

LARGER LANGUAGE MODELS PROVABLY GENERALIZE BETTER

Anonymous authors

Paper under double-blind review

ABSTRACT

Why do larger language models generalize better? To explore this question, we develop generalization bounds on the pretraining objective of large language models (LLMs) in the compute-optimal regime, as described by the Chinchilla scaling laws. We introduce a novel, fully empirical Freedman-type martingale concentration inequality that tightens existing bounds by accounting for the variance of the loss function. The generalization bound can be broken into three contributions: the number of parameters per token, the loss variance, and the quantization error at a fixed bitrate. As language models are scaled up, the number of parameters per data point stays constant; however, both the loss variance and the quantization error decrease, implying that larger models should have *smaller* generalization gaps. We examine why larger models tend to be more quantizable from an information theoretic perspective, showing that the rate at which they can integrate new information grows slower than their capacity on the compute optimal frontier. From these findings we produce a scaling law for the generalization gap, showing that our bounds decrease in a predictable way.

1 INTRODUCTION

Large language models (LLMs) have demonstrated a remarkable general purpose problem solving capacity across a wide range of complex tasks which humans are able to perform, from classical NLU (Brown, 2020), forecasting (Gruber et al., 2023), mathematics (Trinh et al., 2024), spatial reasoning (Patel & Pavlick, 2022), and many other areas. For a large majority of individual tasks, model capabilities increase monotonically as the next token prediction loss from the pretraining objective decreases.

A conceptually useful story about the learning process involves the model accommodating predictive subprograms of progressively larger computational depth and complexity. During pretraining, shallow details are absorbed first: the log likelihood loss is most easily decreased by learning word frequencies, syntax, and grammar. As these details are absorbed, slightly higher level structures such as facts, relations, and idioms then become the next batch of the lowest hanging fruit, giving way to yet higher level structures and pattern matching.

For reasons that are not yet well understood, this process is reflected in the pretraining objective as a power law for LLMs and other generative models on natural data. The frontier of the best achievable performance given a fixed computational budget C obeys a fixed and predictable power law relationship $L(C) \propto C^{-\alpha}$ over many orders of magnitude (Kaplan et al., 2020). This relationship appears to vary considerably with the kind of data (Henighan et al., 2020) (e.g. natural text vs images vs math) and only weakly on the details of the model architecture and training method (Bahri et al., 2021).

Effort in quantifying *what* this relationship is in a given domain and *how* it varies as model size and dataset size are traded off has been extremely valuable in guiding where resources are spent in constructing more capable AI models (Brown, 2020; Besiroglu et al., 2024; OpenAI, 2023; Dubey et al., 2024) and charting a path for the future. In this work, we target the *why* of scaling laws. While mathematically simple toy models or toy data are valuable, we aim to study the *why* of scaling laws on real models and real data by focusing on just one contribution to the scaling law curve: the token-wise generalization gap. Constructing a generalization bound sensitive enough to capture the

small differences between architectures and yet simple enough to write down in a short formula is likely impossible; however, even the broad strokes of behavior such as how generalization scales with compute have not been addressed. Thus, here we focus high level understanding rather than algorithmic intervention.

In order to construct the relevant generalization bounds, we introduce a novel empirical Freedman concentration inequality (Freedman, 1975). Our generalization bound highlights three critical components—the ratio of parameters per token in compute-optimal scaling (which is roughly constant), the token-wise loss variance (which decreases with model size), and the performance gap between quantized and unquantized models (which also decreases with model size). As an alternative to quantization, we bound the information transfer between dataset and the model, showing that the information content in the model grows sublinearly with model size, and thus the complexity decreases with model size. These components collectively contribute to a predictable reduction in the generalization gap as models grow larger.

2 BACKGROUND

2.1 GENERALIZATION BOUNDS

At a high level, we are interested in the expected test error (population risk) $\mathbb{E}_{X' \sim p_{\mathcal{D}}}[R_{h(X)}(X')]$ for a given model (hypothesis) h depending on the training set X but evaluated on a test set X' sampled from the data distribution $p_{\mathcal{D}}$. One conceptually convenient way of breaking down this quantity is into the irreducible error, approximation gap, and generalization gap:¹

$$\mathbb{E}_{X' \sim p_{\mathcal{D}}}[R_{h(X)}(X')] = \underbrace{R_*(X)}_{\text{Irreducible Error } E} + \underbrace{R_{h(X)}(X) - R_*(X)}_{\text{Approximation Gap } A} + \underbrace{\mathbb{E}_{X' \sim p_{\mathcal{D}}}[R_{h(X)}(X')] - R_{h(X)}(X)}_{\text{Generalization Gap } G}.$$

The first term describes the entropy of natural text, e.g. the amount of truly random information content in the data, which cannot be further explained even when knowing the true data generating process. The second term describes the approximation gap, capturing the extent to which the trained model is able to fit the training data. This term combines both model capacity, e.g. as described by universal approximation theorems (Cybenko, 1989), as well as optimization via how well the training algorithm is able to find the given solution. Finally, we have the generalization gap, capturing the extent to which training and testing performance diverge on account of overfitting to the statistically irrelevant regularities in X . Though generalization bounds focus on the last term, all three quantities are of interest for understanding LLM behavior. Empirically, it has been observed that the generalization gap for LLMs (at least in the low epoch regime) tends to be extremely small or even negligible compared to the other two terms.

Among the simplest generalization bounds is the finite hypothesis with prior generalization bound applied to IID data (Shalev-Shwartz & Ben-David, 2014). With probability at least $1 - \delta$,

$$\mathbb{E}_{X' \sim p_{\mathcal{D}}}[R_{h(X)}(X')] - R_{h(X)}(X) \leq \Delta \sqrt{\frac{\log 1/P(h) + \log 1/\delta}{2m}}$$

where m is the number of IID data points, Δ is an upper bound on the range of values the risk can take, and $P(h)$ is a prior distribution over hypotheses in a discrete hypothesis class \mathcal{H} . With a judicious choice of prior, $\log 1/P(h)$ can be related to the compressed size of the model measured in nats (Lotfi et al., 2022).

During text pretraining, the individual tokens are not sampled IID. Thus, a generalization bound requires treating entire documents (often thousands of tokens) as the elements the empirical risk is computed over. Note that modern language models have hundreds of times more parameters than documents they were trained on. With the help of very extreme compression methods and using smoothing to bound Δ , it is possible to construct nonvacuous bounds (Lotfi et al., 2024a). However, the required compression (greater than 100 times) is so severe that it cripples model performance.

¹We note this differs from the commonly referred to estimation-approximation error breakdown (Bottou & Bousquet, 2007) or the bias-variance decomposition (Brown & Ali, 2024); however, the train error-generalization gap is more useful for our purposes.

In a recent work, [Lotfi et al. \(2024c\)](#) explore breaking down generalization into tokenwise generalization, e.g. how the loss varies with each individual predicted token being resampled under the distribution but keeping the context the same. Splitting up the training dataset X into the sequence of tokens $[X_k]_{k=1}^D$, the authors bound

$$T = \frac{1}{D} \sum_{k=1}^D \mathbb{E}[R_h(X_k | X_{<k}) | X_{<k}] - R_h(X),$$

where $R_h(X_k | X_{<k})$ is the negative log likelihood for token k given the context $X_{<k}$, and the expectation is taken with respect to $p(X_k | X_{<k})$ from the data distribution. The authors bound T using Azuma’s inequality to arrive at a bound scaling as $\Delta \sqrt{\frac{\log 1/P(h)}{2D}}$. We improve upon this bound, reducing to a leading term of $\Delta \frac{\log 1/P(h)}{D}$.

2.2 CHINCHILLA SCALING LAWS

A key insight from the current machine learning paradigm is that the dataset should not be considered a fixed quantity. Rather than optimizing to find the best model for a given dataset, one should instead try to find the best performing model and dataset for a given computational budget. [Muennighoff et al. \(2024\)](#) describes the optimal allocation of resources for increasing the size of the model and increasing the size of the dataset under the assumption that data is plentiful relative to the computational budget.

Let N be the number of parameters and D be the number of training tokens. In the one epoch regime of LLM pretraining, the negative log likelihood loss is well-approximated by the power law

$$R(N, D) = E + \frac{A}{N^\alpha} + \frac{B}{D^\beta},$$

where A, B are empirically estimated constants, exponents α, β have similar values, and E is the irreducible error. Optimizing $N(C)$ and $D(C)$ under the constraint of a fixed compute budget $C \approx 6ND$ ([Kaplan et al., 2020](#)), one arrives at

$$N^*(C) = G(C/6)^a \quad D^*(C) = G^{-1}(C/6)^b$$

for constants G, a, b that depend on A, B, α, β ([Hoffmann et al., 2022](#)).

Within the margin of statistical error, we have $a = b = 0.5$ in the optimal allocation of compute ([Besiroglu et al., 2024](#)). Therefore, the ratio of parameters per token, $N^*(C)/D^*(C) = G^2$, is a fixed constant. Evaluating the constants from [Muennighoff et al. \(2024\)](#), we have $G^2 \approx 1/20$. Note that many open source models optimize performance amortized over both training time and inference time compute, which leads to smaller than Chinchilla optimal models, e.g. models with a ratio $N/D < G^2$. In the context of this paper, we will assume the Chinchilla optimal scaling $N/D = G^2$ with the understanding that any generalization bounds we construct would only be tighter if the ratio N/D is smaller.

To test our theory, we use the open source Pythia model family ([Biderman et al., 2023](#)) ranging from 70 million to 12 billion parameters. Unlike other open source LLMs, we have full access to both the Pythia model checkpoints from training and the Pile dataset they were trained on ([Gao et al., 2020](#)). From these intermediate checkpoints, we choose the set of models along the compute optimal frontier to match $N/D = 1/20$, reflecting the choice for number of training steps and model size that one would have made optimizing only for performance at the given computational budget. The chosen checkpoints are plotted in the training frontier of these models in [Figure 1](#).

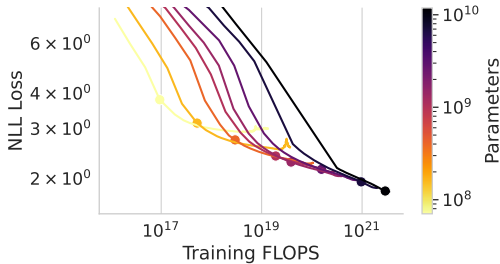


Figure 1: Pythia models and checkpoints chosen along the compute optimal frontier (checkpoints given by the marked values).

3 GENERALIZATION BOUND

In this section, we will build up the components for our generalization bounds. To capture the relevant behavior, we derive a new concentration inequality for martingales. We apply a prior weighted union bound to this concentration inequality so that we can apply it to models in a large hypothesis class, taking advantage of the low complexity inherent in compressible models. Bounding the worst case loss behavior using prediction smoothing, we apply this bound to LLMs.

3.1 AN EMPIRICAL FREEDMAN’S CONCENTRATION INEQUALITY

Theorem 3.1. *Let X_1, \dots, X_n be a sequence of \mathcal{F}_k -measurable random variables. Let Y_1, \dots, Y_n be any other sequence of \mathcal{F}_{k-1} -measurable random variables such that the difference is bounded below: $A_k = (Y_k - X_k)/\Delta > -1$ for some $\Delta \geq 0$. Let K be a finite subset of $(0, 1)$. Then, with probability at least $1 - \delta$,*

$$\frac{1}{n} \sum_{k=1}^n (\mathbb{E}[X_k | \mathcal{F}_{k-1}] - X_k) \leq \Delta C + \Sigma \sqrt{C}, \quad (1)$$

where $C := \frac{1}{n} \log \frac{|K|}{\delta}$, and

$$\Sigma = \Sigma(C, \Delta, \{X_k - Y_k\}_{k=1}^n, K) := \min_{s \in K} \Delta \sqrt{C} (1 - s)/s + \frac{\Delta}{\sqrt{C}} \frac{1}{n} \sum_{k=1}^n v(sA_k)/s$$

and $v(x) = x - \log(1 + x)$.

