38]. **(Right)** Training from scratch on CIFAR-10. The FiLM projector alone is unable to fit the data when training from scratch, and instead a sum of FiLM and Kronecker Product projectors perform the best.

117 length  $l(h)$  of the hypotheses we find that fits the training data, then we can improve the tightness of  
 118 the PAC-Bayes bound.

119 **Model Compression.** *Model compression* aims to find a nearly equivalent model that can be ex-  
 120 pressed fewer bits and is useful for deployments on mobile devices or for improving inference time  
 121 on specialized hardware [8, 9]. For computing PAC-Bayes generalization bounds, however, only  
 122 the model size matters. Therefore, we can employ compression methods which may otherwise be  
 123 unfavorable in practice due to worse computational requirements. *Pruning* and *quantization* are  
 124 among the most widely used methods for model compression. In this work, we rely on quantization  
 125 (Section 4.2) to achieving tighter generalization bounds.

## 126 4 Tighter Generalization Bounds via Adaptive Subspace Compression

127 Training a neural network involves taking many gradient steps in a high-dimensional space  $\mathbb{R}^D$ .  
 128 Although  $D$  may be large, the loss landscape has been, in some sense, found to be simpler than  
 129 typically believed [14, 25, 22]. Analogous to the notion of intrinsic dimensionality more generally,  
 130 Li et al. [38] attempt to find the lowest dimensional subspace in which the network can be trained  
 131 and still fit the training data. The weights of a neural network  $\theta \in \mathbb{R}^D$  are parametrized in terms of  
 132 an initialization  $\theta_0$  and a projection  $w \in \mathbb{R}^d$  to a lower dimensional subspace through a fixed matrix  
 133  $P \in \mathbb{R}^{D \times d}$ ,

$$\theta = \theta_0 + Pw. \quad (3)$$

134 To facilitate favorable conditioning during optimization,  $P$  is chosen to be approximately orthonormal  
 135  $P^\top P \approx I_{d \times d}$  and for scalability, Li et al. [38] utilize random normal matrices of the form  
 136  $P \sim \mathcal{N}(0, 1)^{D \times d} / \sqrt{D}$ , and their sparse approximations [39, 37].

137 In its original form, intrinsic dimensionality (ID) is used as a scientific tool to measure the complex-  
 138 ity of the learning task. Unlike other methods like pruning, ID has the benefit of scaling dramatically  
 139 up or down in size with difficulty of the problem. While ID was not used with quantization as a full  
 140 fledged model compression method, we find that it can be. For our work, the ability to find the  
 141 intrinsic dimension  $d \ll D$  has profound implications for the compressibility of the models, and  
 142 therefore our ability to construct generalization bounds. As demonstrated by Zhou et al. [59], the  
 143 compressibility of neural network models has a direct connection to generalization and immediately  
 144 allows us to compute non-vacuous PAC-Bayes bounds for transfer learning. Building upon ID, we  
 145 describe our approach next.

### 146 4.1 Finding Better Random Subspaces to Train In

147 To further improve upon the scalability and effectiveness of the projections  $P$  used by Li et al. [38],  
 148 we introduce three new methods.

149 **Kronecker Sum Projector:** Using the Kronecker product  $\otimes$ , we construct the matrix  
 150  $P_{\oplus} = (\mathbf{1} \otimes R_1 + R_2 \otimes \mathbf{1}) / \sqrt{2D}$  where  $R_1, R_2 \sim \mathcal{N}(0, 1)^{\sqrt{D} \times d}$  and  $\mathbf{1}$  is the vector of all ones

151 in  $\mathbb{R}^{\sqrt{D}}$ . From the fact that  $R_1 \perp R_2$  and the entries are standard normal, one can verify that  
 152  $P_{\oplus}^{\top} P_{\oplus} = I_{d \times d} + \mathcal{O}(1/\sqrt{D})$ .

153 **Kronecker Product Projector:** Alternatively, we can form the matrix  $P_{\otimes} = Q_1 \otimes Q_2 / \sqrt{D}$  with  
 154 the smaller  $Q_1, Q_2 \sim \mathcal{N}(0, 1)^{\sqrt{D} \times \sqrt{d}}$ , and again this matrix is approximately orthogonal:  $P_{\otimes}^{\top} P_{\otimes} =$   
 155  $(Q_1^{\top} Q_1 / \sqrt{D}) \otimes (Q_2^{\top} Q_2 / \sqrt{D}) = I \otimes I + \mathcal{O}(D^{-1/4}) = I_{d \times d} + \mathcal{O}(D^{-1/4})$ .<sup>2</sup>

156 The matrix vector multiply  $w \mapsto Pw$  for both the above projectors can be performed in time  
 157  $O(d\sqrt{D})$  and  $O(\sqrt{dD})$  respectively, rather than the  $O(dD)$  that is required by the dense random  
 158 matrix. Figure 1 demonstrates the runtime speedup and training performance improvement in com-  
 159 parison to the methods using by Li et al. [38]. Notably, the Kronecker-structured projections retain  
 160 the fidelity of the dense random matrix while being orders of magnitude faster than the alternative  
 161 operators when scaling to larger values of  $d$ .

