
On Information Geometry and Iterative Optimization in Model Compression: Operator Factorization

Zakhar Shumaylov^{1,2*}, Vasileios Tsiaras¹, Yannis Stylianou¹
¹Apple, ²University of Cambridge

Abstract

The ever-increasing parameter counts of deep learning models necessitate effective compression techniques for deployment on resource-constrained devices. This paper explores the application of information geometry, the study of density-induced metrics on parameter spaces, to analyze existing methods within the space of model compression, primarily focusing on operator factorization. Adopting this perspective highlights the core challenge: defining an optimal low-compute submanifold (or subset) and projecting onto it. We argue that many successful model compression approaches can be understood as implicitly approximating information divergences for this projection. We highlight that when compressing a pre-trained model, using information divergences is paramount for achieving improved zero-shot accuracy, yet this may no longer be the case when the model is fine-tuned. In such scenarios, trainability of bottlenecked models turns out to be far more important for achieving high compression ratios with minimal performance degradation, necessitating adoption of iterative methods. In this context, we prove convergence of iterative singular value thresholding for training neural networks subject to a soft rank constraint. To further illustrate the utility of this perspective, we showcase how simple modifications to existing methods through softer rank reduction result in improved performance under fixed compression rates.

1 Introduction

The last decade in deep learning has witnessed a trend towards increasingly larger models, exemplified by the rise of transformer-based architectures for large language models (LLMs), which have achieved state-of-the-art results in various natural language processing tasks. However, the substantial size of these models poses challenges for deployment on resource-constrained devices, in terms of computational cost, memory footprint, and energy consumption, necessitating model compression as a crucial strategy for mitigating deployment costs.

Although theoretical limits suggest the possibility of arbitrarily narrow models [71], practical experience demonstrates the difficulties in training and generalizing such compact architectures [39]. Model compression techniques aim to address this challenge by instead training large networks, utilizing the full power of overparameterization [5], with the aim of reducing the computational and memory footprint of these models post-training, preserving as much performance as possible. Following the categorization presented by [30], existing compression techniques broadly fall into five categories. Quantization methods reduce the bit-width of weights and activations (e.g. [49, 83, 42]). Pruning techniques induce sparsity in weight matrices (e.g. [78, 50, 22, 47, 41]). Distillation methods leverage a teacher-student framework to transfer knowledge to a smaller model (e.g. [27, 56, 65]). Parameter sharing methods re-use various network components [85, 64]. And finally, operator factorization methods decompose weight matrices and tensors into multiple components, which individually are significantly cheaper (e.g. [24, 52, 12, 60, 18, 46, 74, 89]), often resulting in improved scaling

*Work done during an internship at Apple.

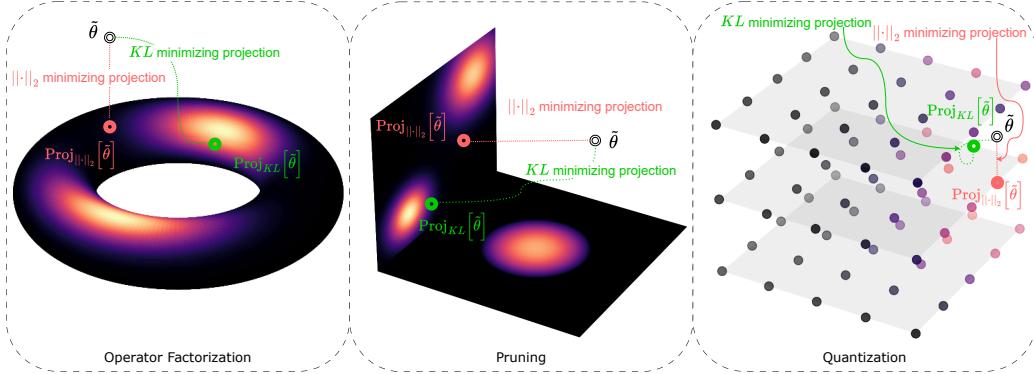


Figure 1: Illustration of the main components of compression methods, highlighted by the main three directions: (1) Operator Factorization, with torus depicting the low rank manifold; (2) Pruning, with two planes depicting the sparse vector space; (3) Quantization, with lattice depicting the level of quantization. On all three, the difference between euclidean versus information projection on the low-compute manifold/set. The color here depicts the loss of the resulting model, highlighting that in a single projection step KL minimizing projection result in significantly improved models.

laws [66, 68]. In this paper we focus on operator factorization, but note that techniques above are often complementary, and combined approaches frequently yield the best compression results [32, 73, 97, 95], although are non-orthogonal, and must be combined mindfully [26]. Recent years have seen increased theoretical attention devoted to pruning techniques, leading to more theoretically grounded methods with improved performance (e.g. [92, 45, 54]). However, within operator factorization, the focus has primarily remained on empirical evaluations (e.g. [33, 34, 35, 63]), with largest focus on language modelling tasks. This disparity raises a critical question: can a more rigorous theoretical understanding be developed for operator factorization methods?

This paper addresses this question directly from two perspectives. Firstly, by establishing a precise mathematical formulation of the problem, we demonstrate that the appropriate framework for understanding model compression lies within information geometry [3]. Similar observations have previously been highlighted in the context of pruning of natural parameters by [51], but has yet to make its way into deep learning. This perspective illuminates why Fisher/Hessian-based formulations are often the most effective, and we explicitly identify the various approximations inherent in existing methods. Secondly, we consider the effect of iterative compression methods in analogy with [92], highlighting that iterative compression implicitly adjusts the optimization problem, highlighting its connection to iterative singular value hard thresholding. To illustrate the utility of this formulation, we analyze a recently proposed method of [16], highlighting limitations of the proposed singular value cutoff criterion, and, leveraging our iterative information geometric framework, derive modifications to methods of [16] and [93], through additions of Fisher information.

2 Preliminaries

We consider a supervised problem of predicting outputs $y \in \mathbb{Y}$ from inputs $x \in \mathbb{X}$. We assume a probabilistic model for the conditional distribution in the form $p_\theta(y|x) = p(y|f(x, \theta))$, where $p(y|\cdot)$ is taken to be an exponential family with natural parameters in \mathbb{F} and $f : \mathbb{X} \times \mathbb{R}^D \rightarrow \mathbb{F}$ is a prediction function parameterized by $\theta \in \mathcal{M} = \mathbb{R}^D$. Given N iid training samples $(x_n, y_n)_{n=1}^N$, we aim to minimize the negative log-likelihood (NLL):

$$\mathcal{L}(\theta) := - \sum_n \log p_\theta(y_n|x_n) = - \sum_n \log p(y_n|f(x_n, \theta)). \quad (1)$$

This framework covers most common scenarios such as least-squares regression with $\mathbb{X} = \mathbb{R}^{\dim \mathbb{X}}$, $\mathbb{Y} = \mathbb{R}^{\dim \mathbb{Y}}$ with $p(y|f) = \mathcal{N}(y; f, I_{\dim \mathbb{Y}})$ or C -class classification with cross-entropy loss $\mathbb{Y} = \{1, \dots, C\}$, $\mathbb{F} = \mathbb{R}^C$ and $p(y=c|f) = \pi_c(f) := \exp(f_c) / \sum_i \exp(f_i)$. Classically, minimizing Equation (1) involves using (stochastic) gradient descent. However, adaptive schemes with preconditioning or Gauss-Newton updates, such as natural gradient descent or its approximations yield faster and stabler convergence [21, 8, 25, 53]. These methods are able to align with

the (information) geometry of the problem by pulling back the distance onbetween parameters via the Kullback-Leibler divergence between the corresponding probability distributions, rather than the Euclidean distance. The intuition can be summarized using the following two facts, implicitly considering $f : \mathbb{X}^N \times \mathbb{R}^D \rightarrow \mathbb{F}^N$ for ease of presentation:

$$\text{KL}(p_\theta | p_{\theta+\Delta\theta}) = \frac{1}{2} \Delta\theta^\top \mathcal{I}(\theta) \Delta\theta + \mathcal{O}(\Delta\theta^3), \quad \text{and} \quad \mathcal{L}(\theta) = \text{KL}(p_{\text{em}} | p_\theta) + \text{const.} \quad (2)$$

for p_{em} the empirical distribution, \mathcal{I} denoting the fisher information metric (FIM), coinciding with a generalized Gauss-Newton [76] approximation of the Hessian [44]:

$$\mathcal{I}(\theta) := \sum_n \mathbb{E}_{p_\theta(y|x_n)} \left[\nabla_\theta \log p_\theta(y|x_n) \nabla_\theta \log p_\theta(y|x_n)^\top \right].$$

Thus, for the problems considered the distance induced by KL divergence, represented infinitesimally using the FIM, appears to be most fitting.

3 Model Compression

The central goal of model compression is to derive a *computationally efficient* and *generalizable* model for a given task in Equation (1). The simplest, and seemingly straightforward approach is to directly train a small model. While theory suggests that even extremely narrow networks can approximate any desired function [71], training such networks in practice proves incredibly challenging, if not impossible [87, 7]. Consequently, the standard practice involves training an overparameterized model, which is then compressed. Though effective, this two-stage process may not actually be the most efficient use of computational resources under a fixed budget [10].

To illustrate the core concepts, we consider a simplified framework for the case of low rank factorization where $\mathcal{M} = \mathbb{R}^{\sum n_i \times n_{i+1}}$ represents the space of weight matrices, and $\mathcal{M}_{<r} = \{\theta \in \mathbb{R}^{\sum n_i \times n_{i+1}} \mid r_i \leq \text{rank } \theta \leq r\}$ denotes the variety of lower-rank matrices, where r_i represents the rank constraint for the i -th layer. For a comprehensive discussion of the properties of these sets, see [29]. Considering $\mathcal{M}_{<r}$ as the low compute subset reduces the number of parameters for each layer from $n_i \times n_{i+1}$ in \mathcal{M} to $r_i \times (n_i + n_{i+1})$, meaning that for sufficiently small r_i these become significantly cheaper, going from a quadratic to a linear number of parameters with respect to dimensions of the matrix. We emphasize that the discussion here generalizes to the setting of pruning and quantization, in which cases the set $\mathcal{M}_{<r}$ would correspond to a linear subspace or lattice correspondingly. Generally, the bilevel problem of interest can be written as

$$\min_r \|r\|_1 \quad \text{subject to} \quad \min_{\theta \in \mathcal{M}_{<r}} \mathcal{L}(\theta) \leq \varepsilon,$$

where $\|r\|_1$ represents a measure of the compression level (e.g., number of parameters), $\mathcal{L}(\theta)$ is the loss function, and ε is an acceptable error threshold. Directly tackling such a bilevel constrained minimization is not computationally feasible by the virtue of r_i being integers and constraint set being non-convex, making the optimization problem not approachable using gradient based methods. Although relaxation based methods exist via masking [23], nuclear norm regularization [93] or alternative formulations [92]. As is common with optimization problems, the problem needs to be approached in an iterative manner. As mentioned above, the most common approach in model compression is to first find a full over-parameterized solution, after which the appropriate $r = r(\tilde{\theta})$ is identified and $\tilde{\theta}$ is projected as $\theta_r = \text{proj}_{\mathcal{M}_{<r}(\tilde{\theta})}(\tilde{\theta})$. As projection may result in performance degradation, the resulting model is fine-tuned:

$$\tilde{\theta} = \arg \min_{\theta \in \mathcal{M}} \mathcal{L}(\theta) \xrightarrow{r(\tilde{\theta})} \theta_r = \arg \min_{\theta \in \mathcal{M}_{<r(\tilde{\theta})}} \mathcal{L}(\theta).$$

Such projections can either be done once, resulting in so called *train-then-sparsify* methods, or iteratively, resulting in so called *sparsify-during-training* methods [30]. As long as the minimum is on the restricted set, it should be possible to find it, as $\mathcal{M}_{<r_i}$ is simply connected [84].