Proof: See [Section A.1](#). This concentration inequality states that the sequence of random variables concentrates on their conditional means with a term Σ depending on the empirical variation of the loss value. We note that Σ can be viewed as a variance term. As we show in [Appendix A](#), using a small K , the variance proxy can be upper bounded: $\Sigma \leq 2\sqrt{\frac{1}{n} \sum_{k=1}^n (X_k - Y_k)^2}$. This is explicitly related to the empirical variance but with the mean replaced by Y_k . Although the minimization form above is unwieldy, it produces significantly tighter estimates of Σ (a factor of $\sim 5x$ smaller). When the loss variation Σ is small, concentration happens at a rate linear in the complexity C rather the slower \sqrt{C} rate.

This concentration inequality provides the core result for our generalization bounds, and to the best of our knowledge it is the first martingale concentration inequality to incorporate a variance term which can be evaluated on the original data. We can view this bound as aiming to achieve the benefits that Freedman’s inequality has over Azuma’s inequality while being fully empirical, replacing the population variance with a fully empirical proxy. Our approach is analogous to the fully empirical Bernstein bound derived in [Maurer & Pontil \(2009\)](#), but in the martingale rather than IID setting. Unfortunately, the proof technique of [Maurer & Pontil \(2009\)](#) does not carry over to the martingale case and instead we take quite a different approach. We derive our concentration inequality in [Theorem A.5](#) making use of a proxy Y_k that is \mathcal{F}_{k-1} -measurable but which can take the place of $\mathbb{E}[X_k | \mathcal{F}_{k-1}]$ in the variance. In practice, we choose this quantity to be the mean of the model NLL under resampling of the given token according to the *model* distribution in place of the data distribution.

3.2 EXTENDING TO A DISCRETE HYPOTHESIS CLASS

From the concentration inequality in [Equation 1](#), we derive the following discrete hypothesis class generalization bound.

Theorem 3.2. *Let X_1, \dots, X_n be a sequence of (possibly dependent) random variables. Let $R_h(X_k | X_{<k})$ denote the risk for element X_k given the previous elements of the sequence for hypothesis h in a countable hypothesis class \mathcal{H} . Let $p_h(X_k | X_{<k})$ be any (hypothesis dependent) distribution over X_k conditioned on $X_{<k}$. Consider a prefix free coding of each $h \in \mathcal{H}$ and let $L(h)$ be the length of that code measured in nats. Let K be a finite subset of \mathbb{R}^+ . Assuming $R_h(X_k | X_{<k}) - \mathbb{E}_{Y_k \sim p_h}[R_h(Y_k | X_{<k})] \leq \Delta$ for some $\Delta > 0$, we have that simultaneously for all $h \in \mathcal{H}$, with probability at least*

$$1 - \delta, \quad \frac{1}{n} \sum_k \mathbb{E}[R_h(X_k | X_{<k}) | X_{<k}] \leq \frac{1}{n} \sum_k R_h(X_k | X_{<k}) + \Delta C + \Sigma \sqrt{C}, \quad (2)$$

where the complexity C is given by

$$C := \frac{L(h) + \log |K|/\delta}{n}$$

and $\Sigma = \Sigma(C, \Delta, \{A_k\}_{k=1}^n, K)$ from [Theorem 3.1](#) for $A_k = R_h(X_k | X_{<k}) - \mathbb{E}_{Y_k \sim p_h}[R_h(Y_k | X_{<k})]$.

Proof: See [Section A.2](#).

3.3 WORST CASE BEHAVIOR AND SMOOTHING

The last component of our bounds is the smoothing to bound the worst case behavior of the model, which in general for the negative log likelihood can be arbitrarily large. We employ the prediction smoothing idea from [Lotfi et al. \(2024a\)](#), where the model is mixed with a uniform distribution over the tokens with a given mixing parameter. Unlike application in previous work, we optimize over this parameter analytically so that we can remove it from the bounds and evaluation entirely, instead merely as a tool for constructing bounds while all evaluations are with the unsmoothed model.

Theorem 3.3. *For the categorical negative log likelihood objective $\hat{R}_h = -\frac{1}{n} \sum_{k=1}^n \log p_h(X_k | X_{<k})$ on V classes and $C \in \mathbb{R}^+$, there exists a prediction smoothed model $p_s(\cdot) = (1 - \alpha)p_h(\cdot) + \alpha/V$ which has a worst case loss $\Delta_s = \sup_{X_k, X_{<k}} -\log p_s(X_k | X_{<k}) \leq \log(V/\alpha)$, and the risk satisfies*

$$\hat{R}_s + C \Delta_s \leq \hat{R}_h + C \log V + \sqrt{2C}, \quad (3)$$

for some value $\alpha \in (0, 1)$ (approximately $C/(1 + C)$).

The proof is provided in [Section A.3](#).

3.4 GENERALIZATION BOUND FOR CHINCHILLA LANGUAGE MODELS

Finally, we assemble these three components into a generalization bound that we can empirically evaluate for language models. Combining the prediction smoothing bound with [Theorem 3.2](#) applied to the smoothed quantized model produces the result of [Theorem 3.4](#). Note that each term in the expression has an interpretable meaning.

Theorem 3.4. *Let X_1, \dots, X_D be the sequence of D (possibly dependent) tokens formed from concatenating each sequence in the dataset together into a single stream of tokens. Let $R_h(X_k | X_{<k}) = -\log p_h(X_k | X_{<k})$ denote the NLL for element X_k given the previous elements for a given model h and with vocabulary size V and N parameters. Let $\hat{R}_h = \frac{1}{D} \sum_{k=1}^D R_h(X_k | X_{<k})$ be the empirical risk and $R_h = \frac{1}{D} \sum_{k=1}^D \mathbb{E}[R_h(X_k | X_{<k}) | X_{<k}]$ be the tokenwise expected risk for that model. Let K be a finite subset of $(0, 1)$. For any given quantization q of h using b bits per parameter with expected risk R_q , there exists a label smoothed model s_q with $R_{s_q}(X_k | X_{<k}) = (1 - \alpha)R_q(X_k | X_{<k}) + \alpha/V$ for fixed $\alpha \in (0, 1)$ which achieves a tokenwise population risk with probability $1 - \delta$*

$$R_{s_q} \leq \hat{R}_h + \underbrace{C \log V}_{\text{Random Guess NLL}} + \underbrace{\Sigma \sqrt{C}}_{\text{Loss Variation}} + \underbrace{\sqrt{2C}}_{\text{Smoothing Cost}} + \underbrace{(\hat{R}_q - \hat{R}_h)}_{\text{Quantization Gap}}, \quad (4)$$

where the complexity C is given by

$$C = \left(\frac{N}{D}\right)b \log 2 + \frac{1}{D} \log \frac{|K|}{\delta},$$

and $\Sigma = \Sigma(C, \Delta, \{A_k\}_{k=1}^n, K)$ (defined in [Theorem 3.2](#)) which can be upper bounded in terms of the empirical loss variance

$$\Sigma \leq 2 \sqrt{\frac{1}{D} \sum_{k=1}^D (R_q(X_k | X_{<k}) - \mathbb{E}_{Y_k \sim p_h}[R_q(Y_k | X_{<k})])^2}.$$

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323

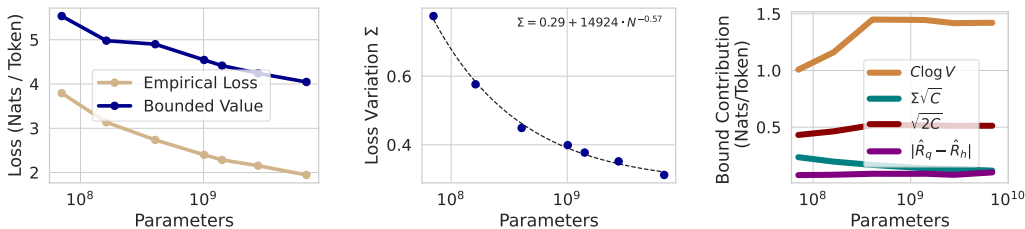


Figure 2: **Left:** A direct comparison of our evaluated generalization bound, and the empirical loss as a function of model scale. As the model is scaled up, our bound improves just like the empirical loss. **Center:** Loss variation Σ entering into the generalization bound. As the loss deviation decreases, so does the largest term in our bound. **Right:** Comparison of the relative scale of the contributions to Theorem 3.4.

To make sense of the bound, let’s consider the various terms. The bounded quantity on the left hand side, R_{sq} , is the expected tokenwise risk of the smoothed and quantized version of hypothesis h . The bound is written in terms of the empirical risk of the original model h , with $\hat{R}_q - \hat{R}_{sq}$ controlled by the **smoothing cost** and **quantization gap**. Typically, the largest contribution to the bound is $C \log V$, e.g. the complexity times the **negative log likelihood of random guessing**. The **loss variation** relates to how spread the empirical loss is and can be seen as a model realizability term. If there existed a 0 loss model in the model class, then this term could be brought to 0. Given the nonzero entropy of natural text, this will not be the case; however, as models improve and approach the irreducible error, so too will the empirical loss variation.

In this setup, the complexity C is just the ratio of parameters to data points, $\frac{N}{D} = G^2$, times the number of bits per parameter used in the quantization b , plus a negligible additional term. The decreased complexity of using fewer bits for b trades off with the quantization gap, and in principle this parameter should be optimized to achieve the best bound. As all terms in the expression can be evaluated empirically, we can determine how much of the empirical observation it can explain and how much remains to be understood.

To get a sense for the scale of the different terms, consider the following typical scenario. For simplicity, suppose $V = 50000$, $\log V \approx 11$, $b = 3$, $G^2 = 1/20$, which yields $C \approx 1/9$. Σ varies with model scale but is of scale $1/10$, and the quantization gap is around $1/10$. Evaluating these terms, $R_{sq} - \hat{R}_h \leq 11/9 + 1/30 + 1.4/3 + 1/10 \approx 1.8$ nats per token, and we see that the random guess and smoothing terms contribute most to the size of the bound. To put this value into perspective, the empirical risk \hat{R}_h itself is around 2 nats per token and the boundary between vacuous and nonvacuous bounds is at $\log V \approx 11$ nats per token.

4 EMPIRICAL EVALUATION

As Theorem 3.4 is fully empirical, we simply need to empirically evaluate the loss variation term Σ along with the quantization gap and we can evaluate the generalization bound. We compute these quantities on the given Pythia checkpoints on the Pile dataset on which they were trained and quantized using GPTQ (Frantar et al., 2023) to $b = 4$ bits per parameter, and we evaluate the bounds with failure probability $\delta = 0.01$. The results are shown in Figure 2 for the Chinchilla compute optimal checkpoints within the Pythia model family. Additional evaluation details are provided below.

We compute Σ with K given by 1000 equally spaced points between $[0, 1]$ excluding the endpoints. We estimate the risk and loss variation on an IID subsample from the collection of token-context pairs in the training dataset of size 10^4 and bound the difference with standard techniques. We note that largest 12B parameter model failed to quantize properly (possibly due to the learning rate drop as it was the only checkpoint taken at the end of training) and so we removed it from the evaluation. As we must pick out the compute optimal checkpoints from training runs that were not designed for this purpose, for the two smallest models (70m and 160m parameters), which the relevant region of

the training curve is more sparsely sampled with checkpoints, the closest models consistently have too small a value of N/D , biasing these two initial data points towards lower values and we urge the reader to keep this in mind while reading the plots.

In [Figure 2](#), we can observe several points. In [Figure 2](#) center, we see that the loss variation decreases with model size in a predictable way, also obeying a power law relationship with a constant offset. As additional compute is spent, the models explain a larger fraction of the variance in the data and the loss variation decreases, but like with the irreducible error there is a minimum value that it is converging towards. Beyond language models, we would expect to see this behavior also, but the predictable improvements with compute give it the simple relationship observed here. In [Figure 2](#) right, we break down the individual contributions to the generalization bound, with the quantization error and loss variance being small and decreasing whereas the $\mathcal{C} \log V$ term and smoothing terms make up the majority and are not decreasing (if the complexity \mathcal{C} is constant). [Figure 2](#) left shows the comparison with the bound value R_{sq} and the empirical risk \hat{R}_h . The fact that the quantization gap at a fixed number of bits is decreasing model size has been observed to an even larger extent in other work ([Chee et al., 2024](#); [Tseng et al., 2024](#)) with more advanced quantization methods. This property suggests that if b were able to be freely optimized in the bound, the complexity \mathcal{C} would actually decrease with model size, and we explore evidence for and consequences of this idea in the following section.