162 **FiLM projector:** It is well known that the BatchNorm parameters have an outsized effect on the  
 163 downstream task performance relative to their size, and this fact has been used in Featurewise in-  
 164 dependent Linear Modulation (FiLM) [46, 16] for efficient control of neural networks in many dif-  
 165 ferent settings. Several authors have explored performing finetuning for transfer learning solely on  
 166 these parameters and the final linear layer [32]. Drawing on these observations, we construct a  
 167 projection matrix  $P_{\text{FiLM}}$  where only columns corresponding with BatchNorm or head parameters  
 168 are non-zero and sampled from  $\mathcal{N}(0, 1)^d / \sqrt{D}$ , which we also show in Figure 1. While the FiLM  
 169 projector is highly effective for transfer learning (shown in Figure 1 left), the performance saturates  
 170 quickly when training from scratch. For this reason, when training from scratch we employ the sum  
 171  $P_{\text{FiLM}+\otimes} = (P_{\text{FiLM}} + P_{\otimes}) / \sqrt{2}$ , which outperforms the two projectors individually as shown in  
 172 Figure 1 right.

## 173 4.2 Quantization Scheme and Training

174 As a consequence of finite-precision computations (e.g. 32-bit floating points), the hypothesis space  
 175  $\mathcal{H}$  of any neural network remains finite in practice. Moreover, the ability of neural networks to be  
 176 compressed [26, 11] to far fewer number of bits per parameter while retaining similar performance  
 177 allows us to compress the learned subspace vector  $w \in \mathbb{R}^d$  from Eq. (3) as well.

178 Given the full precision weights  $w = [w_1, \dots, w_d] \in \mathbb{R}^d$  and a vector  $c = [c_1, \dots, c_L] \in \mathbb{R}^L$  of  $L$   
 179 quantization levels, we construct the quantized vector  $\hat{w} = [\hat{w}_1, \dots, \hat{w}_d]$  such that  $\hat{w}_i = c_{q(i)}$  where  
 180  $q(i) = \arg \min_k |w_i - c_k|$ . The quantization levels  $c$  are learned alongside  $w$ , where the gradients  
 181 are defined using the straight through estimator [3, 56]:

$$\frac{\partial \hat{w}_i}{\partial w_j} = \delta_{ij} \quad \text{and} \quad \frac{\partial \hat{w}_i}{\partial c_k} = \mathbf{1}_{[q(i)=k]} \quad (4)$$

182 We initialize  $c$  with uniform spacing between the minimum and maximum values in parameter  
 183 vector  $w$ , which performs better than k-means initialization, as also observed in Choi et al. [10].  
 184 To further compress the network, we use a variable length code in the form of arithmetic coding  
 185 [41], which takes advantage of the fact that certain quantization levels are more likely than others.  
 186 Given probabilities  $p_k$  (empirical fractions) for cluster  $c_k$ , arithmetic coding of  $w$  takes at most  
 187  $\lceil d \times \mathbb{H}(p) \rceil + 2$  bits, where  $\mathbb{H}(p)$  is the entropy  $\mathbb{H}(p) = -\sum_k p_k \log_2 p_k$ . For a small number of  
 188 quantization levels, we find it makes a small improvement in the compression over more common  
 189 Huffman coding.

190 In summary, we use  $\lceil d \times \mathbb{H}(p) \rceil + 2$  bits for coding the quantized weights  $\hat{w}$ ,  $16L$  bits for the  
 191 codebook  $c$  (represented in half precision), and additional  $L \times \lceil \log_2 d \rceil$  bits for representing the  
 192 probabilities  $p_k$ , arriving at  $l(w) \leq \lceil d \times \mathbb{H}(p) \rceil + L \times (16 + \lceil \log_2 d \rceil) + 2$ .

<sup>2</sup>As neither  $D$  nor  $d$  is typically a perfect square, we concatenate a dense random matrix to pad out the  
 difference between  $D, d$ , and a perfect square. As  $(\sqrt{D} + 1)^2 = D + 2\sqrt{D} + 1$ , we have that the size of  
 this padding is at most  $\sqrt{D} \times \sqrt{d}$ , so it does not increase the asymptotic cost of performing the matrix vector  
 multiplies.

---

**Algorithm 1** Compute PAC-Bayes Bound.

---

```
1: Inputs: Neural network  $f_\theta$ , Training dataset  $\{x_i, y_i\}_{i=1}^n$ , Clusters  $L$ , Intrinsic dimension  $d$ , Confidence  $1 - \delta$ , and Prior distribution  $P$ .
2: function COMPUTE_BOUND( $f_\theta, L, d, (x_i, y_i)_{i=1}^n, \delta, P$ )
3:    $w \leftarrow \text{TRAIN\_ID}(f_\theta, d, (x_i, y_i)_{i=1}^n)$  ▷ (Section 4.1)
4:    $\hat{w} \leftarrow \text{TRAIN\_QUANTIZE}(w, L, (x_i, y_i)_{i=1}^n)$ 
5:   Compute quantized train error  $\hat{R}(\hat{w})$ .
6:    $\mathbb{KL}(Q, P) \leftarrow \text{GET\_KL}(\hat{w}, P)$  ▷ (Section 3)
7:   return GET_CATONI_BOUND( $\hat{R}(\hat{w}), \mathbb{KL}(Q, P), \delta, n$ ) ▷ (Section 3)
8: end function
9: function TRAIN_QUANTIZE( $w, L, (x_i, y_i)_{i=1}^n$ ) ▷ (Section 4.2)
10:   Initialize  $c \leftarrow \text{GET\_CLUSTERS}(w, L)$ 
11:   for  $i = 1$  to quant_epochs do
12:      $c \leftarrow c - \rho \nabla_c \mathcal{L}(w, c)$  and  $w \leftarrow w - \rho \nabla_w \mathcal{L}(w, c)$ 
13:   end for
14:   return  $\hat{w}$ 
15: end function
16: function GET_KL( $(\hat{w}, P)$ )
17:    $c, \text{count} \leftarrow \text{GET\_UNIQUE\_VALS\_COUNTS}(\hat{w})$ 
18:    $\text{message\_size} \leftarrow \text{DO\_ARITHMETIC\_ENCODING}(\hat{w}, c, \text{count})$ 
19:    $\text{message\_size} \leftarrow \text{message\_size} + \text{hyperparam\_search}$  ▷ (Section 4.3)
20:   return  $\text{message\_size} + 2 \times \log(\text{message\_size})$ 
21: end function =0
```