To summarize, any model compression algorithm must address two main questions. That is **rank selection**, and **projection**. In Section 3.2, we summarize existing approaches in the literature surrounding operator factorization, showcasing that their heuristic success can be attributed to improvements in either rank or projection. Based on the discussion in the previous section, we argue that projection must be performed according to information distance, with rank selection according to the resulting distance, in terms of KL divergence, to the projection.

3.1 Information Geometry in Model Compression

Following the framework of [2], we begin by considering a family of probability distributions $\mathcal{D} = \{p_\theta(\cdot) \mid \theta \in \mathcal{M} \subseteq \mathbb{R}^d\}$, parameterized by $\theta \in \mathcal{M}$. This family forms a statistical manifold, with θ serving as a local coordinate system. This manifold is then equipped with a *divergence* $D(p_\theta \| p_{\theta'})$, quantifying dissimilarity between distributions p_θ and $p_{\theta'}$ on \mathcal{D} . In general, this divergence does not need to be symmetric or satisfy the triangle inequality. An important class of divergences in this context are the *Bregman divergences* generated by a strictly convex and differentiable function ϕ :

$$D_\phi(\theta \| \theta') = \phi(\theta) - \phi(\theta') - \nabla\phi(\theta)^\top(\theta - \theta'). \quad (3)$$

For exponential and mixture families, the *Kullback-Leibler (KL) divergence* $\text{KL}(p_\theta \| p_{\theta'}) = \mathbb{E}_{p_\theta}[\log(p_\theta(x)/p_{\theta'}(x))]$ arises naturally and can be viewed as a Bregman divergence, generated by the entropy or log-partition function [3]. The local structure of the statistical manifold arises from this divergence. Specifically, considering an infinitesimal perturbation $d\theta$ around a point θ , the second-order Taylor expansion of the divergence $D(p_\theta \| p_{\theta+d\theta})$ defines the *Riemannian metric tensor* $g(\theta)$ with components: $g_{ij}(\theta) = \frac{\partial^2}{\partial\theta^i\partial\theta^j} D(p_\theta \| p_{\theta'})|_{\theta=\theta'}$. This metric endows \mathcal{D} with a Riemannian structure, allowing us to define notions of distance and orthogonality. As is often considered in information geometry [6], we are also interested in dual connections generated by the divergence [2]:

$$\Gamma_{ijk}(\theta) = -\frac{\partial^3}{\partial\theta_i\partial\theta_j\partial\theta'_k} D(\theta \| \theta')|_{\theta=\theta'}, \quad \Gamma_{ijk}^*(\theta) = -\frac{\partial^3}{\partial\theta'_i\partial\theta'_j\partial\theta_k} D(\theta \| \theta')|_{\theta=\theta'},$$

which are dually coupled with respect to g_{ij} [4]. The dual connections Γ and Γ^* result in two families of geodesics, termed the *e-geodesics* and *m-geodesics*, respectively. These play a crucial role in defining projections onto submanifolds within the statistical manifold.

Iterative m-projection for Pruning in Flat Manifolds A manifold in coordinates θ is *e-flat* (exponential-flat), if $\Gamma_{ijk} = 0$, and *m-flat* (mixture-flat) if $\Gamma_{ijk}^* = 0$. A typical example of an *e-flat* manifold is the exponential family: $p(x, f) = \exp\{\sum f_i k_i(x) - \psi(f)\}$, where $k_i(x)$ are given functions and ψ is the normalizing factor. In such manifolds, coordinate pruning, as in [51], can be viewed as iterative *m-projection* and analyzed as follows. Suppose we are given two flat submanifolds of \mathcal{M} with $p_1 \in E_1 \subset E_2$. Then for p_2^* the *m-projection* of p_1 onto E_2 , $p_2^* = \arg \min_{p \in E_2} D(p_1 \| p)$, we have that p_2^* is given by the *m-geodesic* connecting p_1 and p_2^* orthogonal to E_2 at p_2^* (e.g. by Proposition 3.1 or [51]). Then, for any $p_2 \in E_2$, the Generalized Pythagorean theorem [6, 59] results in: $D(p_1 \| p_2) = D(p_1 \| p_2^*) + D(p_2^* \| p_2)$. This allows for direct analysis of error accumulation under iterative projection onto lower-dimensional (pruned) flat submanifolds.

Iterative m-projection in Curved and Overparameterized Manifolds In the case of overparameterized models, the mapping $f : \mathcal{M} \rightarrow \mathbb{F}$, the image $f(\mathcal{M})$ may not be a simple submanifold, potentially containing non-differentiable structures, making theory for curved exponential families [3, 1] unusable. While the generalized pythagorean theorem may still be used when the image is convex [40], the case of deep neural networks does not fall into any of these categories. Luckily, for the problem of model compression, the global structure of the parameterization is not of interest, only the local structure is. This is further helped by the fact that for operator factorization [48, Example 8.14] and pruning (but not quantization), the low-compute subset is a smooth manifold, implying that the generalized pythagorean theorem can be used as long as locally, the parameterization image results in a submanifold [6, Corollary 4.2]. How can this be achieved?

We consider an open ball around the solution we found $B(\tilde{\theta})$ such that intersection with the low rank manifold \mathcal{M}_r is non-empty. Then, it is open and a submanifold [48]. We can then consider the set $f(\mathcal{M}_r \cap B(\tilde{\theta}))$. Since $\mathcal{M}_r \cap B(\tilde{\theta})$ is an open subset of \mathcal{M}_r , it is a submanifold and therefore a manifold. Thus considering $f|_{B(\tilde{\theta})}$, assuming it is a submersion, by Lemma A.3 below we know that $f(\mathcal{M}_r \cap B(\tilde{\theta})) \in \mathbb{F}$ is a submanifold, establishing information projections through geodesics:

Proposition 3.1 (Corollary 4.2 [6]). *Let N be a differentiable submanifold of \mathcal{M} . Then $q \in N$ is a stationary point of the function $D(p \| \cdot) : N \rightarrow \mathbb{R}, r \mapsto D(p \| r)$ iff the *m-geodesic* from p to q meets N orthogonally.*

3.2 Related Works in Operator Factorization

SVD, TRP, DLRT The first class of approaches are based on Euclidean projection of matrices onto their low-rank variants, which by Eckart Youngs Theorem is equivalent to SVD:

$$\text{proj-}\|\cdot\|_{2,\mathcal{M}_{<r}}(\tilde{\theta}) = \arg \min_{\theta \in \mathcal{M}_{<r}} \|\theta - \tilde{\theta}\|_2^2. \quad (4)$$

There have been a multitude of papers on various variants of such a technique starting with [67] and many since [52, 12, 60, 18, 46, 36, 74]. This approach has been successful in compressing models and sometimes even improving generalization by reducing overfitting during training [91, 62].

These approaches have also been adapted to be used during training via iterative projection [93, 75, 16]. The main benefit of using SVD directly is that minimizers can be found efficiently separately per layer. For classical SVD, the rank has to be provided in advance, or similar to iterative methods can be chosen based on the values directly. However, beyond heuristics it is not clear how to choose the rank. Both the distribution and size of singular values differs significantly amongst layers [24] and in transformers is significantly different between attention and MLP weights [97, 96].

FWSVD, TFWNSVD, RankDyna The next class of approaches we are going to consider is based precisely on the discussion of Section 3, with methods approximating the information projection:

$$\text{proj-KL}_{\mathcal{M}_{<r}}(\tilde{\theta}) = \arg \min_{\theta \in \mathcal{M}_{<r}} \text{KL}(p_{\tilde{\theta}} \mid p_{\theta}) \approx \arg \min_{\theta \in \mathcal{M}_{<r}} \frac{1}{2} (\theta - \tilde{\theta})^\top \mathcal{I}(\tilde{\theta})(\theta - \tilde{\theta}) + \mathcal{O}((\theta - \tilde{\theta})^3). \quad (5)$$

However, computing the full FIM is computationally prohibitive due to its quartic complexity with respect to the weight matrix size. Thus, approximations must be used. Methods like Fisher-Weighted SVD (FWSVD) [33] and its variants like TFWNSVD [34], along with the iterative projection method RankDyna [35], implicitly approximate this KL projection, although they often frame the FIM as representing parameter “importance” [57] related to the loss function. These approaches can be viewed as attempting to minimize the KL divergence under significant approximations: TFWNSVD diagonalizes the fisher information $\mathcal{I}(\theta) \approx \text{diag } \mathcal{I}(\theta) := \tilde{\mathcal{I}}$, decoupling layers and discarding potentially crucial intra-layer interactions. But even with a diagonal FIM, the resulting weighted SVD problem remains computationally expensive [79]. To address this, FWSVD further approximates the diagonal FIM by summing the diagonal elements column-wise, making the approximation identical for each row: $\tilde{\mathcal{I}}(\theta_k)_{ij} \approx \text{diag} \sum_j \tilde{\mathcal{I}}(\theta_k)_{ij} := \tilde{\mathcal{I}}(\theta_k)_{ij}$ for layer k parameters $\theta_k = \text{vec } W^k$. As in the Euclidean case, the rank has to be provided in advance or chosen based on the values of the singular values, although selection criteria do exist [31]. Analogously, importance based weighting [57] can be described through this lens.