5 COMPRESSIBILITY AND THE SUBLINEAR INFORMATION GROWTH IN LLMs

While not obvious from efficient quantization algorithms like GPTQ ([Frantar et al., 2023](#)), there are good reasons to believe that the model complexity term $L(h)/D$ decreases with model scale on the compute optimal frontier.

5.1 QUANTIZABILITY FROM THE HESSIAN SPECTRUM AND QUIP

So far we have split the compressed size of the model $L(h)$ featured in the complexity term into the number of parameters N times the number of bits per parameter used in the quantization b : $L(h) \leq bN \log 2$. In this splitting increased compressibility of a model shows up in terms of requiring a smaller number of bits b to achieve a given quantization gap $\hat{R}_q - \hat{R}_h$. In [Appendix B](#), we provide a theoretical argument using the QuIP quantization framework ([Chee et al., 2024](#)) for why we should expect that larger models can be more easily quantized. If the hessian around the solution weights is PSD and the spectrum decays sufficiently rapidly, then we should expect the quantization error to decrease with model size. In [Section B.1](#), we investigate the Hessian spectrum empirically finding that it indeed decays sufficiently quickly (though not always PSD). Unfortunately, the version of QuIP needed to construct this argument cannot be run in practice due to the impractically large computational constraints. Empirically it has been observed by some that practical quantization algorithms also reveal that larger models are more quantizable ([Tseng et al., 2024](#)), though the effect is not very pronounced with the GPTQ algorithm we use here.

Alternatively, we present a more abstract information theoretic argument to provide evidence for the fact that $L(h)/D$ decreases with model scale even if we do not have an explicit compression scheme that makes it so.

5.2 INFORMATION ACCUMULATION IN LLMs

Despite the many parameters, at initialization the information content in a neural network (and thus the size a compression scheme can achieve) is extremely small. It suffices to specify the model architecture and a random seed used to initialize the weights, both of which can be specified ahead of time without seeing the data. As the loss decreases, information is transferred incrementally from the dataset to the weights of the model with each additional data point. This information can be quantified abstractly using algorithmic information theory and prequential coding ([Rissanen, 1984](#); [Dawid, 1984](#)). Let X be the training dataset with D tokens and h the LLM with N parameters. Let $K(X)$ be the (prefix) Kolmogorov complexity of X , that is the length of a shortest self delimiting program that produces X when run, and we will consider the description of the LLM architecture as well as code for performing arithmetic encoding and decoding as part of the language in question. From

the symmetry of information property (the analog of Bayes rule), $K(X, h) = K(h) + K(X|h) + c$, where c is a small constant and $K(X|h)$ is the length of a shortest (self delimiting) program which takes as input h and produces X . As described in [Blier & Ollivier \(2018\)](#); [Voita & Titov \(2020\)](#) and more specifically in [Zhang et al. \(2020\)](#) for measuring the information transfer, rearranging $K(h) = K(X, h) - K(X|h) - c$ one can estimate an upper bound on the information stored in a given model using prequential coding.

In a prequential code ([Rissanen, 1984](#); [Dawid, 1984](#)), we consider the hypothesis h_k obtained from training after seeing only $k - 1$ of the data points in X . Starting with the initialization at h_0 , one can use h_{k-1} to encode X_k using arithmetic coding (or any entropy stream code) with $-\log_2 p_{h_{k-1}}(X_k|X_{<k})$ bits and transmit this data. After transmitting X_k , the model is trained on X_k and the next data point is transmitted and so forth until the entire dataset has been encoded and transmitted. On the other end, this data allows reconstructing both the entire dataset and the model if the details of the training algorithm are known. One simply uses h_k to decode the given data point, train one step, and repeat so that the exact sequence of models $[h_0, h_1, h_2, \dots, h_D]$ is repeated. The length of the code for X, h is then only $-\sum_{k=1}^D \log_2 p_{h_{k-1}}(X_k|X_{<k})$, the area under the loss curve.

Though used for classification problems in [Zhang et al. \(2020\)](#), we can readily repurpose the approach for the autoregressive unsupervised learning task. For the coding of $X|h$, one can estimate this using the code length of the final model $-\sum_{k=1}^D \log_2 p_{h_D}(X_k|X_{<k})$. Assembling these two components together, one can estimate an upper bound on $K(h)$. Notably, one can convert a regular code of ℓ bits into a prefix free code of $L = \ell + 2 \log_2 \ell + 1$ bits. Written in terms of the empirical risk and up to additive logarithmic factors $\tilde{O}(1)$,

$$K(h) \log(2) \leq \sum_{k=1}^D [R_{h_{k-1}}(X_k|X_{<k}) - R_h(X_k|X_{<k})] + \tilde{O}(1). \tag{5}$$

If we plug in the Chinchilla scaling law $R(N, D) = E + AN^{-\alpha} + BD^{-\beta}$ as an estimate for the risk along the training trajectory, (and noting that $\frac{1}{D} \sum_{k=1}^D f(k) \rightarrow \frac{1}{D} \int_1^D f(x) dx$ as $D \rightarrow \infty$) we would have

$$K(h) \log(2) \leq \left(\sum_{k=1}^D R(N, k) \right) - DR(N, D) + \tilde{O}(1) = \tilde{O}(D^{1-\beta}). \tag{6}$$

Looking at the right hand side $K(h) = \tilde{O}(D^{1-\beta})$ we have an insightful result for LLMs: the information content in the model grows sublinearly in the size of the training dataset, with a coefficient depending on the scaling law.

As shown in [Figure 3](#) left, using the actual loss curves to evaluate $\sum_{k=1}^D [R_{h_{k-1}}(X_k|X_{<k}) - R_h(X_k|X_{<k})]$ and fitting the results to a power law we get a very good agreement with the predicted $\tilde{O}(D^{1-\beta})$. Notably the empirical fit yields $2 \times 10^4 \cdot N^{0.64} \propto D^{0.64}$, and from the constants estimated in [Besiroglu et al. \(2024\)](#), $\beta = 0.37$, which would yield $D^{1-\alpha} = D^{0.63}$. While the empirical values for the upper bound lie above the straightforward value one gets from quantization and parameter counting bN over the range of Pythia models, the curves predict a crossover point at $\approx 30\text{B}$ parameter models. This would mean that despite the large number of parameters, the information stored in each one decreases with scale.

5.3 IMPLICATIONS FOR GENERALIZATION BOUNDS

If we apply this observation to upper bound the complexity featured in our generalization bounds ([Theorem 3.4](#)) $\mathcal{C} = \frac{L(h) + \log |K|/\delta}{D} = K(h) \log(2)/D + \frac{\log |K|/\delta}{D} = \tilde{O}(D^{-\beta})$, we see that the complexity will actually decrease with the size of the dataset even as the ratio with parameters is held constant. With this scaling we can derive a version of the generalization bound [Theorem 3.4](#) without needing to consider quantization or considering the number of explicit parameters in the model, provided that the Chinchilla scaling law holds.

We evaluate the non asymptotic generalization bound of [Theorem 3.4](#) but using the complexity derived from empirical prequential coding bound in [Figure 3](#) (right) and break it down into the scaling

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

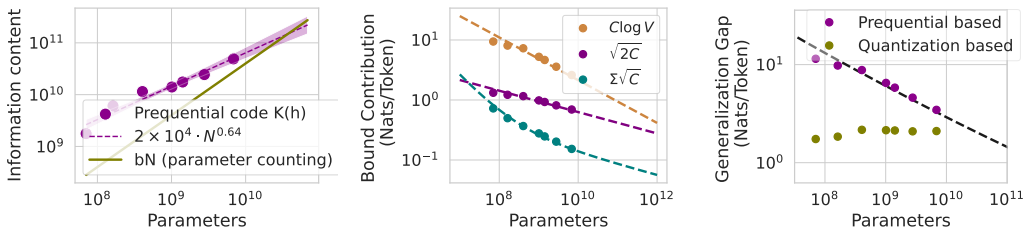


Figure 3: **Left:** Information content contained in the model as upper bounded by $K(h)$ from the information transfer prequential coding approach vs parameter counting and quantization. Fitting a power law to the prequential $K(h)$ yields $2 \times 10^4 \cdot N^{0.64}$. While parameter counting gives a better upper bound over the range of Pythia models, the sublinear scaling of the prequential bound means that it overtakes it eventually, somewhere around 30B sized models. **Center:** The contributions of the various terms to our generalization bounds when using prequential coding complexity, along with their power law fits. **Right:** Comparison of generalization bounds produced by the prequential vs quantization based approaches. While the prequential bounds are worse, they follow a power law and improve substantially with scale.

of the individual terms (center), with the Σ term scaling law extrapolated from the fit in Figure 2. Like before, the $C \log V$ term dominates; however, the $\sqrt{2C}$ smoothing term threatens to overtake it with very large model sizes. We can see that the bounds based on the prequential coding are worse than their quantization counterparts, however the bounds improve with scale and can be extrapolated with scaling laws.

At a high level just considering the asymptotics, the generalization gap of Theorem 3.4 will be dominated by the scaling of the smoothing term $\sqrt{2C}$: $R_s - \hat{R}_h = \tilde{O}(D^{-\beta/2})$. To speculate, it seems likely that with a more sophisticated approach for dealing with the unbounded loss, the $\sqrt{C} = \tilde{O}(D^{-\beta/2})$ could be removed, letting the $\tilde{O}(D^{-\beta})$ shine through. If that were the case, then the tokenwise generalization gap could be hidden within the $D^{-\beta}$ of the original scaling law, and we leave such investigations to future work.

6 ADDITIONAL RELATED WORK

Generalization Bounds. Historically, generalization bounds for neural networks have been limited due to their large parameter count; however, significant progress has been made in explaining the generalization behavior of moderately large machine learning models over the years (Dziugaite & Roy, 2017; Zhou et al., 2018; Arora et al., 2018; Lotfi et al., 2022). PAC-Bayes has proven to be a convenient unifying framework for accommodating many of these techniques (Catoni, 2007). Lotfi et al. (2024a) demonstrate how to construct the first non-vacuous generalization bounds for LLMs by handling the unbounded objective with prediction smoothing and employing extreme compression with subspace LoRA (Hu et al., 2021) training.

While Lotfi et al. (2024a) focus on constructing generalization bounds at the level of documents, Lotfi et al. (2024b) take a different approach by applying Azuma’s inequality to derive a martingale-based bound at the token level, leveraging the much larger number of data points. We adopt this approach here, and improve on the $O(\sqrt{C})$ complexity on the dominating term from Azuma’s inequality to the $O(C)$ of the dominating term of our inequality through the use of the loss variation. Generalization bounds have also been used to constrain the context learning performance of LLMs (Li et al., 2023), and generalization over the hypothesis space of prompts in vision-language models has been explored (Akinwande et al., 2023).

Chugg et al. (2023) develop a broad family of generalization bounds applicable to both IID and martingale settings, generalizing many previous theoretical results. The martingale bounds presented in their work are more similar to ours than most other approaches; however, they do not provide fully empirical bounds suitable for our purposes. In terms of concentration inequalities, the closest work is

486 [Maurer & Pontil \(2009\)](#), where the authors derive a fully empirical Bernstein inequality. However,
487 their result and proof technique do not extend to the non-IID martingale setting.
488

489 **Post-Training Quantization.** For hardware efficiency, there has been significant research into
490 how low the precision of a model can be reduced post-training without substantially degrading its
491 performance ([Hassibi et al., 1993](#); [Hubara et al., 2021](#); [Yao et al., 2022](#); [Dettmers et al., 2022](#)).
492 Empirically, the community has found that 3 or 4 bits provide a reasonable tradeoff between model
493 compression and performance. Further research has also focused on applying quantization at the
494 scale of LLMs, achieving a small number of bits per parameter, with [Ma et al. \(2024\)](#) pushing the
495 limits to 1.58 bits per parameter. With some fine-tuning, even binary networks have shown promise
496 ([Liu et al., 2024](#)).