---

### 193 4.3 Subspace Dimension Optimization and Hyperparameters in the Universal Prior

194 The smaller the chosen intrinsic dimension  $d$ , the more similar  $\theta$  is to the initialization  $\theta_0$  in Eq. (3).  
195 Consequently, that value of  $\theta$  is more likely under the universal prior given the shorter description  
196 length. Note that in this prior, we condition on the random seed used to generate  $\theta_0$  and  $P$ . As  
197 we optimize over different parameters such as the subspace dimension  $d = 1, \dots, D$ , and possibly  
198 other hyperparameters such as the learning rate, or number of quantization levels  $L$ , we must encode  
199 these into our prior and thus pay a penalty for optimizing over them. We can accommodate this very  
200 simply by considering the hypothesis  $h$  as not just specifying the weights, but also specifying these  
201 hyperparameters:  $h = (\theta, d, L, \text{lr})$ , and therefore using the universal prior  $P(h) = 2^{-K(h)}/Z$  we  
202 pay additional bits for each of these quantities:  $K(h) \leq K(\theta|d, L) + K(d) + K(L) + K(\text{lr})$ . If  
203 we optimize over a fixed number  $H$  of distinct values known in advance for a given hyperparameter  
204 such as  $L$ , then we can code  $L$  using this information in only  $\log_2(H)$  bits. In general, we can also  
205 bound the dimensionalities searched over by the maximum  $D$  so that  $K(d) \leq \lceil \log_2 D \rceil$  in any case.

206 For transfer learning, we replace  $\theta_0$  with a learned initialization  $\theta_{\mathcal{D}}$  that is found using the pretraining  
207 task and data  $\mathcal{D}$ . With the ID compression, the universal prior  $P(h|\mathcal{D})2^{-K(h|\mathcal{D})}/Z$  will place higher  
208 likelihood on solutions  $\theta$  that are close to the pretraining solution  $\theta_{\mathcal{D}}$ .

## 209 5 Empirical Non-Vacuous Bounds

### 210 5.1 Non-Vacuous PAC-Bayes Bounds

211 To compute our bounds on benchmark datasets, we use the joint approach of projected subspace  
212 training, parameter quantization, and optimizing the subspace dimension with respect to the bounds  
213 as summarized in Algorithm 1. We also optimize other bound hyperparameters such as the learning  
214 rate of the quantization-aware training. We offer more details about bound optimization in Ap-  
215 pendix D. Our results for *data-independent* priors, i.e., priors that do are not pre-trained on a subset  
216 of the training data, are shown in Table 2. We derive the first non-vacuous bounds on FashionM-  
217 NIST, CIFAR-10, and CIFAR-100 without data-dependent priors. These results are of significance  
218 as we argue in Section 5.2 that using data-dependent priors are not explanatory about the learn-  
219 ing process. In particular, we improve over the compression bound results obtained by Zhou et al.  
220 [59] on MNIST from  $< 46\%$  to 11.55%. We also dramatically improve their compression results  
221 as we reduce the compressed size for the best bound by 94% bringing it down from 6.23 KB to  
222 0.38KB for MNIST with LeNet5. The tightness of our SOTA subspace compression bounds allows

Table 2: Our PAC-Bayesian Subspace Compression Bounds with data-independent priors compared to *state-of-the-art* (SoTA) bounds. All results are obtained with 95% confidence, i.e.  $\delta = .05$ .

Dataset	Data-independent priors		Data-dependent priors	
	Err. Bound (%)	SoTA (%)	Err. Bound (%)	SoTA (%)
MNIST	<b>11.6</b>	21.7 [47]	<b>1.4</b>	1.5 [47]
+ SVHN Transfer	<b>9.0</b>	–	–	–
FashionMNIST	<b>32.8</b>	–	<b>10.1</b>	38 [18]
+ CIFAR-10 Transfer	<b>28.2</b>	–	–	–
CIFAR-10	<b>58.2</b>	–	<b>16.6</b>	16.7 [47]
+ ImageNet Transfer	<b>35.1</b>	–	–	–
CIFAR-100	<b>94.6</b>	–	<b>44.4</b>	–
+ ImageNet Transfer	<b>81.3</b>	–	–	–

223 us to improve the understanding of several deep learning phenomena as discussed in Section 6. See  
 224 Appendix A for additional results.

## 225 5.2 Data-Dependent Priors

226 So far, we demonstrated the strength of our bounds on *data-independent* priors, where we con-  
 227 siderably improve on the state-of-the-art. However, a number of recent papers have considered  
 228 data-dependent priors as a way of achieving tighter bounds [47, 18]. In this setup, the training data  
 229  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$  is partitioned into two parts,  $\mathcal{D}_a$  and  $\mathcal{D}_b$ , with length  $n - m$  and  $m$ . The first  
 230 dataset is used to construct a data dependent prior  $P(h|\mathcal{D}_a)$ , and then the bound is formed over the  
 231 remaining part of the process: the adaptation of the prior  $P(h|\mathcal{D}_a)$  to the posterior  $Q(h)$  using the  
 232 data  $\mathcal{D}_b$ . The empirical risk is computed over  $\mathcal{D}_b$  only. Intuitively, using dataset  $\mathcal{D}_a$  it is possible to  
 233 construct a much tighter prior over the possible neural network solutions.