As the overarching goal is to minimize the loss $\text{KL}(p_{\text{em}} \mid p_{\theta})$, we may attempt to instead expand the expression above around $\tilde{\theta}$, proposed by [57]. While at the minimum the first order term is zero, in practice it is not exactly equal to zero and the approximation itself may not be accurate.

$$\begin{aligned} \text{proj-KL}_{\text{em}, \mathcal{M}_{<r}}(\tilde{\theta}) &= \arg \min_{\theta \in \mathcal{M}_{<r}} \text{KL}(p_{\text{em}} \mid p_{\theta}) \\ &\approx \arg \min_{\theta \in \mathcal{M}_{<r}} (\theta - \tilde{\theta})^\top \nabla_{\theta} \text{KL}(p_{\text{em}} \mid p_{\theta}) \\ &\quad + \frac{1}{2} (\theta - \tilde{\theta})^\top \mathcal{I}_{\text{em}}(\tilde{\theta})(\theta - \tilde{\theta}) + \mathcal{O}((\theta - \tilde{\theta})^3). \end{aligned} \quad (6)$$

ASVD, CALDERA, OATS Due to the inherent cost of evaluating the FIM, many works [94, 90, 97] have instead focused on ensuring agreement of intermediary (post-) activations, turning it into a proxy for the full information distance. Distributions of intermediary activations are quite particular, oftentimes possessing significant outliers [97, 94]. If we denote $\theta_i \in \mathbb{R}^{n_i \times n_{i+1}}$ and X_{i-1} as precomputed post-activations or some statistic of post-activations, the approach is to minimize the average reconstruction error, which from the point of view of information geometry corresponds to finding projection under per-layer gaussian density assumption, with $p_{\theta}^{\mathcal{N}}(\cdot) = \prod_i \mathcal{N}(\cdot | X_i, I_{n_i \times n_i})$:

$$\text{proj-KL}_{\mathcal{M}_{<r}}^{\mathcal{N}}(\tilde{\theta}) = \arg \min_{\theta \in \mathcal{M}_{<r}} \sum_i \|\theta_i X_{i-1} - \tilde{\theta}_i X_{i-1}\|_2^2 = \arg \min_{\theta \in \mathcal{M}_{<r}} \text{KL}(p_{\theta}^{\mathcal{N}} \mid p_{\tilde{\theta}}^{\mathcal{N}}) \quad (7)$$

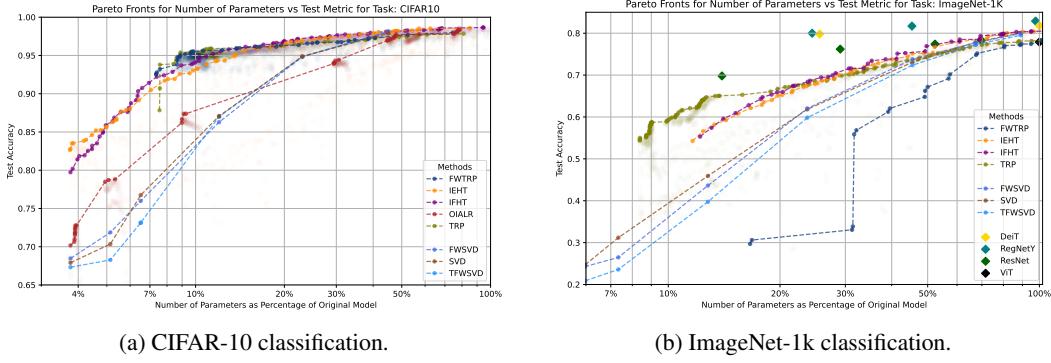


Figure 2: Pareto fronts for compression methods in Section 3 for ViT-B/16 for image classification. Diamonds show performance reported by [82].

Taking on such a view raises a question: is such a gaussian assumption the right one? Oftentimes post-activations cluster around multiple points (e.g. for different classes), or may contain outliers. As before, the rank has to be provided in advance or chosen based on the value of the singular values. Other heuristics also exist, e.g. [94] performs an expensive binary search for each rank depending on the accuracy drop.

4 Sparsify During Training as Iterative Hard Thresholding

In the previous section we discussed information geometric aspects involved in the *projection* step. However, the question *rank selection* is addressed only heuristically. Based on the observation that training/fine-tuning bottlenecked networks to good accuracy is not always possible, we are motivated to consider *sparsify-in-training* methods, iteratively reducing the rank.

The Deep Linear Network Example To illustrate the rationale behind iterative rank reduction, we consider the case of deep linear networks. While simple, these models have proven valuable for studying overparameterization, revealing insights into implicit biases and optimization benefits [55]. A key observation in overdetermined problems is the tendency for solutions to exhibit low-rank structures, a phenomenon captured by analyzing the entropy of low-rank solutions within the solution space, with low-rank solutions possessing significantly higher entropy. However, convergence to such a low-rank minimum is not guaranteed [15], even if they are more likely. Therefore, iterative rank reduction methods can be interpreted as a means of guiding the optimization process towards these desirable low-rank minima, effectively “massaging” the solution towards a compressed representation.

Iterative Hard Thresholding in Low Rank Approximation Based on the observations, that iterative rank reduction [16, 75] is a projection onto the low-rank manifold and that hard thresholding of singular values is a proximal step on the rank function, we consider the resulting minimization problem and proximal gradient method used to optimize it. Similar observations, have been made in the context of pruning [92]. Low-rank approximation is a fundamental problem in various fields with a rich literature of algorithms to address it, e.g. in the context of matrix completion [86]. Because the rank of a matrix is not a continuous function, many methods relax this constraint to allow for continuous optimization, including nuclear norm minimization [11, 69] as a convex relaxation; via low-rank factorization [37, 80], resulting in a non-convex least-squares problem; or considering a rank constrained formulation approached via singular value thresholding [14, 81].

While the methods are effective, they often operate under simplifying assumptions on the objective, which make them non-applicable to the setting of neural networks, necessitating an extension of the existing analyses. The problem we are going to consider here is going to be of the form (for the sake of simplicity here presented as for a single matrix, but extends naturally beyond):

$$\min_{W \in \mathcal{M}} \mathcal{L}(W) + \lambda \text{rank}(W). \quad (8)$$

Now, the classical proximal gradient method for Equation (8) consists of iterative steps of the form $W_{n+1} = \text{prox}_{\alpha_n \lambda \text{rank}(\cdot)}(W_n - \alpha_n \nabla \mathcal{L}(W_n))$. The power of the proximal method comes from the

fact that normally the proximal operator of a function is expensive to compute. Intuitively proximal operator of rank involves projection on the lower rank manifold, which is extremely non-convex and more so has empty interior. However, thanks to Eckart Youngs Theorem we have a closed form way of performing such a projection [29, 28]. More generally however, we are going to consider a (mirror) proximal gradient method, with divergences (Equation (3)), instead of the l_2 norm, with updates of the form

$$W_{n+1} \in \operatorname{argmin}_{W \in \mathcal{M}} \left\{ \frac{1}{\alpha_n} D_{F_n}(W, W_n) + \langle W, \nabla \mathcal{L}(W_n) \rangle + \lambda \operatorname{rank}(W) \right\}. \quad (9)$$

Motivation for such updates is based on discussion of Section 3.2: information induced distances are of more interest. Ideally, instead of $D_{F_n}(W, W_n)$ we would consider $\text{KL}(p(W_n) | p(W))$, however this is captured sufficiently well by the FIM, for which $F_n = \frac{1}{2} \|\cdot\|_{\mathcal{I}(W_n)}^2$.

4.1 Proximal Gradient Algorithm Analysis

In this section, we consider the problem in Equation (8) and analyze the iterative scheme of Equation (9). To show convergence, we assume the following:

- \mathcal{L} is a bounded from below Fréchet differentiable function with Lipschitz continuous gradient, i.e. there exists $L_{\nabla \mathcal{L}} \geq 0$ such that $\|\nabla \mathcal{L}(W) - \nabla \mathcal{L}(W')\| \leq L_{\nabla \mathcal{L}} \|W - W'\|$. Further assume that \mathcal{L} is sub-analytic and coercive.
- Each $F_n : \mathbb{R}^m \rightarrow \mathbb{R}$, assumed to be σ_n -strongly convex, Fréchet differentiable and such that ∇F_n is $L_{\nabla F_n}$ -Lipschitz continuous with $L_{\nabla F_n} \leq \overline{L}_{\nabla F}$.

Note 4.1. The lipschitzness assumption on \mathcal{L} is not necessary, and can be generalized to a local version [38]. Coercivity is not restrictive, as oftentimes weight decay is used. Sub-analyticity condition on \mathcal{L} is not restrictive, as by e.g. [77, Theorem E.3], if \mathcal{N} is a neural network with continuous, piecewise analytic activations with finitely many pieces (e.g., ReLU, sigmoid), then \mathcal{N} is sub-analytic.

Lemma 4.2 ([29]). *The rank function $\operatorname{rank}(W)$ is proper, lower semicontinuous, bounded from above and below. Furthermore, it is semi-algebraic and therefore KL.*

Proposition 4.3 (Proof in Appendix A.1). *Assume that for all n , $0 < \underline{\alpha} \leq \alpha_n \leq \frac{\sigma_n}{L_{\nabla \mathcal{L}}}$ and assumptions above hold. Then, the following are true:*

- $(\mathcal{L} + \lambda \operatorname{rank})(W_n)$ is non-increasing and convergent.
- $\sum_n \|W_{n+1} - W_n\| < +\infty$.
- W_n converges, with $\lim_{n \rightarrow \infty} W_n = W^* \in \operatorname{crit}(\mathcal{L} + \lambda \operatorname{rank})$.
- For $F_n(W) = \frac{1}{2} \|\mathcal{I}_n^{1/2} W\|_2^2$ with $\mathcal{I}_n \rightarrow \mathcal{I}$ positive definite, the smallest non-zero singular value satisfies $\sigma_{\min}(\mathcal{I}^{1/2} W^*) \geq \sqrt{\underline{\alpha} \lambda}$. And thus $\sigma_{\min}(W^*) \geq \sqrt{\sigma_{\max}(\mathcal{I}) \underline{\alpha} \lambda}$.

Proposition 4.3 establishes convergence and a maximal rank of the underlying solution. However, in practice using such a method implies computing SVD of all the weights at each iteration, making it prohibitively expensive in practice for large networks. We can instead provide a different claim for e.g. [16] in Appendix B.1.1.

5 Numerical Illustrations

While theoretical results of Section 4 highlight that the proposed iterative thresholding approach is coherent, it does not provide us with any guidelines on how to select the step-sizes α_n , or provide us with information of whether using information divergences should result in better models. For this reason, we consider what happens numerically, asking the question of what has the largest effect on compressed model accuracy. We summarize all the methods considered for comparison in Table 2, and further experimental details are provided in Appendix C.

5.1 Vision Transformer

For our initial experiments, we focused on compressing a pre-trained Vision Transformer (ViT)-B/16 (86.6M) model [20] trained on CIFAR10 [43] and one trained on ImageNet-21k [70, 17] at a resolution of 224x224. We then finetune during compression on ImageNet 2012 [72].