497 We specifically highlight a few particularly relevant post-training quantization (PTQ) methods. [Frantar](#)
498 [et al. \(2023\)](#) introduce GPTQ (formerly known as OPTQ), which demonstrates *extreme quantization*
499 that scales to billions of parameters without significant degradation while being computationally
500 efficient. This performance is achieved through iterative rounding of the weight columns. In this
501 work, we use GPTQ with 4 bits to post-quantize the models. [Chee et al. \(2024\)](#) propose QuIP, which
502 relies on the insight of *incoherence* in approximate Hessian estimation, ensuring that weights are not
503 too large along a single direction, thereby suppressing outliers, computed successively at each layer.
504 QuIP# further improves upon this incoherence processing ([Tseng et al., 2024](#)). In this work, we adopt
505 the GPTQ quantizer, as it is easier to work with and to adapt to our setting.
506

507 7 DISCUSSION

509 Here we have provided generalization theory to better explain why large language models trained
510 in the compute-optimal regime generalize, with particular attention on how generalization changes
511 with scale. For the term that contributes the most to the generalization bound, we are able to
512 improve the scaling over alternate methods from $\sqrt{\mathcal{C}}$ to \mathcal{C} while still being fully empirical. We
513 explore two approaches for constraining model complexity in the generalization bounds, directly
514 via quantization and parameter counting, and indirectly, via information transfer as quantified by
515 prequential coding. While the quantization approach yields more practical bounds, the information
516 transfer approach yields insights into how generalization scales with model size without requiring an
517 explicit quantization strategy that behaves that way.

518 While we believe that these insights help advance understanding, there are a number of limitations
519 of our approach and many questions that remain unaddressed. As previously mentioned, the $\sqrt{2\mathcal{C}}$
520 smoothing term seems pessimistic and could likely be improved with a different approach. Addition-
521 ally, while the information transfer argument provides evidence *that* the complexity of a model is low
522 based on the training curve, it falls short of explaining *why* the complexity is low. In principle the
523 training curve could look different, leading to a different information transfer rate. Similarly with the
524 Hessian based argument for why larger models are more quantizable as it depends on the empirical
525 spectrum of the Hessian which we have not explained. Furthermore, why does the loss variation term
526 Σ scale in the way that it does?

527 Even more broadly, our generalization bounds constrain only the token-wise generalization gap.
528 While it is intuitive that generalizing well on next token prediction over the training contexts should
529 imply generalization on the full sequences, we are not aware of work that does so and that gap
530 remains to be understood. Similarly, constraining generalization on the NLL objective over data
531 drawn from the natural distribution may be less pertinent. Instead, it may be more relevant to constrain
532 the gap between the quality metrics of model generations and natural data. Farther removed, there
533 is the question of why the training loss scales in the way that it does, and how does that relate to
534 approximation theory and the architecture of the model? Though many questions remain, we hope
535 that the techniques here can yield generalizable insights for tackling this broader set of problems.
536

537 REFERENCES

538 Victor Akinwande, Yiding Jiang, Dylan Sam, and J Zico Kolter. Understanding prompt engineering
539 may not require rethinking generalization. *arXiv preprint arXiv:2310.03957*, 2023.

- 540 Sanjeev Arora, Rong Ge, Behnam Neyshabur, and Yi Zhang. Stronger generalization bounds for deep
541 nets via a compression approach. In *International Conference on Machine Learning*, pp. 254–263.
542 PMLR, 2018.
- 543 Yasaman Bahri, Ethan Dyer, Jared Kaplan, Jaehoon Lee, and Utkarsh Sharma. Explaining neural
544 scaling laws. *Proceedings of the National Academy of Sciences of the United States of America*,
545 121, 2021.
- 546 Tamay Besiroglu, Ege Erdil, Matthew Barnett, and Josh You. Chinchilla scaling: A replication
547 attempt. *arXiv preprint arXiv:2404.10102*, 2024.
- 548 Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O’Brien, Eric
549 Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, et al.
550 Pythia: A suite for analyzing large language models across training and scaling. In *International
551 Conference on Machine Learning*, pp. 2397–2430. PMLR, 2023.
- 552 Léonard Blier and Yann Ollivier. The description length of deep learning models. *Advances in Neural
553 Information Processing Systems*, 31, 2018.
- 554 Léon Bottou and Olivier Bousquet. The tradeoffs of large scale learning. *Advances in neural
555 information processing systems*, 20, 2007.
- 556 Gavin Brown and Riccardo Ali. Bias/variance is not the same as approximation/estimation. *Transac-
557 tions on Machine Learning Research*, 2024.
- 558 Tom B Brown. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.
- 559 Joël Bun, Jean-Philippe Bouchaud, and Marc Potters. Cleaning large correlation matrices: tools from
560 random matrix theory. *Physics Reports*, 666:1–109, 2017.
- 561 Olivier Catoni. Pac-bayesian supervised classification: the thermodynamics of statistical learning.
562 *arXiv preprint arXiv:0712.0248*, 2007.
- 563 Jerry Chee, Yaohui Cai, Volodymyr Kuleshov, and Christopher M De Sa. Quip: 2-bit quantization of
564 large language models with guarantees. *Advances in Neural Information Processing Systems*, 36,
565 2024.
- 566 Ben Chugg, Hongjian Wang, and Aaditya Ramdas. A unified recipe for deriving (time-uniform)
567 pac-bayes bounds. *Journal of Machine Learning Research*, 24(372):1–61, 2023.
- 568 G. Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of Control,
569 Signals and Systems*, 2(4):303–314, 1989. doi: 10.1007/BF02551274. URL [https://doi.
570 org/10.1007/BF02551274](https://doi.org/10.1007/BF02551274).
- 571 A Philip Dawid. Present position and potential developments: Some personal views statistical theory
572 the prequential approach. *Journal of the Royal Statistical Society: Series A (General)*, 147(2):
573 278–290, 1984.
- 574 Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. Llm.int8(): 8-bit matrix
575 multiplication for transformers at scale. *ArXiv*, abs/2208.07339, 2022.
- 576 Kun Dong, David Eriksson, Hannes Nickisch, David Bindel, and Andrew G Wilson. Scalable log
577 determinants for gaussian process kernel learning. *Advances in Neural Information Processing
578 Systems*, 30, 2017.
- 579 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
580 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn,
581 Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston
582 Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Rozière, Bethany Biron,
583 Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris
584 McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Cantón
585 Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David
586 Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes,
587 Egor Lakomkin, Ehab A. AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip

594 Radenovic, Frank Zhang, Gabriele Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme
595 Nail, Grégoire Mialon, Guanglong Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar,
596 Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel M. Kloumann, Ishan Misra, Ivan
597 Evtimov, Jade Copet, Jaewon Lee, Jan Laurens Geffert, Jana Vranes, Jason Park, Jay Mahadeokar,
598 Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng
599 Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park,
600 Joseph Rocca, Joshua Johnstun, Joshua Saxe, Ju-Qing Jia, Kalyan Vasuden Alwala, K. Upasani,
601 Kate Plawiak, Keqian Li, Ken-591 neth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer,
602 Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yearly, Laurens van der Maaten,
603 Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas
604 Blecher, Lukas Landzaat, Luke de Oliveira, Madeline C. Muzzi, Mahesh Babu Pasupuleti, Mannat
605 Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melissa
606 Hall Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal,
607 Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri S. Chatterji, Olivier Duchenne,
608 Onur cCelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasić, Peter Weng, Prajjwal
609 Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao
610 Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert
611 Stojnic, Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan
612 Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa,
613 Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang,
614 Sharath Chandra Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon
615 Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan,
616 Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas
617 Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami,
618 Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti,
619 Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet,
620 Xiaodong Wang, Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuwei Wang, Yaelle Goldschlag,
621 Yashesh Gaur, Yasmine Babaei, Yiqian Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao,
622 Zacharie Delpierre Coudert, Zhengxu Yan, Zhengxing Chen, Zoe Papanikos, Aaditya K. Singh,
623 Aaron Grattafiori, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adi Gangidi, Adolfo Victoria, Ahuva
624 Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie
625 Feinstein, Amanda Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew
626 Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie
627 Franco, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf
628 Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Ben Leonhardi, Bernie Huang, Beth Loyd,
629 Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti,
630 Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton,
631 Catalina Mejia, Changan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu,
632 Chris Cai, Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer,
633 Shang-Wen Li, Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi
634 Parikh, Diana Liskovich, Didem Foss, Dingkan Wang, Duc Le, Dustin Holland, Edward Dowling,
635 Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman,
636 Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel,
637 Francesco Caggioni, Francisco Guzm'an, Frank J. Kanayet, Frank Seide, Gabriela Medina Florez,
638 Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman,
639 Grigory G. Sizov, Guangyi Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah
640 Wang, Han Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman,
641 Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James
642 Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen,
643 Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon
644 Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kaixing(Kai)
645 Wu, U KamHou, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun Zand, Kathy
646 Matosich, Kaushik Veeraraghavan, Kelly Michelen, Keqian Li, Kun Huang, Kunal Chawla, Kushal
647 Lakhotia, Kyle Huang, Lailin Chen, Lakshya Garg, A Lavender, Leandro Silva, Lee Bell, Lei
Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav
Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson, Matthew Lennie,
Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Michael L.
Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike
Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari,

- 648 Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan
649 Singhal, Nick Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong, Ning Zhang,
650 Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth
651 Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollár, Polina
652 Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez,
653 Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li, Rebekkah Hogan, Robin
654 Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott, Sai Jayesh Bondu, Samyak
655 Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Verma, Seiji
656 Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy
657 Zha, Shiva Shankar, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith
658 Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad,
659 Sumit Gupta, Sung-Bae Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman,
660 Tal Remez, Tamar Glaser, Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun
661 Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria
662 Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Andrei Poenaru,
663 Vlad T. Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz,
664 Will Constable, Xia Tang, Xiaofang Wang, Xiaojian Wu, Xiaolan Wang, Xide Xia, Xilun Wu,
665 Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi,
666 Youngjin Nam, Yu Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, Zachary DeVito, Zef
667 Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. The llama 3 herd of models. *ArXiv*,
668 abs/2407.21783, 2024.
- 668 Gintare Karolina Dziugaite and Daniel M Roy. Computing nonvacuous generalization bounds for
669 deep (stochastic) neural networks with many more parameters than training data. *arXiv preprint*
670 *arXiv:1703.11008*, 2017.
- 671 Jack Fitzsimons, Diego Granzio, Kurt Cutajar, Michael Osborne, Maurizio Filippone, and Stephen
672 Roberts. Entropic trace estimates for log determinants. In *Machine Learning and Knowledge Dis-*
673 *covery in Databases: European Conference, ECML PKDD 2017, Skopje, Macedonia, September*
674 *18–22, 2017, Proceedings, Part I 10*, pp. 323–338. Springer, 2017.
- 675 Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and Dan Alistarh. Gptq: Accurate post-training
676 quantization for generative pre-trained transformers, 2023.
- 677 David A. Freedman. On tail probabilities for martingales. 1975.
- 678 Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang,
679 Horace He, Anish Thite, Noa Nabeshima, et al. The pile: An 800gb dataset of diverse text for
680 language modeling. *arXiv preprint arXiv:2101.00027*, 2020.
- 681 Behrooz Ghorbani, Shankar Krishnan, and Ying Xiao. An investigation into neural net optimization
682 via hessian eigenvalue density. In *International Conference on Machine Learning*, pp. 2232–2241.
683 PMLR, 2019.
- 684 Diego Granzio, Xingchen Wan, Timur Garipov, Dmitry Vetrov, and Stephen Roberts. Mlrg deep
685 curvature: An open-source package to analyse and visualise neural network curvature and loss
686 surface. *stat*, 2018.
- 687 Diego Granzio, Timur Garipov, Stefan Zohren, Dmitry Vetrov, Stephen Roberts, and Andrew Gordon
688 Wilson. The deep learning limit: are negative neural network eigenvalues just noise? In *ICML*
689 *2019 workshop on theoretical physics for deep learning*, 2019.
- 690 Diego Granzio, Stefan Zohren, and Stephen Roberts. Learning rates as a function of batch size: A
691 random matrix theory approach to neural network training. *Journal of Machine Learning Research*,
692 23(173):1–65, 2022.
- 693 Nate Gruver, Marc Finzi, Shikai Qiu, and Andrew G Wilson. Large language models
694 are zero-shot time series forecasters. In A. Oh, T. Naumann, A. Globerson,
695 K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Information*
696 *Processing Systems*, volume 36, pp. 19622–19635. Curran Associates, Inc.,
697 2023. URL [https://proceedings.neurips.cc/paper_files/paper/2023/
698 file/3eb7ca52e8207697361b2c0fb3926511-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/3eb7ca52e8207697361b2c0fb3926511-Paper-Conference.pdf).