234 In our setting, similar to transfer learning, we use the prior  $P_{\mathcal{D}_a}(\theta) = 2^{-K(\theta|\theta_{\mathcal{D}_a})}/Z$  where for  
 235 compression we use  $\theta = \theta_{\mathcal{D}_a} + Pw$ , and  $\theta_{\mathcal{D}_a}$  is the solution found by training the model (without  
 236 random projections) on the data  $\mathcal{D}_a$  rather than initializing randomly. With these data-dependent  
 237 priors, we achieve the best bounds in Table 2.

238 However, our adaptive approach exposes a significant downside of data-dependent priors. To the  
 239 extent that PAC-Bayes bounds can be used for explanation, in this case they only yield insights into  
 240 the procedure used to adapt  $P_{\mathcal{D}_a}(\theta)$  to  $Q$  using  $\mathcal{D}_b$ , any learning that is done in finding  $P_{\mathcal{D}_a}(\theta)$  is  
 241 not constrained or explained by the bound (we have no information on the compressibility of  $\theta_{\mathcal{D}_a}$ ).  
 242 Given the ability to adapt the size of the KL to the difficulty of the problem, it is possible to squeeze  
 243 all of the learning into  $P_{\mathcal{D}_a}(\theta)$  and none in this adaption to  $Q$ . This happens as the  $KL \rightarrow 0$ ,  
 244 which we find happens empirically (or very nearly so) for many splits of the data where  $n - m$  is  
 245 large. Setting  $Q(\theta) = \mathbf{1}_{[\theta=\theta_{\mathcal{D}_a}]}$ , the  $KL$  has only the contribution from the optimization over  $d$ :  
 246  $\mathbb{K}L(Q||P_{\mathcal{D}_a}) \leq \log D$  we find that the bound is nothing more than a variant of the simple Hoeffding  
 247 bound where  $\mathcal{D}_b$  is the validation set  $R(\theta_{\mathcal{D}_a}) \leq \hat{R}_{\mathcal{D}_b}(\theta_{\mathcal{D}_a}) + \sqrt{\frac{\log(Dm/\delta)+2}{2m-1}}$ .

248 This phenomenon can be seen in Appendix C, where optimization over  $d = 1, \dots, D$  leads to very  
 249 small values of the KL when  $n - m$  is a large fraction of  $n$ , and hence highly simple solutions  
 250 approaching the Hoeffding bound or confidence interval. This demonstrates that data-dependent  
 251 priors based PAC-Bayes bounds are not particularly meaningful, and attempts to optimize these  
 252 bounds will lead inevitably to the statistical bounds formed from the validation set such as a simple  
 253 Hoeffding bound or confidence interval. While these bounds are correct, they are only informative  
 254 about the difference in performance of the pretrained model on validation data and test data.

255 Note that this is in contrast with the data-independent bounds that constrain and explain generaliza-  
 256 tion for the entirety of the learning process. Similarly, our transfer learning bounds meaningfully  
 257 constrain what happens in the finetuning on the downstream task, but they do not constrain the prior  
 258 determined from the upstream task.

259 **5.3 Non-Vacuous PAC-Bayes Bounds for Transfer Learning**

260 By directly interpreting PAC-Bayes bounds through the lens of compression, we immediately see the  
 261 benefits of employing an upstream dataset for transfer learning. In fact, transfer learning allows us to  
 262 constrain the prior around parameters we know are consistent with the upstream task, and hopefully  
 263 the downstream task as well, reducing the KL-divergence between the prior and the posterior and  
 264 leading to even tighter bounds as we show in Table 2. We obtain unequivocally tighter bounds  
 265 using transfer learning for data-independent priors, providing a theoretical certification that transfer  
 266 learning can improve generalization. Our PAC-Bayes transfer learning approach also indicates that  
 267 transfer learning can boost generalization whenever codings optimized on a pre-training task are  
 268 more efficient for encoding a downstream posterior than an a priori guess made before seeing data. In  
 269 contrast, downstream tasks which differ greatly from the upstream task may only be consistent with  
 270 models that are not compressible under the learned prior, a scenario which may describe negative  
 271 transfer. As pre-training becomes the default in deep learning, our bounds open the door to providing  
 272 learning guarantees for state-of-the-art techniques. See Appendix B for more details.

273 **6 Understanding Generalization through PAC-Bayes Bounds**

274 Unlike the classical statistical learning theory viewpoint of uniform convergence, where the focus is  
 275 on the properties of the hypothesis class as a whole, the ability of a given network to be compressed  
 276 is highly dependent on the actual dataset that the model was trained on. In this respect, the PAC-  
 277 Bayes perspective shows us that the ability to compress and therefore generalize is not a function of  
 278 the hypothesis space but rather of individual functions: there are many networks in that space that  
 279 fit corrupted data and provably cannot be compressed (see Appendix A.1). But these models are not  
 280 the ones that we find when training on real world datasets, as real world datasets have a tremendous  
 281 amount of structure. If they did not, it would be pointless trying to learn from them as famously  
 282 argued by Hume [29] and the No Free Lunch theorem [54, 24]. When training on these structured  
 283 datasets that have low Kolmogorov complexity, the compressibility of the dataset is reflected in the  
 284 compressibility of the model. Breaking the structure in the dataset also breaks the structure in the  
 285 model and leads to worse generalization bounds for shuffled labels, shuffled pixels, and equivariance  
 286 more generally.