Table 1: Showcasing compression algorithms applied to MLP layers based on information and standard projections for the setting of Section 5.1. Using information projections is necessary for zero-shot performance, while fine-tuned performance is nearly identical for the two methods.

	% of parameters	35.7%	36.6%	37.4%	40.9%	47.8%	61.5%	Full Model
Zero-shot	SVD	0.0997	0.1005	0.1054	0.1302	0.4126	0.8601	0.9562
	FWSVD	0.1223	0.1320	0.1199	0.3383	0.6982	0.9147	
Fine-tuned	SVD	0.9237	0.9279	0.9450	0.9656	0.9585	0.9809	0.9878
	FWSVD	0.9336	0.9290	0.9541	0.9446	0.9736	0.9581	

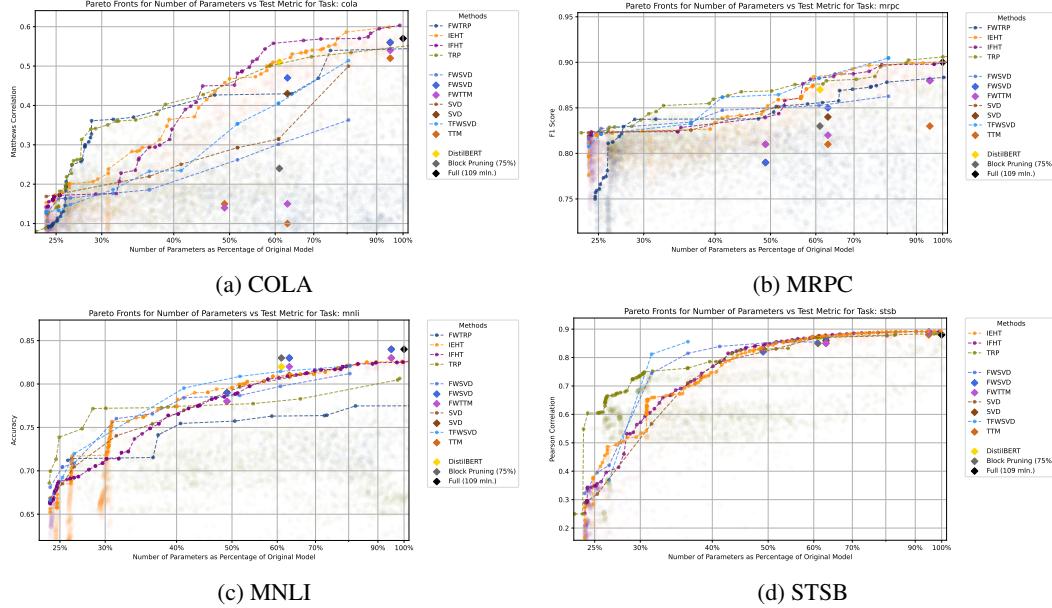


Figure 3: Performance of various compression algorithms on BERT for tasks in the GLUE benchmark. Diamonds show performance reported by [63].

Information Projection: Zero-shot vs Fine-tuned Table 1 shows the zero-shot and fine-tuned accuracy results for compressing the MLP layers of the ViT-B/16 model using fixed rank SVD and FWSVD [33], for various ranks. The results indicate that employing information-based projections (FWSVD) is crucial for preserving zero-shot performance post-compression. However, once the models are fine-tuned, the performance gap between SVD and FWSVD becomes negligible, although FWSVD in more cases than not still exhibits a slight advantage. This suggests that fine-tuning effectively mitigates the errors from Euclidean projection. We similarly observed this trend across variable rank methods, depicted in Figure 2, including the IEHT and IFHT discussed in Section 4. This finding underscores that in training information projections have an insignificant effect.

Trainability as the deciding factor for Extreme Compression Based on the discussion above, information projection has a significant effect only in zero-shot settings. Previous literature has established that initiating training with low-rank models from the outset hinders effective training [39]. We hypothesize that this is the most limiting factor for achieving good performance of compressed models. As illustrated in Figure 2, there is a significant gap between methods employing iterative compression, and those not. To highlight the role of trainability, we conducted two further ablations:

Sparsify-During-Training Methods: We employed OIALR [16], as shown in Figure 2a. While OIALR can recover improved models for the same compression ratios as constant-rank methods, it suffers from non-smooth rank reduction, often involving a large initial cut followed by several smaller ones. This limits the recovery potential of the models. To address this, we implemented the IFHT and IEHT methods from Section 4, which exhibit less severe rank cutoffs.

Different Cutoff Schedules: To further emphasize the connection between trainability and compressed performance, we performed an ablation study on various cutoff schedules for the optimal parameter settings. Three schedules were considered: (1) constant with depth, (2) increasing with depth, (3) decreasing with depth.

It is widely known [13] that earlier layers stabilize more rapidly in training than later layers. Consequently, at each hard thresholding step, earlier layers require more significant adjustments and thus necessitate greater flexibility. We anticipated that the second strategy (increasing with depth) would yield the best performance, while the third strategy (decreasing with depth) would perform the worst. Figure 4 presents the results of this ablation study, corroborating our hypothesis regarding the dominant influence of trainability on achieving optimal performance for compressed models. In analogy with [35] we also compared global vs local rank selection strategies in Appendix D.

5.2 BERT on GLUE

We follow [33, 34, 35, 63] and test on a subset of tasks in the GLUE benchmark [88] by compressing a pre-trained BERT model [19] (104M parameters) using various approaches, with a visualization in Figure 3. We argue that the setting considered here is more akin to zero-shot comparison, as performance of models either quickly diverges, or converges to a non-zero validation loss even without compression when trained for more than 10 epochs, a phenomenon well-documented in BERT finetuning [58]. This limitation in training duration significantly impacts iterative compression techniques. Specifically, the minimal difference observed between the Pareto fronts of iterative and non-iterative methods in Figure 3 arises because iterative approaches, which rely on refining the model over multiple training iterations, effectively “massaging” into a reduced form, become ineffective when prolonged training leads to a decline in test metrics. Consequently, the superior performance of non-iterative methods like FWSVD and TFWFSD in this specific context is primarily attributable to this constraint on training duration. While in scenarios where extended and stable training is feasible, iterative compression methods generally outperform their non-iterative counterparts.

6 Outlooks

This paper explored model compression via information geometry and iterative optimization. Our analysis revealed that many compression techniques can be interpreted as approximating information divergences when projecting models onto lower-cost ones, with better approximations resulting in improved zero-shot performance. Our findings suggest however that the choice of projection matters less than the optimization process itself: iterative compression during training plays the dominant role in achieving strong performance post-compression. We note however that in this paper, we have primarily focused on low rank factorization methods applied to medium-sized models, with more complex structures, and approaches overcoming the problem of trainability left for future work.

References

- [1] S.-i. Amari. Information geometry of the em and em algorithms for neural networks. *Neural networks*, 8(9):1379–1408, 1995.
- [2] S.-i. Amari. Information geometry in optimization, machine learning and statistical inference. *Frontiers of Electrical and Electronic Engineering in China*, 5:241–260, 2010.
- [3] S.-i. Amari. *Information geometry and its applications*, volume 194. Springer, 2016.
- [4] S.-i. Amari and H. Nagaoka. *Methods of information geometry*, volume 191. American Mathematical Soc., 2000.

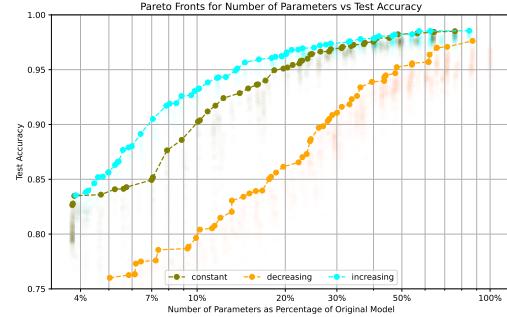


Figure 4: Ablation study of different bound schedules for CIFAR-10. We find that increasing the amount of error with depth tends to perform best due to simpler trainability.

- [5] S. Arora, N. Cohen, and E. Hazan. On the optimization of deep networks: Implicit acceleration by overparameterization. In *International conference on machine learning*, pages 244–253. PMLR, 2018.
- [6] N. Ay, J. Jost, H. Vn L, and L. Schwachhfer. *Information geometry*, volume 64. Springer, 2017.
- [7] M. M. Bejani and M. Ghatee. Adaptive low-rank factorization to regularize shallow and deep neural networks, 2020.
- [8] J. Bernstein and L. Newhouse. Old optimizer, new norm: An anthology, 2024.
- [9] R. I. Bo, E. R. Csetnek, and S. C. Lszl. An inertial forward–backward algorithm for the minimization of the sum of two nonconvex functions. *EURO Journal on Computational Optimization*, 4:3–25, 2016.
- [10] D. Busbridge, A. Shidani, F. Weers, J. Ramapuram, E. Littwin, and R. Webb. Distillation scaling laws, 2025.
- [11] E. Candes and B. Recht. Exact matrix completion via convex optimization. *Communications of the ACM*, 55(6):111–119, 2012.
- [12] P. Chen, S. Si, Y. Li, C. Chelba, and C.-J. Hsieh. Groupreduce: Block-wise low-rank approximation for neural language model shrinking. *Advances in Neural Information Processing Systems*, 31, 2018.
- [13] Y. Chen, A. Yuille, and Z. Zhou. Which layer is learning faster? a systematic exploration of layer-wise convergence rate for deep neural networks. In *The Eleventh International Conference on Learning Representations*, 2023.
- [14] E. Chunikhina, R. Raich, and T. Nguyen. Performance analysis for matrix completion via iterative hard-thresholded svd. In *2014 IEEE Workshop on Statistical Signal Processing (SSP)*, pages 392–395. IEEE, 2014.
- [15] N. Cohen, G. Menon, and Z. Verasztro. Deep linear networks for matrix completion—an infinite depth limit. *SIAM Journal on Applied Dynamical Systems*, 22(4):3208–3232, 2023.
- [16] D. Coquelin, K. Flgel, M. Weiel, N. Kiefer, C. Debus, A. Streit, and M. Gtz. Harnessing orthogonality to train low-rank neural networks. *arXiv preprint arXiv:2401.08505*, 2024.
- [17] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009.
- [18] E. L. Denton, W. Zaremba, J. Bruna, Y. LeCun, and R. Fergus. Exploiting linear structure within convolutional networks for efficient evaluation. *Advances in neural information processing systems*, 27, 2014.
- [19] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding, 2019.
- [20] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uzskoreit, and N. Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale, 2021.
- [21] J. Duchi, E. Hazan, and Y. Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 12(7), 2011.
- [22] E. Frantar, S. P. Singh, and D. Alistarh. Optimal brain compression: A framework for accurate post-training quantization and pruning, 2023.
- [23] S. Gao, T. Hua, Y.-C. Hsu, Y. Shen, and H. Jin. Adaptive rank selections for low-rank approximation of language models. In *NAAACL-HLT*, pages 227–241, 2024.