- 702 Babak Hassibi, David G. Stork, and Gregory J. Wolff. Optimal brain surgeon and general network
703 pruning. *IEEE International Conference on Neural Networks*, pp. 293–299 vol.1, 1993.
704
- 705 Tom Henighan, Jared Kaplan, Mor Katz, Mark Chen, Christopher Hesse, Jacob Jackson, Heewoo
706 Jun, Tom B Brown, Prafulla Dhariwal, Scott Gray, et al. Scaling laws for autoregressive generative
707 modeling. *arXiv preprint arXiv:2010.14701*, 2020.
- 708 Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza
709 Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom
710 Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy,
711 Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and L. Sifre. Training
712 compute-optimal large language models. *ArXiv*, abs/2203.15556, 2022.
- 713 Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
714 and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint*
715 *arXiv:2106.09685*, 2021.
716
- 717 Itay Hubara, Yury Nahshan, Yair Hanani, Ron Banner, and Daniel Soudry. Accurate post training
718 quantization with small calibration sets. In *International Conference on Machine Learning*, 2021.
719
- 720 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child,
721 Scott Gray, Alec Radford, Jeff Wu, and Dario Amodei. Scaling laws for neural language models.
722 *ArXiv*, abs/2001.08361, 2020.
- 723 Yingcong Li, Muhammed Emrullah Ildiz, Dimitris Papailiopoulos, and Samet Oymak. Transformers
724 as algorithms: Generalization and stability in in-context learning. In *International Conference on*
725 *Machine Learning*, pp. 19565–19594. PMLR, 2023.
- 726 James Liu, Guangxuan Xiao, Kai Li, Jason D. Lee, Song Han, Tri Dao, and Tianle Cai. Bitdelta:
727 Your fine-tune may only be worth one bit, 2024.
728
- 729 Sanae Lotfi, Marc Finzi, Sanyam Kapoor, Andres Potapczynski, Micah Goldblum, and Andrew G
730 Wilson. Pac-bayes compression bounds so tight that they can explain generalization. *Advances in*
731 *Neural Information Processing Systems*, 35:31459–31473, 2022.
- 732 Sanae Lotfi, Marc Finzi, Yilun Kuang, Tim G. J. Rudner, Micah Goldblum, and Andrew Gordon
733 Wilson. Non-vacuous generalization bounds for large language models, 2024a.
734
- 735 Sanae Lotfi, Yilun Kuang, Brandon Amos, Micah Goldblum, Marc Finzi, and Andrew Gordon Wilson.
736 Unlocking tokens as data points for generalization bounds on larger language models. *ArXiv*,
737 abs/2407.18158, 2024b.
- 738 Sanae Lotfi, Yilun Kuang, Brandon Amos, Micah Goldblum, Marc Finzi, and Andrew Gordon
739 Wilson. Unlocking tokens as data points for generalization bounds on larger language models.
740 *arXiv preprint arXiv:2407.18158*, 2024c.
741
- 742 Shuming Ma, Hongyu Wang, Lingxiao Ma, Lei Wang, Wenhui Wang, Shaohan Huang, Li Dong,
743 Ruiping Wang, Jilong Xue, and Furu Wei. The era of 1-bit llms: All large language models are in
744 1.58 bits, 2024.
- 745 Andreas Maurer and Massimiliano Pontil. Empirical bernstein bounds and sample variance penaliza-
746 tion. *arXiv preprint arXiv:0907.3740*, 2009.
747
- 748 Gérard Meurant and Zdeněk Strakoš. The lanczos and conjugate gradient algorithms in finite precision
749 arithmetic. *Acta Numerica*, 15:471–542, 2006.
- 750 Niklas Muennighoff, Alexander Rush, Boaz Barak, Teven Le Scao, Nouamane Tazi, Aleksandra
751 Piktus, Sampo Pyysalo, Thomas Wolf, and Colin A Raffel. Scaling data-constrained language
752 models. *Advances in Neural Information Processing Systems*, 36, 2024.
753
- 754 Markus Nagel, Rana Ali Amjad, Mart Van Baalen, Christos Louizos, and Tijmen Blankevoort. Up or
755 down? adaptive rounding for post-training quantization. In *International Conference on Machine*
Learning, pp. 7197–7206. PMLR, 2020.

- 756 OpenAI. Gpt-4 technical report. 2023.
757
- 758 Vardan Papyan. Measurements of three-level hierarchical structure in the outliers in the spectrum of
759 deepnet hessians. *arXiv preprint arXiv:1901.08244*, 2019.
- 760 Roma Patel and Ellie Pavlick. Mapping language models to grounded conceptual spaces. In
761 *International Conference on Learning Representations*, 2022.
762
- 763 Barak A Pearlmutter. Fast exact multiplication by the hessian. *Neural computation*, 6(1):147–160,
764 1994.
- 765 Andres Potapczynski, Marc Finzi, Geoff Pleiss, and Andrew Gordon Wilson. CoLA: Exploiting
766 Compositional Structure for Automatic and Efficient Numerical Linear Algebra. *arXiv preprint*
767 *arXiv:2309.03060*, 2023.
768
- 769 Jorma Rissanen. Universal coding, information, prediction, and estimation. *IEEE Transactions on*
770 *Information theory*, 30(4):629–636, 1984.
- 771 Farbod Roosta-Khorasani and Uri Ascher. Improved bounds on sample size for implicit matrix trace
772 estimators. *Foundations of Computational Mathematics*, 15(5):1187–1212, 2015.
773
- 774 Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning: From theory to*
775 *algorithms*. Cambridge university press, 2014.
- 776 Trieu H. Trinh, Yuhuai Wu, Quoc V. Le, He He, and Thang Luong. Solving olympiad geometry
777 without human demonstrations. *Nature*, 625:476 – 482, 2024.
778
- 779 Albert Tseng, Jerry Chee, Qingyao Sun, Volodymyr Kuleshov, and Christopher De Sa. Quip: Even
780 better llm quantization with hadamard incoherence and lattice codebooks, 2024.
- 781 Shashanka Ubaru, Jie Chen, and Yousef Saad. Fast estimation of $\text{tr}(f(a))$ via stochastic lanczos
782 quadrature. *SIAM Journal on Matrix Analysis and Applications*, 38(4):1075–1099, 2017.
783
- 784 Jean Ville. *Étude critique de la notion de collectif*. 1939. URL [http://eudml.org/doc/](http://eudml.org/doc/192893)
785 [192893](http://eudml.org/doc/192893).
- 786 Elena Voita and Ivan Titov. Information-theoretic probing with minimum description length. *arXiv*
787 *preprint arXiv:2003.12298*, 2020.
788
- 789 Zhewei Yao, Reza Yazdani Aminabadi, Minjia Zhang, Xiaoxia Wu, Conglong Li, and Yuxiong He.
790 Zeroquant: Efficient and affordable post-training quantization for large-scale transformers. *ArXiv*,
791 [abs/2206.01861](https://arxiv.org/abs/2206.01861), 2022.
- 792 Xiao Zhang, Xingjian Li, Dejing Dou, and Ji Wu. Measuring information transfer in neural networks.
793 *arXiv preprint arXiv:2009.07624*, 2020.
794
- 795 Wenda Zhou, Victor Veitch, Morgane Austern, Ryan P Adams, and Peter Orbanz. Non-vacuous
796 generalization bounds at the imagenet scale: a pac-bayesian compression approach. *arXiv preprint*
797 *arXiv:1804.05862*, 2018.
798
799
800
801
802
803
804
805
806
807
808
809

810 A PROOFS

811 A.1 A FULLY EMPIRICAL MARTINGALE FREEDMAN CONCENTRATION INEQUALITY.

812 In this section, unless otherwise specified, \log is used to denote the natural logarithm. We start with
813 three technical lemmas.

814 **Lemma A.1.** *Consider the function $v(a) = a - \log(1 + a)$. Let $\mu \in \mathbb{R}$. For any random variable Z
815 with $\mathbb{E}[Z] = 0$ that satisfies $Z - \mu > -1$, we have*

$$816 \mathbb{E}[\exp(Z - v(Z - \mu))] = (1 - \mu)e^\mu \leq 1.$$

817 *Proof.* We see that

$$818 \mathbb{E}[\exp(Z - v(Z - \mu))] = e^\mu \mathbb{E}[\exp(Z - \mu - v(Z - \mu))] = e^\mu \mathbb{E}[1 + (Z - \mu)].$$

819 As $\mathbb{E}[Z] = 0$, we know $\mathbb{E}[\exp(Z - v(Z - \mu))] = (1 - \mu)e^\mu$. Additionally as $1 - \mu \leq e^{-\mu}$, we have
820 $\mathbb{E}[\exp(Z - v(Z - \mu))] \leq 1$, as desired. \square

821 **Lemma A.2.** *Consider the function $v(a) = a - \log(1 + a)$. For all $a \in (-1, \infty)$, we have*

$$822 v(a) \leq a^2/(1 + a)$$

823 *Proof.* Let $f(a) = a^2/(1 + a) - v(a)$. By direct calculation, we see that $f'(a) = a/(1 + a^2)$,
824 which is a strictly negative function passing through 0 at $a = 0$. As $f(0) = 0$, we have $f(a) \geq 0$
825 for all $a \geq 0$. Note that $\lim_{a \rightarrow -1^-} f(a) = +\infty$, so $f(a)$ must be positive on $(-1, 0)$. The claim
826 follows. \square

827 For the following lemma, consider on the filtered probability space $(\Omega, \mathcal{F}, \{\mathcal{F}_k\}_{k \in \mathbb{N}}, \mathbb{P})$, and we
828 consider expectations with respect to \mathbb{P} .