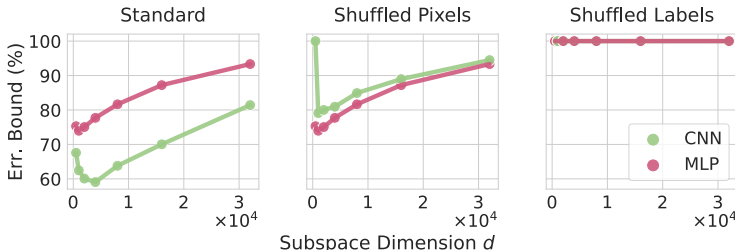


Figure 2: Our PAC-Bayes Bound vs subspace dimension for a fixed size convolutional network and MLP both with 500K parameters. With training on (Left) Standard CIFAR-10 (Middle) CIFAR-10 with shuffled pixels (Right) CIFAR-10 with shuffled labels. Structure in the dataset induces structure in the model. As structure is removed from the dataset, the models become much less compressible hence generalize worse. When the image structure is destroyed with pixel shuffling, we find that the CNN generalizes worse than an MLP.

287 **MLPs vs CNNs:** It is well known that convolutional neural networks (CNNs) generalize much  
 288 better than standard multilayer perceptrons (MLPs) even when controlling for the number of param-  
 289 eters. We find that this fact is reflected in the improved compressibility of CNNs and our general-  
 290 ization bounds, shown in Figure 2 (a) on the image dataset CIFAR-10 using the same number of  
 291 parameters for the two models. The CNN is able find a more structured and lower description length  
 292 explanation of the training data.

293 **Shuffled Pixels:** However, when this image structure is broken by shuffling the pixels we find that  
 294 CNNs are no better at generalizing than MLPs, and this is reflected in the compressed sizes. The  
 295 CNNs become substantially less compressible when trained on this dataset, and hence our bounds  
 296 show them generalizing worse than MLPs, see Figure 2 (b). Not using this image structure (and

297 invariant to a global permutation of the features), MLPs do not suffer when this structure is broken  
298 since they never used it in the first place.

299 **Shuffled Labels:** When the structure of the dataset is entirely broken by shuffling the labels, the  
300 compressibility of the models (both for CNNs and MLPs) are lost. Regardless of the subspace  
301 dimension used, our generalization bounds are all at 100% error. Using low subspace dimensions it  
302 is not possible to fit the training data, and when using a large enough dimension to fit the data the  
303 compressed size of the model is larger than the training data and hence the generalization bounds  
304 are vacuous.

305 **Equivariance:** Designing models which are *equivariant* to certain symmetry transformations has  
306 been a guiding principle for the development of data efficient neural networks in many domains  
307 [12, 13, 52, 53, 21, 31]. While intuitively it is clear that respecting dataset symmetries severely  
308 improves generalization, relatively little has been proven for neural networks [20, 60, 4, 19]. We  
309 compress and evaluate rotationally equivariant and non-equivariant Wide ResNets [53, 57] trained  
310 on MNIST and a rotated version of MNIST. As shown in [Appendix E](#), the rotationally equivariant  
311 models are more compressible and provably generalize better than their non-equivariant counterparts  
312 when paired with a dataset that also has the rotational symmetry.

313 **Is Stochasticity Necessary for Generalization?** It is widely hypothesized that the implicit biases  
314 of SGD help find solutions which generalize better. For example, Arora et al. [1] argue that there  
315 is no regularizing function with a gradient that replicates the benefits of gradient noise. Wu et al.  
316 [55], Smith et al. [49], and Li et al. [40] advocate that gradient noise is necessary to achieve state-  
317 of-the-art performance. In contrast, recent work by Geiping et al. [23] shows that full-batch gradient  
318 descent can match state-of-the-art performance.

319 We train ResNet-18 and LeNet5 models on CIFAR-10 and MNIST, respectively, using full-batch  
320 and stochastic gradient descent with different values of the intrinsic dimensionality. We provide the  
321 training details in [Appendix F](#), but it is worth noting that, unlike Geiping et al. [23], we do not use  
322 flatness inducing regularization for full batch training. For MNIST with LeNet5, the best general-  
323 ization bounds that we obtain are 11.55% and 11.20% using stochastic gradient descent (SGD) and  
324 full batch training respectively. The best generalization bounds that we obtain for CIFAR-10 with  
325 ResNet-18 are 74.68% and 76.25% using SGD and full batch training, respectively. These close  
326 theoretical guarantees on the generalization error for both SGD and full batch training suggest that  
327 the implicit biases of SGD may be helpful, but are not at all crucial to understanding why neural  
328 networks generalize well. We expand this analysis in [Appendix F](#).

329 **Double Descent [44]:** In [Appendix H](#), we show that our bounds are tight and fine-grained enough  
330 to predict the double descent phenomenon, and we discuss the significance of this result.

## 331 7 Discussion

332 In this work, we have constructed a new method for compressing deep learning models which is  
333 highly adaptive to the level of structure in the model and the training dataset. According to the  
334 Occam’s prior which considers shorter compressed length models to be more likely, we construct  
335 state-of-the-art generalization bounds across a variety of settings. Through our bounds and the  
336 viewpoint of Occam’s razor and compressibility, we show how generalization relates to the structure  
337 in the dataset and the structure in the model, and we are able to explain aspects of neural network  
338 generalization for natural image datasets, shuffled pixels, shuffled labels, equivariant models, and  
339 for stochastic training.