- [24] I. Garg, C. Koguchi, E. Verma, and D. Ulbricht. Revealing the utilized rank of subspaces of learning in neural networks. *CoRR*, abs/2407.04797, 2024.
- [25] V. Gupta, T. Koren, and Y. Singer. Shampoo: Preconditioned stochastic tensor optimization, 2018.
- [26] S. B. Harma, A. Chakraborty, E. Kostenok, D. Mishin, D. Ha, B. Falsafi, M. Jaggi, M. Liu, Y. Oh, S. Subramanian, and A. Yazdanbakhsh. Effective interplay between sparsity and quantization: From theory to practice, 2025.
- [27] G. Hinton, O. Vinyals, and J. Dean. Distilling the knowledge in a neural network, 2015.
- [28] J.-B. Hiriart-Urruty and H. Y. Le. From eckart and young approximation to moreau envelopes and vice versa. *RAIRO-Operations Research-Recherche Opérationnelle*, 47(3):299–310, 2013.
- [29] J.-B. Hiriart-Urruty and H. Y. Le. A variational approach of the rank function. *Top*, 21:207–240, 2013.
- [30] T. Hoefler, D. Alistarh, T. Ben-Nun, N. Dryden, and A. Peste. Sparsity in deep learning: Pruning and growth for efficient inference and training in neural networks. *Journal of Machine Learning Research*, 22(241):1–124, 2021.
- [31] S. Hofstee. Fisher information aware dynamic compression of language transformer networks using svd. Master’s thesis, University of Twente, 2024.
- [32] F. Hohman, M. B. Kery, D. Ren, and D. Moritz. Model compression in practice: Lessons learned from practitioners creating on-device machine learning experiences. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*, pages 1–18, 2024.
- [33] Y.-C. Hsu, T. Hua, S. Chang, Q. Lou, Y. Shen, and H. Jin. Language model compression with weighted low-rank factorization. In *International Conference on Learning Representations*, 2022.
- [34] T. Hua, Y.-C. Hsu, F. Wang, Q. Lou, Y. Shen, and H. Jin. Numerical optimizations for weighted low-rank estimation on language models. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 1404–1416, 2022.
- [35] T. Hua, X. Li, S. Gao, Y.-C. Hsu, Y. Shen, and H. Jin. Dynamic low-rank estimation for transformer-based language models. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 9275–9287, 2023.
- [36] M. Jaderberg, A. Vedaldi, and A. Zisserman. Speeding up convolutional neural networks with low rank expansions. *arXiv preprint arXiv:1405.3866*, 2014.
- [37] P. Jain, P. Netrapalli, and S. Sanghavi. Low-rank matrix completion using alternating minimization. In *Proceedings of the forty-fifth annual ACM symposium on Theory of computing*, pages 665–674, 2013.
- [38] X. Jia, C. Kanzow, and P. Mehlitz. Convergence analysis of the proximal gradient method in the presence of the kurdyka–łojasiewicz property without global lipschitz assumptions. *SIAM Journal on Optimization*, 33(4):3038–3056, 2023.
- [39] M. Khodak, N. A. Tenenholz, L. Mackey, and N. Fusi. Initialization and regularization of factorized neural layers. In *International Conference on Learning Representations*, 2021.
- [40] J. Kieffer. Elements of information theory (thomas m. cover and joy a. thomas). *SIAM Review*, 36(3):509–511, 1994.
- [41] B.-K. Kim, G. Kim, T.-H. Kim, T. Castells, S. Choi, J. Shin, and H.-K. Song. Shortened LLaMA: A simple depth pruning for large language models. In *ICLR 2024 Workshop on Mathematical and Empirical Understanding of Foundation Models*, 2024.
- [42] R. Krishnamoorthi. Quantizing deep convolutional networks for efficient inference: A whitepaper. *arXiv preprint arXiv:1806.08342*, 2018.

- [43] A. Krizhevsky, V. Nair, G. Hinton, et al. The cifar-10 dataset. *online: http://www.cs.toronto.edu/kriz/cifar.html*, 55(5):2, 2014.
- [44] F. Kunstner, L. Balles, and P. Hennig. Limitations of the empirical fisher approximation for natural gradient descent, 2020.
- [45] D. Kuznedelev, E. Kurtic, E. Frantar, and D. Alistarh. CAP: Correlation-aware pruning for highly-accurate sparse vision models. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [46] V. Lebedev, Y. Ganin, M. Rakhuba, I. Oseledets, and V. Lempitsky. Speeding-up convolutional neural networks using fine-tuned cp-decomposition. *arXiv preprint arXiv:1412.6553*, 2014.
- [47] Y. LeCun, J. Denker, and S. Solla. Optimal brain damage. *Advances in neural information processing systems*, 2, 1989.
- [48] J. M. Lee. *Smooth manifolds*. Springer, 2003.
- [49] J. Lin, J. Tang, H. Tang, S. Yang, W.-M. Chen, W.-C. Wang, G. Xiao, X. Dang, C. Gan, and S. Han. Awq: Activation-aware weight quantization for on-device llm compression and acceleration. *Proceedings of Machine Learning and Systems*, 6:87–100, 2024.
- [50] T. Lin, S. U. Stich, L. Barba, D. Dmitriev, and M. Jaggi. Dynamic model pruning with feedback. In *International Conference on Learning Representations*, 2020.
- [51] Y.-H. Liu, S.-W. Luo, A.-J. Li, and H.-B. Yu. Information geometry on pruning of neural network. In *Proceedings of 2004 International Conference on Machine Learning and Cybernetics (IEEE Cat. No. 04EX826)*, volume 6, pages 3479–3483. IEEE, 2004.
- [52] L. Maison, H. d. M. d. Bourboux, and T. Courtat. Compression of recurrent neural networks using matrix factorization. *arXiv preprint arXiv:2310.12688*, 2023.
- [53] J. Martens. New insights and perspectives on the natural gradient method. *Journal of Machine Learning Research*, 21(146):1–76, 2020.
- [54] J. McGowan, W. S. Lai, W. Chen, H. Aldridge, J. Clarke, J. R. Garcia, R. Xia, Y. Liang, G. Hennequin, and A. Bernacchia. Efficient model compression techniques with fishleg. In *Workshop on Machine Learning and Compression, NeurIPS 2024*, 2024.
- [55] G. Menon. The geometry of the deep linear network, 2024.
- [56] A. Mishra and D. Marr. Apprentice: Using knowledge distillation techniques to improve low-precision network accuracy. In *International Conference on Learning Representations*, 2018.
- [57] P. Molchanov, A. Mallya, S. Tyree, I. Frosio, and J. Kautz. Importance estimation for neural network pruning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11264–11272, 2019.
- [58] M. Mosbach, M. Andriushchenko, and D. Klakow. On the stability of fine-tuning {bert}: Misconceptions, explanations, and strong baselines. In *International Conference on Learning Representations*, 2021.
- [59] F. Nielsen. On geodesic triangles with right angles in a dually flat space, 2021.
- [60] M. B. Noach and Y. Goldberg. Compressing pre-trained language models by matrix decomposition. In *Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing*, pages 884–889, 2020.
- [61] I. V. Oseledets. Tensor-train decomposition. *SIAM Journal on Scientific Computing*, 33(5):2295–2317, 2011.

- [62] A.-H. Phan, K. Sobolev, K. Sozykin, D. Ermilov, J. Gusak, P. Tichavský, V. Glukhov, I. Osuledets, and A. Cichocki. Stable low-rank tensor decomposition for compression of convolutional neural network. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXIX 16*, pages 522–539. Springer, 2020.
- [63] S. Pletnev, V. Chekalina, D. Moskovskiy, M. Seleznev, S. Zagoruyko, and A. Panchenko. A computational study of matrix decomposition methods for compression of pre-trained transformers. In *Proceedings of the 37th Pacific Asia Conference on Language, Information and Computation*, pages 723–742, 2023.
- [64] B. A. Plummer, N. Dryden, J. Frost, T. Hoefler, and K. Saenko. Neural parameter allocation search. In *International Conference on Learning Representations*, 2022.
- [65] A. Polino, R. Pascanu, and D. Alistarh. Model compression via distillation and quantization. In *International Conference on Learning Representations*, 2018.
- [66] A. Potapczynski, S. Qiu, M. A. Finzi, C. Ferri, Z. Chen, M. Goldblum, C. B. Bruss, C. D. Sa, and A. G. Wilson. Searching for efficient linear layers over a continuous space of structured matrices. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- [67] D. C. Psichogios and L. H. Ungar. Svd-net: An algorithm that automatically selects network structure. *IEEE Transactions on Neural Networks*, 5(3):513–515, 1994.
- [68] S. Qiu, A. Potapczynski, M. A. Finzi, M. Goldblum, and A. G. Wilson. Compute better spent: Replacing dense layers with structured matrices. In *Forty-first International Conference on Machine Learning*, 2024.
- [69] N. Rao, P. Shah, and S. Wright. Forward–backward greedy algorithms for atomic norm regularization. *IEEE Transactions on Signal Processing*, 63(21):5798–5811, 2015.
- [70] T. Ridnik, E. Ben-Baruch, A. Noy, and L. Zelnik-Manor. Imagenet-21k pretraining for the masses, 2021.
- [71] D. Rochau, H. Gottschalk, and R. Chan. New advances in universal approximation with neural networks of minimal width. *arXiv preprint arXiv:2411.08735*, 2024.
- [72] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. Imagenet large scale visual recognition challenge, 2015.
- [73] R. Saha, N. Sagan, V. Srivastava, A. J. Goldsmith, and M. Pilanci. Compressing large language models using low rank and low precision decomposition. *arXiv preprint arXiv:2405.18886*, 2024.
- [74] T. N. Sainath, B. Kingsbury, V. Sindhwani, E. Arisoy, and B. Ramabhadran. Low-rank matrix factorization for deep neural network training with high-dimensional output targets. In *2013 IEEE international conference on acoustics, speech and signal processing*, pages 6655–6659. IEEE, 2013.
- [75] S. Schottländer, E. Zangrando, J. Kusch, G. Ceruti, and F. Tudisco. Low-rank lottery tickets: finding efficient low-rank neural networks via matrix differential equations. *Advances in Neural Information Processing Systems*, 35:20051–20063, 2022.
- [76] N. N. Schraudolph. Fast curvature matrix-vector products for second-order gradient descent. *Neural computation*, 14(7):1723–1738, 2002.
- [77] Z. Shumaylov, J. Budd, S. Mukherjee, and C.-B. Schönlieb. Weakly convex regularisers for inverse problems: Convergence of critical points and primal-dual optimisation. In *Forty-first International Conference on Machine Learning*, 2024.
- [78] S. P. Singh and D. Alistarh. Woodfisher: Efficient second-order approximation for neural network compression. *Advances in Neural Information Processing Systems*, 33:18098–18109, 2020.