829 **Theorem A.3.** *Let X_1, \dots, X_n be a sequence of \mathcal{F}_k -measurable random variables. Let Y_1, \dots, Y_n be
830 any other sequence of \mathcal{F}_{k-1} -measurable random variables such that the difference is bounded above:
831 $A_k = Y_k - X_k > -\Delta$ for some $\Delta \geq 0$. Define $C = \frac{1}{n} \log \frac{1}{\epsilon}$, and $B = \frac{1}{n} \sum_k (\mathbb{E}[X_k | \mathcal{F}_{k-1}] - X_k)$.
832 For any $0 < t < 1/\Delta$ and simultaneously for all n , we have*

$$833 \mathbb{P}[tB \leq C + \frac{1}{n} \sum_{k=1}^n v(tA_k)] \geq 1 - \epsilon. \quad (7)$$

834 *Proof.* Let $v(a) = a - \log(1 + a)$. Define the random variable

$$835 M_k = \exp(t(\mathbb{E}[X_k | \mathcal{F}_{k-1}] - X_k) - v(t(Y_k - X_k))).$$

836 Consider $Z = t(\mathbb{E}[X_k | \mathcal{F}_{k-1}] - X_k)$ and $\mu = t(\mathbb{E}[X_k | \mathcal{F}_{k-1}] - Y_k)$. By construction, we have
837 that $\mathbb{E}[Z | \mathcal{F}_{k-1}] = 0$ and $Z - \mu = t(Y_k - X_k) > -t\Delta > -1$. Thus, applying Lemma A.1, we have

$$838 \mathbb{E}[M_k | \mathcal{F}_{k-1}] \leq 1.$$

839 Therefore $U_n = \prod_{k=1}^n M_k$ is a supermartingale. By Ville's inequality (Ville, 1939), we have

$$840 \sup_n U_n \leq \frac{\mathbb{E}[U_0]}{\epsilon} \leq 1/\epsilon$$

841 with probability at least $1 - \epsilon$. Using the definition of U and taking the log of both sides, we have for
842 all n ,

$$843 t \sum_{k=1}^n (\mathbb{E}[X_k | \mathcal{F}_{k-1}] - X_k) - \sum_{k=1}^n v(t(Y_k - X_k)) \leq \log \frac{1}{\epsilon}. \quad (8)$$

844 Defining $B = \frac{1}{n} \sum_k (\mathbb{E}[X_k | \mathcal{F}_{k-1}] - X_k)$, $C = \frac{1}{n} \log \frac{1}{\epsilon}$, $A_k = Y_k - X_k$, rearrange Equation 8 to
845 obtain

$$846 tB \leq C + \frac{1}{n} \sum_{k=1}^n v(tA_k). \quad (9)$$

847 \square

Corollary A.4. Let X_1, \dots, X_n be a sequence of \mathcal{F}_k -measurable random variables. Let Y_1, \dots, Y_n be any other sequence of \mathcal{F}_{k-1} -measurable random variables such that the difference is bounded below: $A_k = (Y_k - X_k)/\Delta > -1$ for some $\Delta \geq 0$. Let K be a finite subset of $(0, 1)$. Then, with probability at least $1 - \delta$,

$$\frac{1}{n} \sum_{k=1}^n (\mathbb{E}[X_k | \mathcal{F}_{k-1}] - X_k) \leq \Delta C + \Sigma \sqrt{C}, \quad (10)$$

where $C := \frac{1}{n} \log |K|/\delta$, and

$$\Sigma(C, \Delta, \{X_k - Y_k\}_{k=1}^n, K) := \min_{s \in K} \Delta \sqrt{C} (1 - s)/s + \frac{\Delta}{\sqrt{C}} \frac{1}{n} \sum_{k=1}^n v(sA_k)/s$$

Proof. Let $s = t\Delta$. Apply a union bound to Theorem A.3 over the different values of s in K , taking the one that minimizes the bound. Rearrange and isolate terms. \square

Theorem A.5. Let X_1, \dots, X_n be a sequence of \mathcal{F}_{i-1} -measurable random variables. Let Y_0, \dots, Y_{n-1} be any other sequence of \mathcal{F}_{i-1} measurable sequence of random variables such that the difference is bounded above: $X_k - Y_k \leq \Delta$ for some $\Delta \geq 0$. Define $V = \frac{1}{n} \sum_k (X_k - Y_k)^2$ and let $\delta \in (0, 1)$. Then, with probability at least $1 - \delta$, we have

$$\frac{1}{n} \sum_k (\mathbb{E}[X_k | \mathcal{F}_{k-1}] - X_k) \leq \Delta C + 2\sqrt{VC}, \quad (11)$$

where $C \leq n^{-1}(\log 1/\delta + 4 \log \log n/\delta + 6)$.

Proof. Starting from Theorem A.3, we apply Theorem A.2 of $v(a) \leq a^2(1+a)$. For our convenience, here we will define $A_k = Y_k - X_k$.

$$\begin{aligned} tB &\leq C + \frac{1}{n} \sum_k \frac{t^2 A_k^2}{1 + tA_k} \\ &\leq C + \frac{1}{n} \sum_k \frac{t^2 A_k^2}{1 - t\Delta}, \end{aligned} \quad (12)$$

where the second inequality follows from the assumption that $A_k \geq -\Delta$.

Finally, by defining a variance term $V = \frac{1}{n} \sum_k A_k^2 = \frac{1}{n} \sum_k (X_k - Y_k)^2$ and rearranging Equation 12, we see

$$0 \leq t^2(V + B\Delta) - (B + \Delta C)t + C, \quad (13)$$

which we recall holds with probability at least $1 - \epsilon$.

Inequality Sketch: This inequality is very close to what we need. The approach would be to optimize over t and then read off the constraint on B . The minimizer of the quadratic is at $t^* = \frac{B + \Delta C}{2(V + \Delta B)}$. Plugging in this value one would arrive at,

$$(1/4) \frac{(B + \Delta C)^2}{(V + \Delta B)} - (1/2) \frac{(B + \Delta C)^2}{(V + \Delta B)} + C \geq 0$$

Rearranging,

$$\begin{aligned} \frac{1}{4}(B + \Delta C)^2 - B\Delta C &\leq VC \\ \frac{1}{4}(B - \Delta C)^2 &\leq VC \end{aligned}$$

and finally,

$$B \leq \Delta C + 2\sqrt{VC}.$$

At a high level, this determines the overall form of Theorem A.5, however some technical complications arise from the fact that t must be deterministic and chosen ahead of time (it must not depend on the random variables B and C). Instead we will consider optimizing t over a discrete set of

possibilities (not depending on B or C), and consider a union bound over the different possibilities.
Full derivation: Consider the quadratic, Equation 13. Its minimizer is given by

$$t^* = \frac{B + \Delta C}{2(V + \Delta B)}.$$

Consider two cases: $t^* \geq \frac{1}{\Delta}$ and $t^* < \frac{1}{\Delta}$. If $t^* \geq \frac{1}{\Delta}$, then rearranging and solving for B , we see

$$B \leq \Delta C - 2V/\Delta,$$

which is strictly less than the value $\Delta C + 2\sqrt{VC}$, and we are done.

Therefore it suffices to consider the case $t^* < 1/\Delta$, where we can apply Lemma A.3. Note that this result applies for a single t , so it cannot be directly applied to t^* . Instead, we will turn to quantization and apply a union bound. Note that if $B > 0$, using that $V \leq \Delta^2$, we have $t^* \geq \frac{\Delta C}{2\Delta^2} = \frac{C}{2\Delta}$. Therefore we only need to consider the range: $t^* \in (\frac{C}{2\Delta}, \frac{1}{\Delta}) =: T$.

Drawing inspiration from floating point numbers, consider a discrete set Q defined as

$$Q = \left\{ \frac{1}{\Delta} 2^{-b} \left(1 + \frac{k}{K} \right) \mid k = 0, 1, \dots, K-1, b \in \mathbb{N}^+ \right\}$$

for some $K \in \mathbb{N}$. Let

$$q(a) = \arg \min_{q \in Q} |q - a|.$$

From this we can determine that the quantization error is bounded by

$$\sup_{a \in T} \frac{|q(a) - a|}{a} \leq \frac{1}{K}.$$

Define a prior over the values of Q :

$$P(q_{k,b}) = P(k)P(b) = \frac{1}{K} \frac{1}{Z(b+2)(\log_2(b+2))^2}.$$

By direct calculation, we see that $1 = \sum_{k,b} P(k)P(b) = (\sum_{b=0}^{\infty} \frac{1}{(b+2)(\log_2(b+2))^2})/Z \leq 1/Z$, therefore $Z \leq 1$.

Now we apply a union bound for Lemma A.3 over values of $t \in Q$. For each $t \in Q$, we set $\epsilon(t) = \delta P(t)$. We have

$$\begin{aligned} & \mathbb{P} [\forall t \in Q : t^2(V + B\Delta) - (B + \Delta C_{\epsilon(t)})t + C_{\epsilon(t)} < 0] \\ & \leq \sum_{t \in Q} \mathbb{P} [t^2(V + B\Delta) - (B + \Delta C_{\epsilon(t)})t + C_{\epsilon(t)} < 0] \\ & \leq \sum_{t \in Q} \epsilon(t) = \delta \sum_{t \in Q} P(t) = \delta. \end{aligned}$$

Therefore, uniformly for all $t \in Q$,

$$t^2(V + B\Delta) - (B + \Delta C_{\epsilon(t)})t + C_{\epsilon(t)} \geq 0 \quad (14)$$

with probability at least $1 - \delta$.

Now we plug in $t = q(t^*)$. Note that $0 \leq b \leq \log_2 \frac{2}{C}$, so we have

$$\begin{aligned} \log \frac{1}{\epsilon(t)} &= \log \frac{1}{\delta P(t)} \leq \log \frac{K}{\delta} + \log(3 + \log_2 1/C) + 2 \log \log_2(3 + \log_2 1/C) \\ &\leq \log \frac{K}{\delta} + 2 + 2 \log \log 1/C \leq \log \frac{K}{\delta} + 2 + 2 \log \log n \end{aligned} \quad (15)$$

Plugging in this quantized value of t to Equation 14 and using the quantized error bound, Equation 15, we have

$$\begin{aligned} \frac{(B + \Delta C)^2}{4(V + \Delta B)} &\leq C + \frac{3}{K} \left(1 + \frac{1}{K} \right) \frac{1}{4} \frac{(B + \Delta C)^2}{(V + \Delta B)} \\ &\leq C(1 + 4/K), \end{aligned} \quad (16)$$

where in the second line we chose a $K \geq 6$ so that we have $(1 - \frac{3}{K}(1 + \frac{1}{K}))^{-1} \leq 1 + \frac{4}{K}$. Solving Equation 16 for B , we have the inequality

$$\begin{aligned} B &\leq \Delta C(1 + 8/K) + \sqrt{\Delta^2 C^2 (8/K)^2 + 4CV(1 + 4/K)} \\ &\leq \Delta C(1 + 16/K) + 2\sqrt{VC(1 + 16/K)}, \end{aligned} \quad (17)$$

where the second line follows from the fact that $\sqrt{x+y} \leq \sqrt{x} + \sqrt{y}$. Define $\mathcal{C} = C(1 + 16/K)$. Choosing $K = \lceil 16 \log 1/\delta \rceil$ (which is > 6), we have

$$nC \leq \log 1/\delta + 1 + \lceil \log(\lceil 16 \log 1/\delta \rceil) + 2 + 2 \log \log n \rceil (1 + 1/\log(1/\delta)) \quad (18)$$

Applying some simplifications to Equation 17 and Equation 18, we obtain

$$B \leq \Delta C + 2\sqrt{VC}.$$

with

$$C \leq \frac{1}{n} (\log 1/\delta + 4 \log \log n / \delta + 6).$$

□

A.2 GENERALIZATION BOUND

Converting to Prefix Free Codes Let $\ell(h)$ be the length of a prefix free code and $L(h)$ be the length of a non prefix free code. We wish to make a normalized prior over the hypotheses. We would like to use approximately $2^{-L(h)}$ where $L(h)$ is the codelength of hypothesis h , but this sum would diverge. Instead consider $\ell(L) = L + 2 \log_2(L) + 1$. Computing the sum, $\sum_h 2^{-\ell(h)} = \sum_{L=1}^{\infty} 2^{L-\ell(L)} = \sum_{L=1}^{\infty} \frac{1}{2L^2} = \pi^2/12 < 1$. Thus with the prior $P(h) = \frac{12}{\pi^2} 2^{-\ell(h)}$ can be used for any countable hypothesis class, placing higher mass on elements with shorter descriptions, which is closely related to Kolmogorov complexity as explored in Lotfi et al. (2022).

If we know the length of the object ahead of time, then we are free to use a regular code in place of a prefix free code. For a fixed number of parameters N and bits per parameter b , we know the length of the code, and thus we can use $\ell(h) \leq bN$.

Applying Theorem 3.1 to the sequence $R_h(X_k | X_{<k})$ with $\delta(h) = \epsilon P(h)$ for each hypothesis individually, the probability that the bound is violated for an arbitrary hypothesis constrained with a union bound $\sum_h \epsilon P(h) = \epsilon$, and therefore Theorem 3.2 holds with probability at least $1 - \epsilon$ (replacing δ with ϵ in its expression). The $\log 1/\delta$ in Theorem A.5 becomes $\log 1/\delta + \log 1/P(h) \leq \log 1/\delta + \ell(h) \log 2$.