340 **Limitations** We note that despite the power of this compression scheme and our bounds, there are  
341 certain aspects of deep learning generalization that remain unexplained. Our compression bounds  
342 still show a preference towards models with a smaller number of base parameters as shown in [Ap-  
343 ppendix G](#), despite the fact that larger models tend to generalize better. While we have achieved  
344 better model compression than previous works, it is unlikely that we are close to theoretical limits.  
345 We hope that with a further improved compression scheme (possibly using nonlinear parameter pro-  
346 jections), we might find that larger deep learning model solutions are compressible to smaller sizes  
347 than smaller models. Additionally, while our bounds show that the compressibility of our models  
348 implies generalization, we have no statements in the reverse direction. We believe however that the  
349 compression approach and Occam’s razor have yet untapped explanatory power in deep learning.

## References

- 350
- 351 [1] Sanjeev Arora, Nadav Cohen, and Elad Hazan. On the Optimization of Deep Networks: Im-  
352 plicit Acceleration by Overparameterization. *Proceedings of the 35th International Conference*  
353 *on Machine Learning (ICML)*, 2018. 9
- 354 [2] Sanjeev Arora, Rong Ge, Behnam Neyshabur, and Yi Zhang. Stronger generalization bounds  
355 for deep nets via a compression approach. In *International Conference on Machine Learning*,  
356 pages 254–263. PMLR, 2018. 1, 2
- 357 [3] Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or Propagating Gradients  
358 through Stochastic Neurons for Conditional Computation. *Preprint arXiv:1308.3432v1*, 2013.  
359 5
- 360 [4] Alberto Bietti, Luca Venturi, and Joan Bruna. On the sample complexity of learning under  
361 geometric stability. *Advances in Neural Information Processing Systems*, 34, 2021. 9
- 362 [5] Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncer-  
363 tainty in neural network. In *International conference on machine learning*, pages 1613–1622.  
364 PMLR, 2015. 2
- 365 [6] Olivier Catoni. PAC-Bayesian Supervised Classification: the Thermodynamics of Statistical  
366 Learning. *Institute of Mathematical Statistics Lecture Notes*, 2007. 3, 15
- 367 [7] Pratik Chaudhari, Anna Choromanska, Stefano Soatto, Yann LeCun, Carlo Baldassi, Christian  
368 Borgs, Jennifer Chayes, Levent Sagun, and Riccardo Zecchina. Entropy-sgd: Biasing gradient  
369 descent into wide valleys. *Journal of Statistical Mechanics: Theory and Experiment*, 2019  
370 (12):124018, 2019. 2
- 371 [8] Yu Cheng, Duo Wang, Pan Zhou, and Tao Zhang. A survey of model compression and accel-  
372 eration for deep neural networks. *arXiv preprint arXiv:1710.09282*, 2017. 4
- 373 [9] Yu Cheng, Duo Wang, Pan Zhou, and Tao Zhang. Model compression and acceleration for deep  
374 neural networks: The principles, progress, and challenges. *IEEE Signal Processing Magazine*,  
375 35(1):126–136, 2018. 2, 4
- 376 [10] Yoojin Choi, Mostafa El-Khamy, and Jungwon Lee. Towards the limit of network quantization.  
377 *arXiv preprint arXiv:1612.01543*, 2016. 5
- 378 [11] Yoojin Choi, Mostafa El-Khamy, and Jungwon Lee. Towards the Limit of Network Quantiza-  
379 tion. *The 6th International Conference on Learning Representations (ICLR)*, 2017. 5
- 380 [12] Taco Cohen and Max Welling. Group equivariant convolutional networks. In *International*  
381 *conference on machine learning*, pages 2990–2999. PMLR, 2016. 9
- 382 [13] Taco S Cohen, Mario Geiger, Jonas Köhler, and Max Welling. Spherical cnns. *arXiv preprint*  
383 *arXiv:1801.10130*, 2018. 9
- 384 [14] Yann Dauphin, Razvan Pascanu, Çağlar Gülçehre, Kyunghyun Cho, Surya Ganguli, and  
385 Yoshua Bengio. Identifying and attacking the saddle point problem in high-dimensional non-  
386 convex optimization. In *NIPS*, 2014. 4
- 387 [15] Nan Ding, Xi Chen, Tomer Levinboim, Beer Changpinyo, and Radu Soricut. Pactran: Pac-  
388 bayesian metrics for estimating the transferability of pretrained models to classification tasks.  
389 *arXiv preprint arXiv:2203.05126*, 2022. 3
- 390 [16] Vincent Dumoulin, Ethan Perez, Nathan Schucher, Florian Strub, Harm de Vries, Aaron  
391 Courville, and Yoshua Bengio. Feature-wise transformations. *Distill*, 3(7):e11, 2018. 5
- 392 [17] Gintare Karolina Dziugaite and Daniel M Roy. Computing nonvacuous generalization bounds  
393 for deep (stochastic) neural networks with many more parameters than training data. *arXiv*  
394 *preprint arXiv:1703.11008*, 2017. 2
- 395 [18] Gintare Karolina Dziugaite, Kyle Hsu, Waseem Gharbieh, Gabriel Aprino, and Daniel M. Roy.  
396 On the Role of Data in Pac-Bayes Bounds. *The 24th International Conference on Artificial*  
397 *Intelligence and Statistics (AISTATS)*, 2021. 2, 7
- 398 [19] Bryn Elesedy. Group symmetry in pac learning. In *ICLR 2022 Workshop on Geometrical and*  
399 *Topological Representation Learning*, 2022. 9
- 400 [20] Bryn Elesedy and Sheheryar Zaidi. Provably strict generalisation benefit for equivariant mod-  
401 els. In *International Conference on Machine Learning*, pages 2959–2969. PMLR, 2021. 9