- [79] N. Srebro and T. Jaakkola. Weighted low-rank approximations. In *Proceedings of the 20th international conference on machine learning (ICML-03)*, pages 720–727, 2003.
- [80] R. Sun and Z.-Q. Luo. Guaranteed matrix completion via non-convex factorization. *IEEE Transactions on Information Theory*, 62(11):6535–6579, 2016.
- [81] J. Tanner and K. Wei. Normalized iterative hard thresholding for matrix completion. *SIAM Journal on Scientific Computing*, 35(5):S104–S125, 2013.
- [82] H. Touvron, M. Cord, M. Douze, F. Massa, A. Sablayrolles, and H. Jégou. Training data-efficient image transformers and distillation through attention, 2021.
- [83] A. Tseng, J. Chee, Q. Sun, V. Kuleshov, and C. D. Sa. QuIP\$: Even better LLM quantization with hadamard incoherence and lattice codebooks. In *Forty-first International Conference on Machine Learning*, 2024.
- [84] A. Uschmajew and B. Vandereycken. *Geometric methods on low-rank matrix and tensor manifolds*. Springer, 2020.
- [85] C. Üyük, M. Lasby, M. Yassin, U. Evci, and Y. Ioannou. Learning parameter sharing with tensor decompositions and sparsity. *arXiv preprint arXiv:2411.09816*, 2024.
- [86] T. Vu, E. Chunikhina, and R. Raich. On local linear convergence rate of iterative hard thresholding for matrix completion. *IEEE Transactions on Signal Processing*, 70:5940–5953, 2022.
- [87] R. Waleffe and T. Rekatsinas. Principal component networks: Parameter reduction early in training, 2020.
- [88] A. Wang, A. Singh, J. Michael, F. Hill, O. Levy, and S. Bowman. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In T. Linzen, G. Chrupała, and A. Alishahi, editors, *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 353–355, Brussels, Belgium, Nov. 2018. Association for Computational Linguistics.
- [89] H. Wang, S. Agarwal, and D. Papailiopoulos. Pufferfish: Communication-efficient models at no extra cost. *Proceedings of Machine Learning and Systems*, 3:365–386, 2021.
- [90] X. Wang, Y. Zheng, Z. Wan, and M. Zhang. Svd-llm: Truncation-aware singular value decomposition for large language model compression. *arXiv preprint arXiv:2403.07378*, 2024.
- [91] G. I. Winata, S. Cahyawijaya, Z. Lin, Z. Liu, and P. Fung. Lightweight and efficient end-to-end speech recognition using low-rank transformer. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6144–6148. IEEE, 2020.
- [92] D. Wu, I.-V. Modoranu, M. Safaryan, D. Kuznedelev, and D. Alistarh. The iterative optimal brain surgeon: Faster sparse recovery by leveraging second-order information. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- [93] Y. Xu, Y. Li, S. Zhang, W. Wen, B. Wang, W. Dai, Y. Qi, Y. Chen, W. Lin, and H. Xiong. Trained rank pruning for efficient deep neural networks. In *2019 Fifth Workshop on Energy Efficient Machine Learning and Cognitive Computing-NeurIPS Edition (EMC2-NIPS)*, pages 14–17. IEEE, 2019.
- [94] Z. Yuan, Y. Shang, Y. Song, Q. Wu, Y. Yan, and G. Sun. Asvd: Activation-aware singular value decomposition for compressing large language models. *arXiv preprint arXiv:2312.05821*, 2023.
- [95] C. Zhang, J. Cheng, G. A. Constantinides, and Y. Zhao. LQER: Low-rank quantization error reconstruction for LLMs. In *Forty-first International Conference on Machine Learning*, 2024.
- [96] S. Zhang and V. Petyan. Low-rank is required for pruning LLMs. In *Sparsity in LLMs (SLLM): Deep Dive into Mixture of Experts, Quantization, Hardware, and Inference*, 2025.
- [97] S. Zhang and V. Petyan. OATS: Outlier-aware pruning through sparse and low rank decomposition. In *The Thirteenth International Conference on Learning Representations*, 2025.

A Appendix: Mathematical Step

Here, we recap some of the basic notions, primarily following [9]. We denote the Euclidean scalar product and norm as $\langle \cdot, \cdot \rangle$ and $\|\cdot\|$. The finite-dimensional spaces considered here are endowed with the topology induced by the Euclidean norm.

The domain of the function $f : \mathbb{R}^m \rightarrow (-\infty, +\infty]$ is defined by $\text{dom } f = \{x \in \mathbb{R}^m : f(x) < +\infty\}$. We say that f is proper if $\text{dom } f \neq \emptyset$. Consider $f : \mathbb{R}^m \rightarrow (-\infty, +\infty]$ proper and lower semicontinuous function. If $x \in \text{dom } f$, we consider the Fréchet subdifferential of f at x as the set

$$\hat{\partial}f(x) = \left\{ v \in \mathbb{R}^m : \liminf_{y \rightarrow x} \frac{f(y) - f(x) - \langle v, y - x \rangle}{\|y - x\|} \geq 0 \right\}.$$

For $x \notin \text{dom } f$ we set $\partial f(x) := \emptyset$. The limiting (Mordukhovich) subdifferential is defined at $x \in \text{dom } f$ by

$$\partial f(x) = \{v \in \mathbb{R}^m : \exists x_n \rightarrow x, f(x_n) \rightarrow f(x) \text{ and } \exists v_n \in \partial f(x_n), v_n \rightarrow v \text{ as } n \rightarrow +\infty\}$$

while for $x \notin \text{dom } f$, one takes $\partial f(x) := \emptyset$. For f convex, both coincide with the convex subdifferential $\partial f(x) = \{v \in \mathbb{R}^m : f(y) \geq f(x) + \langle v, y - x \rangle \forall y \in \mathbb{R}^m\}$ for all $x \in \text{dom } f$.

If $x \in \mathbb{R}^m$ is a local minimizer of f , then $0 \in \partial f(x)$. Notice that in case f is continuously differentiable around $x \in \mathbb{R}^m$ we have $\partial f(x) = \{\nabla f(x)\}$. Let us denote by

$$\text{crit}(f) = \{x \in \mathbb{R}^m : 0 \in \partial f(x)\}$$

the set of (limiting)-critical points of f . We now define functions satisfying the Kurdyka-Łojasiewicz property, playing a crucial role in proving the convergence of algorithms in the nonconvex setting. For $\eta \in (0, +\infty]$, we denote by Θ_η the class of concave and continuous functions $\varphi : [0, \eta] \rightarrow [0, +\infty)$ such that $\varphi(0) = 0$, φ is continuously differentiable on $(0, \eta)$, continuous at 0 and $\varphi'(s) > 0$ for all $s \in (0, \eta)$. We denote the distance function to a set, defined for $A \subseteq \mathbb{R}^m$ as $\text{dist}(x, A) = \inf_{y \in A} \|x - y\|$ for all $x \in \mathbb{R}^m$.

Definition A.1 (Kurdyka-Łojasiewicz property). Let $f : \mathbb{R}^m \rightarrow (-\infty, +\infty]$ be a proper and lower semicontinuous function. We say that f satisfies the Kurdyka-Łojasiewicz (KL) property at $\bar{x} \in \text{dom } \partial f = \{x \in \mathbb{R}^m : \partial f(x) \neq \emptyset\}$ if there exists $\eta \in (0, +\infty]$, a neighborhood U of \bar{x} and a function $\varphi \in \Theta_\eta$ such that for all x in the intersection

$$U \cap \{x \in \mathbb{R}^m : f(\bar{x}) < f(x) < f(\bar{x}) + \eta\}$$

the following inequality holds

$$\varphi'(f(x) - f(\bar{x})) \text{dist}(0, \partial f(x)) \geq 1$$

If f satisfies the KL property at each point in $\text{dom } \partial f$, then f is called a KL function.

To the class of KL functions belong semi-algebraic, real sub-analytic, semiconvex, uniformly convex and convex functions satisfying a growth condition. See [9] for a full discussion. The following lemma turns out to be useful characterization of the property:

Lemma A.1. *Let $\Omega \subseteq \mathbb{R}^m$ be a compact set and let $f : \mathbb{R}^m \rightarrow (-\infty, +\infty]$ be a proper and lower semicontinuous function. Assume that f is constant on Ω and f satisfies the KL property at each point of Ω . Then there exist $\varepsilon, \eta > 0$ and $\varphi \in \Theta_\eta$ such that for all $\bar{x} \in \Omega$ and for all x in the intersection*

$$\{x \in \mathbb{R}^m : \text{dist}(x, \Omega) < \varepsilon\} \cap \{x \in \mathbb{R}^m : f(\bar{x}) < f(x) < f(\bar{x}) + \eta\}$$

the following inequality holds

$$\varphi'(f(x) - f(\bar{x})) \text{dist}(0, \partial f(x)) \geq 1$$

A.1 Proof of Proposition 4.3

Proposition A.2. *Assume that for all n , $0 < \underline{\alpha} \leq \alpha_n \leq \frac{\sigma_n}{L_{\nabla \mathcal{L}}}$ and assumptions above hold. Then, the following are true:*

- $(\mathcal{L} + \lambda \text{rank}) (W_n)$ is non-increasing and convergent.
- $\sum_n \|W_{n+1} - W_n\| < +\infty$
- W_n converges, with $\lim_{n \rightarrow \infty} W_n = W^* \in \text{crit}(\mathcal{L} + \lambda \text{rank})$.
- For $F_n(W) = \frac{1}{2} \|\mathcal{I}_n^{1/2} W\|_2^2$ with $\mathcal{I}_n \rightarrow \mathcal{I}$ positive definite, the smallest non-zero singular value satisfies $\sigma_{\min}(\mathcal{I}^{1/2} W^*) \geq \sqrt{\underline{\alpha}\lambda}$. And thus $\sigma_{\min}(W^*) \geq \sqrt{\sigma_{\max}(\mathcal{I}) \underline{\alpha}\lambda}$

Proof. Using Lemma 4.2, rank is semi-algebraic and therefore sub-analytic. Thus $(\mathcal{L} + \lambda \text{rank})$ is sub-analytic and therefore KŁ. As rank is not continuous, the overall function may not admit a KŁ exponent. For a constant F_n this would be a direct application of theory in [9, Theorem 12] with $\beta = 0$. In general however, the main descent result is as follows:

$$D_{F_n}(W_{n+1}, W_n) + \langle W_{n+1} - W_n, \alpha_n \nabla \mathcal{L}(W_n) \rangle + \alpha_n \lambda \text{rank}(W_{n+1}) \leq \alpha_n \lambda \text{rank}(W_n).$$

On the other hand, by descent lemma we have

$$\langle \nabla \mathcal{L}(W_n), W_{n+1} - W_n \rangle \geq \mathcal{L}(W_{n+1}) - \mathcal{L}(W_n) - \frac{L_{\nabla \mathcal{L}}}{2} \|W_n - W_{n+1}\|^2.$$

At the same time for any $\mu > 0$,

$$\langle W_{n+1} - W_n, W_{n-1} - W_n \rangle \geq - \left(\frac{\mu}{2} \|W_n - W_{n+1}\|^2 + \frac{1}{2\mu} \|W_{n-1} - W_n\|^2 \right)$$

while by assumption

$$\frac{\sigma_n}{2} \|W_{n+1} - W_n\|^2 \leq D_{F_n}(W_{n+1}, W_n).$$

This implies that above becomes

$$(\lambda \text{rank} + \mathcal{L})(W_{n+1}) + \frac{\sigma_n - L_{\nabla \mathcal{L}} \alpha_n}{2\alpha_n} \|W_{n+1} - W_n\|^2 \leq (\lambda \text{rank} + \mathcal{L})(W_n).$$

By the step size assumption, we have $\sigma_n \geq \underline{\alpha}/L_{\nabla \mathcal{L}}$. By the lipshitzness assumption, we must further have $\sigma_n \leq \overline{L}_{\nabla F}$. Thanks to these bounds, the rest of the results from [9] apply directly.