A.3 PREDICTION SMOOTHING

Theorem A.6. *For the categorical negative log likelihood objective $\hat{R}_h = -\frac{1}{n} \sum_{k=1}^n \log p_h(X_k | X_{<k})$ on V classes and $C \in \mathbb{R}^+$, there exists a prediction smoothed model $p_s(\cdot) = (1 - \alpha)p_h(\cdot) + \alpha/V$ with worst case loss $\sup_{X_k, X_{<k}} -\log p_s(X_k | X_{<k}) \leq \Delta_s$ that satisfies*

$$\hat{R}_s + C\Delta_s \leq \hat{R}_h + C \log V + \sqrt{2C}, \quad (19)$$

for some value $\alpha(C, V) \in (0, 1)$ (approximately $C/(1+C)$).

Proof. We have

$$\begin{aligned} -\log p_s &= -\log((1 - \alpha)p_h + \alpha/V) \\ &\leq -\log p_h - \log(1 - \alpha + \alpha/V). \end{aligned}$$

Noting that, $-\log p_s(X_k | X_{<k}) \leq \Delta_s = \log(V/\alpha)$, so adding $C\Delta$ to both sides yields

$$\hat{R}_s + C\Delta_s \leq \hat{R}_h - \log(1 - \alpha + \alpha/V) + C \log(V/\alpha), \quad (20)$$

where the right-hand side is minimized at

$$\alpha = \frac{VC}{(V-1)(1+C)}.$$

Note that α is a *deterministic quantity* that we can compute ahead of time based on the model we are bounding. Therefore we need not pay additional bits for a union bound over values of α . Substituting α into Equation 20, we have

$$\begin{aligned}\hat{R}_s - \hat{R}_h + C\Delta_s &\leq \log(1+C) + C \log \frac{(V-1)(1+C)}{C} \\ &\leq (1+C) \log(1+C) + C \log(V/C) \\ &\leq C \log V + \sqrt{2C},\end{aligned}$$

where the last line follows from $(1+x) \log(1+x) - x \log x \leq \sqrt{2x}$ for $x > 0$. The claim follows. \square

B QUANTIZABILITY FROM THE HESSIAN

We have shown that the quantization gap tends to be quite small in practice, but why is this the case? A more complete explanation of why LLMs generalize would need to explain why they are readily quantizable, not just why they should achieve a small generalization gap if they are quantizable. In this section we attempt to shed light on this question on the existence of a quantized solution for a model which achieves low quantization error at a small bitrate, regardless of how computationally efficient it is to actually produce said quantized model.

As a starting point in the analysis of many quantization schemes (Nagel et al., 2020), consider the Lagrange remainder form of the quadratic Taylor expansion of the loss around a given solution of the weights θ , with $\hat{\theta}$ being our desired quantization,

$$L(\hat{\theta}) = L(\theta) + g^\top (\hat{\theta} - \theta) + (\hat{\theta} - \theta)^\top H (\hat{\theta} - \theta)$$

holding with equality for g evaluated at θ and the Hessian H evaluated at an unknown but fixed point ξ on the linear path between θ and $\hat{\theta}$. If we use a stochastic rounding algorithm that is unbiased, then the first order term can be neglected as

$$\mathbb{E}[g^\top (\hat{\theta} - \theta)] = 0,$$

and a high dimensional vector θ ensures the sum will concentrate around the expectation. This leaves the quadratic form with the Hessian.

A key property for low precision quantization of the weights (while minimizing the quadratic quantization error) is that the scale of the components of the eigenvectors of H do not differ by a large extent. If they do, then the quantization range must simultaneously provide coverage over a large range of values. This criterion is formalized through the notion of incoherence, introduced in Chee et al. (2024), which we briefly present below with a simplification of their analysis.

B.0.1 INCOHERENCE

A Hessian is μ -incoherent if the eigenvectors in the decomposition $Q\Lambda Q^\top = H \in \mathbb{R}^{N \times N}$ satisfy

$$\forall i, j : |Q_{ij}| \leq \mu/\sqrt{N},$$

and a parameter vector is μ -incoherent if it satisfies $\forall j : |\theta_j| \leq \mu \|\theta\|/\sqrt{N}$. Following QuIP (Chee et al., 2024), rather than quantizing the weights θ directly we will instead quantize the weights after applying a random orthogonal transformation matrix $P \in \mathbb{R}^{N \times N}$. Let $w = P^\top \theta$ and likewise $\theta = Pw$. Applying this rotation, the Hessian is transformed: $H_w = P^\top H_\theta P$ and likewise the eigenvectors Q from $H_\theta = Q\Lambda Q^\top$ are also multiplied $Q^\top \mapsto Q^\top P$.

If we choose P as a random Gaussian matrix: $\mathcal{N}(0, 1/N)^{N \times N}$, applying a rotation by Q^\top preserves the spherically symmetric distribution. Therefore the eigenvectors $Q^\top P$ of H_w are $\mathcal{N}(0, 1/N)^{N \times N}$ distributed. Applying a union bound over the Gaussian tail probability of the N^2 elements, the maximum absolute value entry of Q is at most $\sqrt{\frac{2 \log(2N^2/\delta)}{N}}$ with probability $1 - \delta$ and therefore incoherent with $\mu = \sqrt{2 \log(2N^2/\delta)}$.

B.0.2 SCALAR LDLQ

QuIP introduces the LDLQ quantization algorithm which quantizes weights sequentially and autoregressively taking into account how previous quantized values impact the quadratic Taylor expansion of the loss. Applying LDLQ to the entire vector of weights w rather than block by block, one has the following relation on the quantized weights \hat{w} . Let $L^\top DL = H_w$ be the LDL decomposition of $H_w = P^\top H_\theta P$, then we can express the quantization of the weights as

$$\hat{w} = \mathcal{Q}(w + (L - I)(w - \hat{w}))$$

where \mathcal{Q} quantizes the weights elementwise with nearest or unbiased stochastic rounding. As $L - I$ is a lower triangular matrix, the full \hat{w} can be quantized sequentially in an autoregressive manner. With this quantization scheme, Tseng et al. (2024) prove that the error of the quadratic in the Taylor expansion satisfies

$$(\hat{w} - w)^\top H(w - \hat{w}) \leq \frac{\mu^2 \sigma^2}{n} \text{Tr}(H^{1/2})^2, \quad (21)$$

where the pointwise quantization error of the scalar quantizer is assumed to be $\mathbb{E}[(\mathcal{Q}(x) - x)^2] \leq \sigma^2$ (see Theorem 4.1 applied to block size 1 and a single $N \times 1$ weight matrix), where σ^2 is a function of the bitrate. For example $x \in [0, 1]$, then a uniform grid would achieve $\sigma^2 = 2^{-2b-2}$ for b bits per parameter. From incoherence processing μ grows only logarithmically with the dimension N .

The achievable bits per parameter for a fixed quantization error depends crucially on the spectrum of the Hessian H , and how it scales with the dimension. We now pursue Krylov subspace routines with Hessian vector products to evaluate these quantities numerically for LLMs.

B.1 ESTIMATING $\text{Tr}(H^{1/2})$

To estimate the trace of the square root of the Hessian matrix, $\text{Tr}(H^{1/2})$, we begin by assuming that the Hessian is positive semi-definite (i.e., it contains no negative eigenvalues). The square root of the Hessian, denoted as S , can be expressed as:

$$S = \sum_{i=1}^P \sqrt{\lambda_i} \phi_i \phi_i^\top,$$

where λ_i and ϕ_i represent the eigenvalues and corresponding eigenvectors of the Hessian, respectively. Consequently, the trace of the square root of the Hessian is:

$$\text{Tr}(H^{1/2}) = \sum_{i=1}^P \sqrt{\lambda_i} = n \int_0^\infty p(\lambda) \sqrt{\lambda} d\lambda,$$

where $p(\lambda)$ is the spectral density function associated with the Hessian’s eigenvalues.

A direct computation of the full eigendecomposition to obtain $\text{Tr}(H^{1/2})$ has a computational complexity of $\mathcal{O}(n^3)$, which is infeasible for large models. Instead, we employ stochastic spectral density estimation techniques (Granzio et al., 2018; Pappas, 2019; Ghorbani et al., 2019), which scale linearly with the number of parameters. The key idea involves using the Pearlmutter trick (Pearlmutter, 1994) to efficiently compute Hessian-vector products:

$$\nabla(\nabla L^\top v) = H v,$$

where v is a random vector. This allows us to approximate the trace by leveraging the identity:

$$\text{Tr}(H) = \mathbb{E}[\text{Tr}(v v^\top H)] = \mathbb{E}[v^\top H v],$$

assuming v has zero mean and unit variance. These stochastic methods are well-established in machine learning (Fitzsimons et al., 2017; Dong et al., 2017).

Building upon the work of Ubaru et al. (2017), we can derive an explicit bound on the estimation of $\text{Tr}(H^{1/2})$ using stochastic Lanczos quadrature (SLQ).

Theorem B.1. Let $\mathbf{H} \in \mathbb{R}^{n \times n}$ be a symmetric positive definite matrix with eigenvalues ordered as $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$ and condition number $\kappa = \frac{\lambda_1}{\lambda_n}$. For any $\epsilon, \eta \in (0, 1)$, if the SLQ parameters satisfy

$$m \geq \frac{\sqrt{\kappa}}{4} \log \frac{K}{\epsilon} \quad (\text{Lanczos steps}),$$

$$n_v \geq \frac{24}{\epsilon^2} \log \frac{2}{\eta} \quad (\text{Rademacher vectors}),$$

where $K = (\lambda_{\max} - \lambda_{\min})(\sqrt{\kappa} - 1)^2$, then the output Γ of stochastic Lanczos quadrature satisfies:

$$\Pr \left[\left| \frac{\text{Tr}(\sqrt{\mathbf{H}}) - \Gamma}{\text{Tr}(\sqrt{\mathbf{H}})} \right| \leq \epsilon \right] \geq 1 - \eta.$$

The proof of this theorem is provided in [Appendix D](#). However, we observe that the bound on the trace provided here is overly conservative for practical purposes. Therefore, we also establish a result demonstrating self-averaging, which shows that the estimator converges to the true value based on a single random vector.

Theorem B.2. For a single random vector v , the signal-to-noise ratio of the trace estimator for a matrix $\mathbf{H} \in \mathbb{R}^{n \times n}$, where the spectral moments of \mathbf{H} do not depend on the matrix dimension, scales as:

$$\frac{\sqrt{\text{Var}(v^T \mathbf{H} v)}}{\mathbb{E}(v^T \mathbf{H} v)} = \mathcal{O}(n^{-\frac{1}{2}}).$$

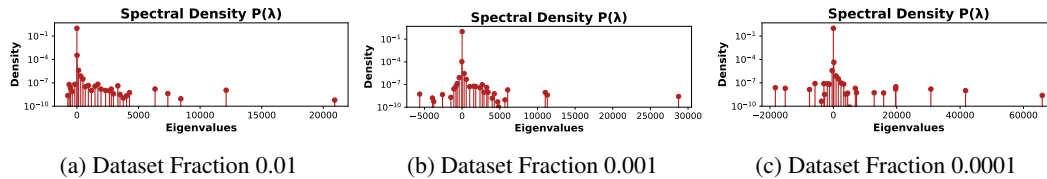


Figure 4: Spectral density plots of the 70M parameter Pythia model trained on varying fractions of the Pile dataset using the same data and random vector seed.

We utilize the CoLA ([Potapczynski et al., 2023](#)) library to compute the spectral approximation of the Hessian. This involves leveraging the relationship between the Lanczos T matrix and Gaussian quadrature ([Meurant & Strakoš, 2006](#); [Granzio et al., 2019](#)). However, these concepts are highly specialized and may not be familiar to all readers. Therefore, we provide a high-level overview without delving into the intricate mathematical details.