- 402 [21] Marc Finzi, Samuel Stanton, Pavel Izmailov, and Andrew Gordon Wilson. Generalizing con-  
403 volutional neural networks for equivariance to lie groups on arbitrary continuous data. In  
404 *International Conference on Machine Learning*, pages 3165–3176. PMLR, 2020. 9
- 405 [22] T. Garipov, Pavel Izmailov, Dmitrii Podoprikin, Dmitry P. Vetrov, and Andrew Gordon Wil-  
406 son. Loss surfaces, mode connectivity, and fast ensembling of dnns. In *NeurIPS*, 2018. 4
- 407 [23] Jonas Geiping, Micah Goldblum, Phillip E. Pope, Michael Moeller, and Tom Goldstein.  
408 Stochastic Training Is Not Necessary For Generalization. *The 10th International Conference*  
409 *on Learning Representations (ICLR)*, 2022. 9
- 410 [24] Christophe Giraud-Carrier and Foster Provost. Toward a justification of meta-learning: Is the  
411 no free lunch theorem a show-stopper. In *Proceedings of the ICML-2005 Workshop on Meta-*  
412 *learning*, pages 12–19, 2005. 8
- 413 [25] Ian J. Goodfellow and Oriol Vinyals. Qualitatively characterizing neural network optimization  
414 problems. *CoRR*, abs/1412.6544, 2015. 4
- 415 [26] Song Han, Huizi Mao, and William J. Dally. Deep Compression: Compressing Deep Neural  
416 Networks with Pruning, Trained Quantization and Huffman Coding. *The 4th International*  
417 *Conference on Learning Representations (ICLR)*, 2016. 2, 5
- 418 [27] Geoffrey E Hinton and Drew Van Camp. Keeping the neural networks simple by minimiz-  
419 ing the description length of the weights. In *Proceedings of the sixth annual conference on*  
420 *Computational learning theory*, pages 5–13, 1993. 2
- 421 [28] Sepp Hochreiter and Jürgen Schmidhuber. Flat minima. *Neural computation*, 9(1):1–42, 1997.  
422 1, 2
- 423 [29] David Hume. *A Treatise of Human Nature*. Oxford University Press, Oxford, 1978. revised  
424 P.H. Nidditch. 8
- 425 [30] Marcus Hutter et al. Algorithmic complexity. 2008. 3
- 426 [31] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ron-  
427 neberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Židek, Anna Potapenko, et al.  
428 Highly accurate protein structure prediction with alphafold. *Nature*, 596(7873):583–589, 2021.  
429 9
- 430 [32] Fahdi Kanavati and Masayuki Tsuneki. Partial transfusion: on the expressive influence of  
431 trainable batch norm parameters for transfer learning. In *Medical Imaging with Deep Learning*,  
432 pages 338–353. PMLR, 2021. 5
- 433 [33] Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping  
434 Tak Peter Tang. On large-batch training for deep learning: Generalization gap and sharp min-  
435 ima. *arXiv preprint arXiv:1609.04836*, 2016. 2
- 436 [34] John Langford. *Quantitatively tight sample complexity bounds*. PhD thesis, Carnegie Mellon  
437 University, 2002. 2
- 438 [35] John Langford and Rich Caruana. (not) bounding the true error. *Advances in Neural Informa-*  
439 *tion Processing Systems*, 14, 2001. 2
- 440 [36] John Langford and Matthias Seeger. Bounds for Averaging Classifiers. *Tech. rep CMU-CS-*  
441 *01-102, Carnegie Mellon University*, 2001. 2, 3, 15
- 442 [37] Quoc V. Le, Tamás Sarlós, and Alex Smola. Fastfood: Approximate kernel expansions in  
443 loglinear time. *ArXiv*, abs/1408.3060, 2013. 4
- 444 [38] Chunyuan Li, Heerad Farkhoor, Rosanne Liu, and Jason Yosinski. Measuring the intrinsic  
445 dimension of objective landscapes. *ArXiv*, abs/1804.08838, 2018. 4, 5
- 446 [39] Ping Li, Trevor J. Hastie, and Kenneth Ward Church. Very sparse random projections. In *KDD*  
447 *'06*, 2006. 4
- 448 [40] Zhiyuan Li, Sadhika Malladi, and Sanjeev Arora. On the Validity of Modeling SGD with  
449 Stochastic Differential Equations (SDEs). *Preprint arXiv:2102.12470*, 2021. 9
- 450 [41] David JC MacKay, David JC Mac Kay, et al. *Information theory, inference and learning*  
451 *algorithms*. Cambridge university press, 2003. 5
- 452 [42] Andreas Maurer. A Note on the PAC Bayesian Theorem. *Preprint arXiv 041099v1*, 2004. 3,  
453 15