For the last point, this is a direct consequence of [28, Theorem 3]. To be precise, we must have that $\sigma_{\min}(\mathcal{I}_n^{1/2} W_{n+1}) \geq \sqrt{\alpha_n \lambda}$. By continuity of singular values, using e.g. Weyl's inequality, we have the desired result as then the minimum non-zero singular value is continuous outside of zero. Lastly, $\sigma_{\min}(W^*) = \sigma_{\min}(\mathcal{I}^{-1/2} \mathcal{I}^{1/2} W^*) \geq \sigma_{\min}(\mathcal{I}^{-1/2}) \sigma_{\min}(\mathcal{I}^{1/2} W^*)$ \square

A.2 Proof of Lemma A.3

Lemma A.3. *Let $f : X \rightarrow Y$ a submersion between two manifolds, and let $M \in Y$ a subset of Y . Then, M is an embedded submanifold of Y if and only if $f^{-1}(M)$ is an embedded submanifold of X .*

Proof. One way comes directly from [48, Corollary 6.31]. For the other direction, we can use the constant rank theorem, which gives us charts $U \subseteq X \rightarrow V \subseteq Y$, such that $f : U \rightarrow V$ is $(x_1, \dots, x_m) \rightarrow (x_1, \dots, x_n)$, where $m = \dim X, n = \dim Y$ with $m \geq n$. Then $f^{-1}(M) \subseteq X$, so $f^{-1}(M) \cap U \subseteq V$ embedded. Then find slice chart for $f^{-1}(M) \cap U$, which gives us a slice chart for $M \cap V$. \square

B Methods Overview

B.1 OIALR

[16] observed that the orthogonal basis of a network's weights stabilizes during the training process. Based on this, they propose to reduce the number of trainable parameters by iteratively orthogonalizing as follows:

1. The network is initially trained with a traditional full-rank scheme.

Table 2: Summary of Compression Methods

Projection	Method	Origin	Iterative	Rank Selection Criterion
Euclidean	OIALR	[16]	Yes	Layer Maximal Singular Value
	IEHT	(Ours Section 4 and Appendix B.2)	Yes	Layer Energy
	TRP	[93]	Yes	Layer Energy
	globalIEHT	(Ours Section 4 and Appendix B.2)	Yes	Total Energy
	SVD	[67]	No	Fixed
Information	TTM	[61]	No	Fixed
	IFHT	(Ours Section 4 and Appendix B.2)	Yes	Layer 'Fisher' Energy
	FWTRP	(Ours Appendix B.3)	Yes	Layer 'Fisher' Energy
	globalIFHT	(Ours Section 4 and Appendix B.2)	Yes	Global 'Fisher' Energy
	FWSVD	[33]	No	Fixed
	TFWSVD	[34]	No	Fixed
	FWTTM	[63]	No	Fixed

2. After a number of iterations, the network's matrix weights are transitioned to their $U\Sigma V^T$ representation using singular value decomposition (SVD).
3. The orthogonal bases U and V^T are frozen, and only the square matrix Σ is trained using backpropagation.
4. After a specified number of training steps, the bases U and V^T are updated by extracting the new bases from the trained Σ matrix using SVD.
5. A new inner rank is found by removing the singular values from Σ whose absolute magnitude is less than β times the largest singular value in the current Σ . β is a hyperparameter that defaults to 0.1.
6. This process is repeated until the end of training.

We summarize the method in Algorithm 1, and attempt to formalize the approach below.

B.1.1 Orthogonality Informed Low Rank Training

Proposition 4.3 establishes convergence and a maximal rank of the underlying solution. However, in practice using such a method implies computing SVD of all the weights at each iteration, making it prohibitively expensive. Below we instead consider the specific algorithm of [16]. For $U^0 = V^0 = \text{Id}$ and $S^0 = W^0$, we can summarize it as:

$$\begin{aligned} S_*^{k+1} &\in \underset{S}{\text{crit}} \mathcal{L}(U_k S V_k^\top) \\ S^{k+1} &\in \text{Prox}_{\lambda_k}(\text{rank})(S_*^{k+1}) \\ U_{k+1}, V_{k+1}^\top &\in \text{SVD}(U_k S^{k+1} V_k^\top) \\ W^{k+1} &= U_k S^{k+1} V_k \end{aligned}$$

This effectively means that we iteratively project via the prox operator, followed by finding a critical point on the low-rank submanifold. Note however, success of this method relies on the fact that during training the basis elements stabilize quite quickly. This allows us to choose a local chart on the low-rank submanifold on which we can perform the desired projection.

In general it would be rather difficult to say anything about convergence of such an algorithm, as even in the convex case, given a small misalignment of the linear chart U_k, V_k with the chart at the global minimum, the global minimum becomes unattainable.

However, we can instead control directly the loss value jump under projection. So, say that instead of $S_*^{k+1} \in \underset{S}{\text{crit}} \mathcal{L}(U_k S V_k^\top)$, we perform only a single step of gradient descent. In this case we have that $S_*^{k+1} = S^k - \frac{1}{L_{\nabla \mathcal{L}}} \nabla_S \mathcal{L}(U_k S^k V_k^\top)$. However we know that $\nabla_S \mathcal{L}(U_k S V_k^\top) = U_k^\top \nabla_W \mathcal{L}(U_k S V_k^\top) V_k$. By descent lemma we have that:

$$\mathcal{L}(U_k S_*^{k+1} V_k^\top) - \mathcal{L}(U_k S^k V_k^\top) \leq -\frac{L_{\nabla \mathcal{L}}}{2} \|S_*^{k+1} - S^k\|^2 = -\frac{1}{2L_{\nabla \mathcal{L}}} \|\nabla_S \mathcal{L}(U_k S^k V_k^\top)\|^2$$

Thus, gradient descent guarantees decrease of the objective function. What about the projection step? Assuming now again that $S_*^{k+1} \in \underset{S}{\text{crit}} \mathcal{L}(U_k S V_k^\top)$, we know that $\nabla_S \mathcal{L}(U_k S_*^{k+1} V_k^\top) = 0$. This

however, does not guarantee criticality in the full space $\nabla_W \mathcal{L}(U_k S_*^{k+1} V_k^\top) \neq 0$. Despite this, the descent lemma still provides us with a bound on the total decrease:

$$\begin{aligned}
& \mathcal{L}(U_k S^{k+1} V_k^\top) - \mathcal{L}(U_k S_*^{k+1} V_k^\top) \\
& \leq \langle \nabla_W \mathcal{L}(U_k S_*^{k+1} V_k^\top), U_k S_*^{k+1} V_k^\top - U_k S^{k+1} V_k^\top \rangle \\
& \quad + \frac{L_{\nabla \mathcal{L}}}{2} \|U_k S_*^{k+1} V_k^\top - U_k S^{k+1} V_k^\top\|^2 \\
& = \langle U_k^\top \nabla_W \mathcal{L}(U_k S_*^{k+1} V_k^\top) V_k, S_*^{k+1} - S^{k+1} \rangle + \frac{L_{\nabla \mathcal{L}}}{2} \|U_k S_*^{k+1} V_k^\top - U_k S^{k+1} V_k^\top\|^2 \\
& = \langle \nabla_S \mathcal{L}(U_k S_*^{k+1} V_k^\top), S_*^{k+1} - S^{k+1} \rangle + \frac{L_{\nabla \mathcal{L}}}{2} \|S_*^{k+1} - S^{k+1}\|^2 = \\
& = \frac{L_{\nabla \mathcal{L}}}{2} \|S_*^{k+1} - S^{k+1}\|^2 \\
& \leq \lambda_k L_{\nabla \mathcal{L}} \{\text{rank}(S_*^{k+1}) - \text{rank}(S^{k+1})\} \\
& < \sigma_{\min}^2(S^{k+1}) L_{\nabla \mathcal{L}} \{\text{rank}(S_*^{k+1}) - \text{rank}(S^{k+1})\}
\end{aligned}$$

Thus, as long as $\|S_*^{k+1} - S^{k+1}\|^2$ is bounded, the objective increase is not large. This is directly controlled through λ_k , as by definition of prox, we have $\frac{1}{2\lambda_k} \|S^{k+1} - S_*^{k+1}\|^2 \leq \text{rank}(S_*^{k+1}) - \text{rank}(S^{k+1})$. We also know that λ_k controls the largest singular value cutoff, and thus we have that $\lambda_k < \sigma_{\min}^2(S^{k+1})$. Although it is worth emphasizing that normally in the projection step we have explicit control on the size $\|S_*^{k+1} - S^{k+1}\|^2$.

B.2 IEHT and IFHT

Based on the observations in Section 5.1, and section above we realize that the cutoff selection proposed by [16] is detrimental to downstream performance, especially given that compatibility the initial choice of chart is fundamental for establishing convergence. We instead propose to perform cutoffs based on the total energy (sum of squares of singular values), which turns out to be much milder.