[Figure 4](#) illustrates the spectral density of a 70M parameter Pythia model trained on different subsets of the Pile dataset ([Gao et al., 2020](#)). Specifically, as we decrease the number of training samples—from 1% ([Figure 4a](#)) to 0.1% ([Figure 4b](#)) and further to 0.01% ([Figure 4c](#))—we observe an increase in the largest eigenvalue and an increase in the mass of negative spectral density. These phenomena are consistent with previous studies on ResNets and VGGs, where spiked Wigner random matrix theory models have been employed to understand such behaviors ([Granzio et al., 2022](#)).

Future work aimed at establishing a tighter empirical bound could explore advanced random matrix theory techniques ([Bun et al., 2017](#)), potentially utilizing the variance of the Hessian ([Granzio et al., 2022](#)). In this study, we adopt a simpler approach by shifting the Hessian spectrum by the magnitude of the largest negative eigenvalue, thereby ensuring a positive semi-definite Hessian and providing a trivial upper bound.

From [Figure 4b](#), we observe that the variance of each estimator remains low and that convergence is achieved with relatively few Lanczos iterations. Additionally, [Figure 5](#) demonstrates that varying the random vector introduces minimal variance, while different data subsets do exhibit some variance, as indicated by the error bars computed over three different seeds (see [Figure 5b](#)). For clarity, [Figure 5a](#) provides another example of the spectrum with a different seed vector on the same subsampled dataset, showing negligible differences compared to [Figure 4a](#).

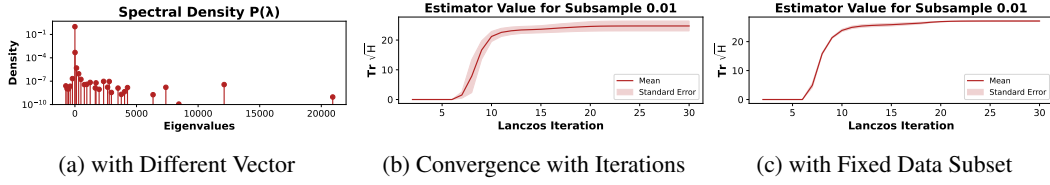


Figure 5: Comparison of spectral density and $\text{Tr}(\sqrt{\mathbf{H}})$ estimations for different subsample sizes and configurations.

With confidence in the accuracy of our estimations for $\text{Tr}(\mathbf{H}^{1/2})$ and the Hessian spectrum, we can interpret the implications for model quantization. Despite the high dimensionality and the presence of many distinct eigenvalues, the Hessian spectrum decays rapidly in density. This indicates that $\text{Tr}(\mathbf{H}^{1/2})$ grows sublinearly with the model dimension, rather than exhibiting the worst-case linear scaling. Consequently, as model size increases, the ratio $L(h)/D$ is expected to decrease, allowing for a more favorable tradeoff between the bitrate and the quantization gap. This supports the hypothesis that larger models on the compute-optimal frontier are more easily quantizable, thereby contributing to their improved generalization performance.

C STOCHASTIC TRACE ESTIMATION IMPROVEMENT WITH MODEL SIZE

Here, we provide the proof that for a spectrum independent of model dimension, the stochastic trace estimator has a bigger signal to noise ratio as a function of dimension.

Lemma C.1. *Let $\mathbf{u} \in \mathbb{R}^{P \times 1}$ random vector, where \mathbf{u}_i is zero mean and unit variance and finite 4th moment $\mathbb{E}[\mathbf{u}_i^4] = m_4$. Then for $\mathbf{H} \in \mathbb{R}^{P \times P}$, then*

- (i) $\mathbb{E}[\mathbf{u}^T \mathbf{H} \mathbf{u}] = \text{Tr} \mathbf{H}$,
- (ii) $\text{Var}[\mathbf{u}^T \mathbf{H} \mathbf{u}] \leq (2 + m_4) \text{Tr}(\mathbf{H}^T \mathbf{H})$.

Proof. For the expectation, we see

$$\mathbb{E}[\mathbf{u}^T \mathbf{H} \mathbf{u}] = \sum_{i,j=1}^P \mathbf{H}_{i,j} \mathbb{E}[\mathbf{u}_i \mathbf{u}_j] = \sum_{i=1}^P \mathbf{H}_{i,i} = \text{Tr} \mathbf{H}.$$

For the variance, we have

$$\begin{aligned} \mathbb{E}[|\mathbf{u}^T \mathbf{H} \mathbf{u}|^2] &= \sum_{i,j} \sum_{k,l} \mathbf{H}_{i,j} \mathbf{H}_{k,l}^T \mathbb{E}[\mathbf{u}_i \mathbf{u}_j^T \mathbf{u}_k \mathbf{u}_l^T] \\ &= \sum_{i,j} \sum_{k,l} \mathbf{H}_{i,j} \mathbf{H}_{k,l}^T [\delta_{i,j} \delta_{k,l} + \delta_{i,l} \delta_{j,k} + \delta_{i,k} \delta_{j,l} + m_4 \delta_{i,j,k,l}] \\ &= (\text{Tr} \mathbf{H})^2 + (2 + m_4) \text{Tr}(\mathbf{H}^2), \end{aligned}$$

whence (ii) follows.

Let us consider the signal to noise ratio for some positive definite $\mathbf{H} \succ c\mathbf{I}$

$$\frac{\sqrt{\text{Var}[\mathbf{u}^T \mathbf{H} \mathbf{u}]}}{\mathbb{E}[\mathbf{u}^T \mathbf{H} \mathbf{u}]} = \sqrt{2 + m_4} \sqrt{\frac{\text{Tr} \mathbf{H}^2}{\text{Tr}^2 \mathbf{H}}} = \sqrt{\frac{2 + m_4}{P}} \sqrt{\frac{\langle \lambda^2 \rangle}{\langle \lambda \rangle^2}} \quad (22)$$

where we denote the mean eigenvalue $\langle \lambda \rangle$ and the mean square eigenvalue similarly. \square

Remark C.2. Note that m_4 is 3 for the Gaussian case and 1 for the Hutchinson trace estimator where the entries are ± 1 with probability half, which justifies its use.

C.1 DERIVING THE IMPACT OF LOW PRECISION LANCZOS

Consider a number taken from our Hessian matrix $a_{i,j}$, which can be represented as $(-1)^s 2^e s$. As the exponent for FP16 has 5 bits, it has a range of $2^5 - 1$. Since the exponent is always integer, there is no loss of information in the range. This means the error is in the significand, which has 6 bits after the 1. Thus, we have $\epsilon = 10^{-7}$.

Then, we see that

$$\begin{aligned} \tilde{\mathbf{H}} &= \begin{bmatrix} a_{1,1}(1 + \mathcal{N}(0, \epsilon)) & a_{1,2}(1 + \mathcal{N}(0, \epsilon)) & \cdots & a_{1,n}(1 + \mathcal{N}(0, \epsilon)) \\ a_{2,1}(1 + \mathcal{N}(0, \epsilon)) & a_{2,2}(1 + \mathcal{N}(0, \epsilon)) & \cdots & a_{2,n}(1 + \mathcal{N}(0, \epsilon)) \\ \vdots & \vdots & \ddots & \vdots \\ a_{m,1}(1 + \mathcal{N}(0, \epsilon)) & a_{m,2}(1 + \mathcal{N}(0, \epsilon)) & \cdots & a_{m,n}(1 + \mathcal{N}(0, \epsilon)) \end{bmatrix} \\ &= \mathbf{H} + \begin{bmatrix} a_{1,1}\mathcal{N}(0, \epsilon) & a_{1,2}\mathcal{N}(0, \epsilon) & \cdots & a_{1,n}\mathcal{N}(0, \epsilon) \\ a_{2,1}\mathcal{N}(0, \epsilon) & a_{2,2}\mathcal{N}(0, \epsilon) & \cdots & a_{2,n}\mathcal{N}(0, \epsilon) \\ \vdots & \vdots & \ddots & \vdots \\ a_{m,1}\mathcal{N}(0, \epsilon) & a_{m,2}\mathcal{N}(0, \epsilon) & \cdots & a_{m,n}\mathcal{N}(0, \epsilon) \end{bmatrix}. \end{aligned}$$

Now under certain assumptions on the elements of the perturbation matrix (essentially the $a_{i,j}$ does not vary too wildly or have wild dependencies) this becomes a Gaussian orthogonal ensemble (GOE) again. Then using the Frobenius Norm, we see that the spectral width will be of order $\epsilon \sqrt{\langle \lambda^2 \rangle}$, which depends on the square root of the average eigenvalue squared of \mathbf{H} . Anything within this will be noise. This is because $\sum_{i,j} a_{i,j}^2 = P \langle \lambda^2 \rangle$. An obvious upper bound of this would be $\epsilon \lambda_1$ but this will likely be super loose. Note that the vast majority of the already broadened spectrum is already very close to zero, so we would expect this to be even more extreme for the unbroadened version. A better strategy might be to sample the noisy version of $a_{i,j}^2$ perhaps using the diagonal approximation, and note that in expectations we expect the square to be $(1 + \epsilon^2)$ the size of its non noisy counter part, which gives an estimation equation

$$\sqrt{\frac{P\epsilon^2 \sum_k a_{i,j}^2}{N(1 + \epsilon^2)}}.$$

D STOCHASTIC LANCZOS QUADRATURE PROOF

Theorem D.1. Consider a symmetric positive definite matrix $\mathbf{A} \in \mathbf{R}^{n \times n}$ with eigenvalues enumerated in reverse order of size $\lambda_1 \geq \lambda_2 \cdots \geq \lambda_n$ and condition number $\kappa = \frac{\lambda_1}{\lambda_n}$. For $\epsilon, \eta \in (0, 1)$ and SLQ parameters satisfying

$$(i) \quad m \geq \frac{\log \frac{\kappa}{\epsilon}}{2 \log \frac{\sqrt{\kappa} + 1}{\sqrt{\kappa} - 1}} \geq \frac{\sqrt{\kappa}}{4} \log \frac{\kappa}{\epsilon} \text{ Lanczos steps}$$

$$(ii) \quad n_v \geq \frac{24}{\epsilon^2} \log \frac{2}{\eta} \text{ Rademacher vectors,}$$

where $K = (\lambda_{max} - \lambda_{min})(\sqrt{\kappa} - 1)^2$. The output Γ of stochastic lanczos quadrature is such that

$$\Pr \left[\left| \frac{\text{Tr}(\sqrt{\mathbf{A}}) - \Gamma}{\text{Tr}(\sqrt{\mathbf{A}})} \right| \leq \epsilon \right] \geq 1 - \eta \quad (23)$$

Proof. The proof follows trivially from Ubaru et al. (2017), where we simply take the more general proof and instead of the general function $f(\mathbf{A})$, we take $f(x) = \sqrt{x}$. The second inequality for m is directly from the paper, but the tighter bound is also available just buried. \square

The proof sketch goes as follows. We bound the error from the Gauss quadrature rule. We start with a function analytic in the interval $[-1, 1]$. Knowing that the Gauss quadrature rule is exact for any

1296 polynomial up to degree $2m + 1$, we bound the sum from $2m + 1$ to infinity using Cauchy-Schwarz.
1297 We use results from Chebyshev coefficients, symmetry and the interval boundaries to get

$$1298 \quad |I - I_m| \leq \frac{4\sqrt{\lambda_1}}{(\rho^2 - 1)\rho^{2m}},$$

1301 where ρ is the sum of the major and minor axis of the Bernstein ellipse. We shift the spectrum so that it
1302 is in the interval $[-1, 1]$, e.g this implies the factor of $\frac{\lambda_1 - \lambda_n}{2}$. The shifted function is not analytic for
1303 $\alpha = -\frac{-\kappa+1}{\kappa-1}$, so this will serve as our major axis. Now as $\frac{x^2}{a^2} + \frac{y^2}{b^2} = 1$ and the focus is $1 = \sqrt{a^2 - b^2}$,
1304 where we take our major axis a in this case to be α . We then have our rate of convergence $\rho = a + b$
1305 through some algebra to be $\frac{\sqrt{\kappa+1}}{\sqrt{\kappa-1}}$. This gives us the value of K . This is combined with the error of
1306 the trace estimator from [Roosta-Khorasani & Ascher \(2015\)](#) and Cauchy-Schwarz to obtain the final
1307 result.
1308

1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349