- 454 [43] David A McAllester. PAC-Bayesian Model Averaging. *Proceedings of the 12th Annual Con-*  
455 *ference on Learning Theory (COLT)*, pages 164–170, 1999. 3, 15
- 456 [44] Preetum Nakkiran, Gal Kaplun, Yamini Bansal, Tristan Yang, Boaz Barak, and Ilya Sutskever.  
457 Deep double descent: where bigger models and more data hurt. *Journal of Statistical Mechan-*  
458 *ics: Theory and Experiment*, 2021, 2020. 9
- 459 [45] Behnam Neyshabur, Srinadh Bhojanapalli, and Nathan Srebro. A pac-bayesian approach to  
460 spectrally-normalized margin bounds for neural networks. *arXiv preprint arXiv:1707.09564*,  
461 2017. 1, 2
- 462 [46] Ethan Perez, Florian Strub, Harm De Vries, Vincent Dumoulin, and Aaron Courville. Film:  
463 Visual reasoning with a general conditioning layer. In *Proceedings of the AAAI Conference on*  
464 *Artificial Intelligence*, volume 32, 2018. 5
- 465 [47] María Pérez-Ortiz, Omar Rivasplata, John Shawe-Taylor, and Csaba Szepesvári. Tighter Risk  
466 Certificates for Neural Networks. *Journal of Machine Learning Research 22 (2021) 1-40*,  
467 2021. 2, 7
- 468 [48] Omar Rivasplata, Vikram M Tankasali, and Csaba Szepesvári. Pac-bayes with backprop. *arXiv*  
469 *preprint arXiv:1908.07380*, 2019. 2
- 470 [49] Samuel L. Smith, Erich Elsen, and Soham De. On the Generalization Benefit of Noise in  
471 Stochastic Gradient Descent. *Proceedings of the 37th International Conference on Machine*  
472 *Learning (ICML)*, 2020. 9
- 473 [50] Ray J Solomonoff. A formal theory of inductive inference. part i. *Information and control*, 7  
474 (1):1–22, 1964. 3
- 475 [51] Niklas Thiemann, Christian Igel, Oliver Wintenberger, and Yevgeny Seldin. A Strongly Qua-  
476 siconvex PAC-Bayesian Bound. *28th Annual Conference on Learning Theory (COLT)*, 2017.  
477 2, 3, 15
- 478 [52] Nathaniel Thomas, Tess Smidt, Steven Kearnes, Lusann Yang, Li Li, Kai Kohlhoff, and Patrick  
479 Riley. Tensor field networks: Rotation-and translation-equivariant neural networks for 3d point  
480 clouds. *arXiv preprint arXiv:1802.08219*, 2018. 9
- 481 [53] Maurice Weiler and Gabriele Cesa. General e (2)-equivariant steerable cnns. *Advances in*  
482 *Neural Information Processing Systems*, 32, 2019. 9
- 483 [54] David H Wolpert and William G Macready. No free lunch theorems for optimization. *IEEE*  
484 *transactions on evolutionary computation*, 1(1):67–82, 1997. 8
- 485 [55] Jingfeng Wu, Wenqing Hu, Haoyi Xiong, Jun Huan, Vladimir Braverman, and Zhanxing Zhu.  
486 On the Noisy Gradient Descent that Generalizes as SGD. *Proceedings of the 37th International*  
487 *Conference on Machine Learning (ICML)*, 2020. 9
- 488 [56] Penghang Yin, Jiancheng Lyu, Shuai Zhang, Stanley Osher, Yingyong Qi, and Jack Xin. Un-  
489 derstanding Straight-Through Estimator in Training Activation Quantized Neural Nets. *The*  
490 *7th International Conference on Learning Representations (ICLR)*, 2019. 5
- 491 [57] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv preprint*  
492 *arXiv:1605.07146*, 2016. 9
- 493 [58] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understand-  
494 ing deep learning (still) requires rethinking generalization. *Communications of the ACM*, 64  
495 (3):107–115, 2021. 1
- 496 [59] Wenda Zhou, Victor Veitch, Morgane Austern, Ryan P. Adams, and Peter Orbanz. Non-  
497 Vacuous Generalization Bounds at the ImageNet Scale: A PAC-Bayesian Compression Ap-  
498 proach. *7th International Conference on Learning Representations (ICLR)*, 2019. 1, 2, 3, 4,  
499 6
- 500 [60] Sicheng Zhu, Bang An, and Furong Huang. Understanding the generalization benefit of model  
501 invariance from a data perspective. *Advances in Neural Information Processing Systems*, 34,  
502 2021. 9

503 **Checklist**

- 504 1. For all authors...
- 505 (a) Do the main claims made in the abstract and introduction accurately reflect the pa-  
506 per’s contributions and scope? [Yes] We provide evidence for the claims made in the  
507 abstract and introduction in Section 4, Section 5.1, and Section 6.
- 508 (b) Did you describe the limitations of your work? [Yes] We describe limitations in Sec-  
509 tion 7.
- 510 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 511 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
512 them? [Yes]
- 513 2. If you are including theoretical results...
- 514 (a) Did you state the full set of assumptions of all theoretical results? [Yes] See Ap-  
515 pendix I.
- 516 (b) Did you include complete proofs of all theoretical results? [Yes] See Appendix I.
- 517 3. If you ran experiments...
- 518 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
519 mental results (either in the supplemental material or as a URL)? [Yes]
- 520 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
521 were chosen)? [Yes] See Appendix D.
- 522 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
523 ments multiple times)? [Yes] See Appendix.
- 524 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
525 of GPUs, internal cluster, or cloud provider)? [Yes] See Appendix D.
- 526 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 527 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 528 (b) Did you mention the license of the assets? [Yes] See Appendix J.
- 529 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]  
530 Our code is submitted in the supplemental material.
- 531 (d) Did you discuss whether and how consent was obtained from people whose data  
532 you’re using/curating? [N/A]
- 533 (e) Did you discuss whether the data you are using/curating contains personally identifi-  
534 able information or offensive content? [N/A]
- 535 5. If you used crowdsourcing or conducted research with human subjects...
- 536 (a) Did you include the full text of instructions given to participants and screenshots, if  
537 applicable? [N/A]
- 538 (b) Did you describe any potential participant risks, with links to Institutional Review  
539 Board (IRB) approvals, if applicable? [N/A]
- 540 (c) Did you include the estimated hourly wage paid to participants and the total amount  
541 spent on participant compensation? [N/A]