Algorithm 1: OIALR, IEHT

```

Data: Model  $M$ 
parameters: Training steps  $t_{\max}$ , delay steps  $d$ , low-rank update frequency  $\nu$ , singular value
cutoff fraction  $\beta$ 
Result: Compressed trained model.
for  $t = 0 \dots t_{\max}$  do
  if  $t < d$  then
    | Train full-rank ;
  else if  $t = d$  then
    | Convert network to  $USV^\top$  representation ;
    | Freeze  $U$  and  $V$  ;
  else if  $t \equiv 0 \pmod{\nu}$  then
    | for  $i = 0 \dots L$  do
      | |  $U'_i, S'_i, V'_i^\top \leftarrow \text{svd}(S_i)$  ;
      | |  $U_i \leftarrow U_i U'_i$  ;
      | |  $V_i \leftarrow V_i V'_i$  ;
      | |  $S_i \leftarrow S'_i$  ;
      | | Remove singular values  $< \beta \cdot \max(S_i)$  ;
      | | For IEHT: Remove lower singular values contributing  $< \beta\%$  of the total energy;
      | | Reshape  $U_i, V_i, S_i$ , and optimizer states;
    end
  else
    | Train network in  $USV^\top$  representation, with  $U$  and  $V$  frozen;
  end
end
# Now compile into smaller model
for  $i = 0 \dots L$  do
  |  $U'_i, S'_i, V'_i^\top \leftarrow \text{svd}(S_i)$  ;
  |  $U_i \leftarrow U_i U'_i \sqrt{S'_i}$  ;
  |  $V_i \leftarrow V_i V'_i \sqrt{S'_i}$  ;
  | Remove  $S_i$  as parameters.
end
return Trained factorized model  $M$ 

```

Algorithm 2: IFHT

```

Data: Model  $M$ 
parameters: Training steps  $t_{\max}$ , delay steps  $d$ , low-rank update frequency  $\nu$ , energy cutoff
fraction  $\beta$ , dataset for fisher estimation  $D$ 
Result: Compressed trained model.
for  $t = 0 \dots t_{\max}$  do
  if  $t < d$  then
    | Train full-rank ;
  else if  $t = d$  then
    | Convert network to  $USV^\top$  representation ;
    | Freeze  $U$  and  $V$  ;
  else if  $t \equiv 0 \pmod{\nu}$  then
    |  $\mathcal{T} \leftarrow \text{fisher}(M, D)$ ;
    | for  $i = 0 \dots L$  do
      | |  $U'_i, S'_i, V'_i^\top \leftarrow \text{svd}(\tilde{\mathcal{T}} S_i)$  ;
      | |  $U_i \leftarrow U_i \tilde{\mathcal{T}}^{-1} U'_i$  ;
      | |  $V_i \leftarrow V_i V'_i$  ;
      | |  $S_i \leftarrow S'_i$  ;
      | | Remove lower singular values contributing less than  $< \beta\%$  of the total energy;
      | | Re-factorize to keep  $U_i, V_i$  semi-orthogonal;
      | | Reshape  $U_i, V_i, S_i$ , and optimizer states;
    end
  else
    | Train network in  $USV^\top$  representation, with  $U$  and  $V$  frozen;
  end
end
# Now compile into smaller model
for  $i = 0 \dots L$  do
  |  $U'_i, S'_i, V'_i^\top \leftarrow \text{svd}(S_i)$  ;
  |  $U_i \leftarrow U_i U'_i \sqrt{S'_i}$  ;
  |  $V_i \leftarrow V_i V'_i \sqrt{S'_i}$  ;
  | Remove  $S_i$  as parameters.
end
return Trained factorized model  $M$ 

```

B.3 TRP and FWTRP

Description of TRP: In trained rank pruning, every `trp_frequency` steps, the networks singular values are hard thresholded to 0, with nuclear norm gradients saved, and subsequently each `nuclear_norm_frequency`, gradient steps are performed.

Description of FWTRP: The overall description is the same as that of TRP, but with SVD now being performed with respect to approximated fisher information.

Algorithm 3: TRP (Trained Rank Pruning)

Data: Model M
parameters: Training steps t_{\max} , low-rank update frequency trp_frequency , singular value cutoff fraction β

Result: Compressed trained model.

```

for  $t = 0 \dots t_{\max}$  do
    if  $t \bmod \text{trp\_frequency} \neq 0$  then
        | Train full-rank ;
    else
        | for  $i = 0 \dots L$  do
            |    $U_i, S_i, V_i^T \leftarrow \text{svd}(W_i)$  ;
            |   Remove lower singular values contributing less than  $< \beta\%$  of the total energy at  $k$ :
            |    $W_i \leftarrow U_i[:, :k]S_i[k, :k]V_i[:, :k]^T$  ;
            |    $G_i \leftarrow U_i[:, :k]V_i[:, :k]^T$ 
        end
        if  $t \bmod \text{nuclear\_norm\_frequency} = 0$  then
            | Make gradient step with  $\text{nuclear\_norm\_weight} * G_i$ 
        end
    end
    # Now compile into smaller model
    for  $i = 0 \dots L$  do
        |    $U_i, S_i, V_i^T \leftarrow \text{svd}(W_i)$  ;
        |    $U_i \leftarrow U_i\sqrt{S_i}$  ;
        |    $V_i \leftarrow V_i\sqrt{S_i}$  ;
        |   Remove  $S_i$  as parameters.
    end
return Trained factorized model  $M$ 

```

Algorithm 4: FWTRP (Trained Rank Pruning)

Data: Model M
parameters: Training steps t_{\max} , low-rank update frequency trp_frequency , regularization frequency $\text{nuclear_norm_frequency}$, singular value cutoff fraction β

Result: Compressed trained model.

```

for  $t = 0 \dots t_{\max}$  do
    if  $t \bmod \text{trp\_frequency} \neq 0$  then
        | Train full-rank ;
    else
        |  $\tilde{T} = \text{fisher}(M, D)$ :
        | for  $i = 0 \dots L$  do
            |    $U_i, S_i, V_i^T \leftarrow \text{svd}(\tilde{T}S_i)$ :
            |   Remove lower singular values contributing less than  $< \beta\%$  of the total energy at  $k$ :
            |    $W_i \leftarrow \tilde{T}^{-1}U_i[:, :k]S_i[k, :k]V_i[:, :k]^T$  ;
            |    $G_i \leftarrow U_i[:, :k]V_i[:, :k]^T$ 
        end
        if  $t \bmod \text{nuclear\_norm\_frequency} = 0$  then
            | Make gradient step with  $\text{nuclear\_norm\_weight} * G_i$ 
        end
    end
    # Now compile into smaller model
    for  $i = 0 \dots L$  do
        |    $U_i, S_i, V_i^T \leftarrow \text{svd}(W_i)$  ;
        |    $U_i \leftarrow U_i\sqrt{S_i}$  ;
        |    $V_i \leftarrow V_i\sqrt{S_i}$  ;
        |   Remove  $S_i$  as parameters.
    end
return Trained factorized model  $M$ 

```

C Further Experimental Details

For CIFAR10, the models were trained for 250 epochs, while for ImageNet for 50. For BERT, models were trained for 100 epochs, even though convergence was observed very quickly. Pre-trained models were acquired from the `timm` and `huggingface` libraries and utilized as the baseline, with various compression techniques applied.

In order to create the visualizations on Figures 2a, 2b and 3, all the algorithms of Appendix B were run on various parameter settings. For each run, test accuracy was evaluated at the end of each epoch and plotted on the figure. Then, the resulting pareto front was plotted to showcase best possible performance amongst all the parameters. The parameter values for all methods are shown in Tables 3 to 5. Each experimental setting was run on a single 24GB GPU for up to 1 week of execution.

Table 3: Parameter variations for ViT on CIFAR10

Method	Parameter	Explored Values
SVD, FWSVD, TFWSVD	rank	8, 16, 24, 32, 64, 128, 256
TRP, FWTRP	trp_threshold	0.9, 0.95, 0.99
	trp_frequency	50, 100, 500, 1000
	nuclear_norm_weight	0.0003, 0.00005, 0.005
	nuclear_norm_frequency	trp_frequency // 2
OIALR	oialr_threshold	0.25, 0.2, 0.15
	oialr_frequency	10, 15, 25
	oialr_depth_schedule	'constant'
	oialr_type	'epoch'
	oialr_min_rank_percent	0.02
IEHT, IFHT	oialr_threshold	0.9, 0.925, 0.95, 0.97
	oialr_frequency	10, 15, 25
	oialr_depth_schedule	'constant', 'increasing'
	oialr_type	'epoch'
	oialr_min_rank_percent	0.02
Global IHT	oialr_threshold	0.8, 0.9, 0.95, 0.99
	oialr_frequency	10, 25, 50
	oialr_type	'epoch'
Base Training	epochs	250
	learning_rate	0.0001
	batch_size	64
	optimizer	AdamW

Table 4: Parameter variations for ViT on ImageNet

Method	Parameter	Explored Values
SVD, FWSVD, TFWSVD	rank	24, 32, 64, 128, 256, 384, 512
TRP, FWTRP	trp_threshold	0.9, 0.95, 0.97, 0.99
	trp_frequency	5000, 10000, 50000
	nuclear_norm_weight	0.003, 0.05
	nuclear_norm_frequency	trp_frequency // 2
IEHT, IFHT	oialr_threshold	0.9, 0.925, 0.95, 0.97
	oialr_frequency	2, 4, 8, 10
	oialr_depth_schedule	'constant', 'increasing'
	oialr_type	'epoch'
	weight_decay	5e-4
	oialr_min_rank_percent	0.02
Base Training	epochs	50
	learning_rate	0.0001
	batch_size	128
	optimizer	AdamW

Table 5: Parameter variations for BERT on GLUE

Method	Parameter	Explored Values
SVD, FWSVD, TFWSVD	rank	16, 24, 32, 64, 96, 128, 196, 256, 384
IEHT, IFHT	oialr_threshold	0.95, 0.97, 0.98, 0.99
	oialr_frequency	10, 25, 50, 100, 500, 1000
	oialr_depth_schedule	'increasing'
	oialr_type	'step'
	oialr_min_rank_percent	0.02
TRP, FWTRP	trp_threshold	0.95, 0.97, 0.98, 0.99
	trp_frequency	10, 50, 100, 500
	nuclear_norm_weight	0.0003, 0.005
	nuclear_norm_frequency	trp_frequency // 2
Base Training	epochs	100
	learning_rate	0.00002
	batch_size	32
	optimizer	AdamW

D Global vs Local

Given the results of [35], we wanted to see whether a move towards global rank selection in IFHT would result in improved performance in the setting of CIFAR10 with vision transformers. We illustrate the joint plot in Figure 5. For the sake of consistency, we illustrate a global modification of IEHT as well. As expected, the globalIEHT method does not perform well, as the singular values differ significantly amongst the layers, resulting in a global selection that is likely to create bottlenecks very early on. However, unlike [35], globalIFHT also underperforms, though not significantly so, when compared to the purely local IFHT. We hypothesize that the reason for this is analogous to the observation that BERT finetuning is more analogous to zero-shot performance comparison, as finetuning tends to not generalize, as discussed in Section 5.2.

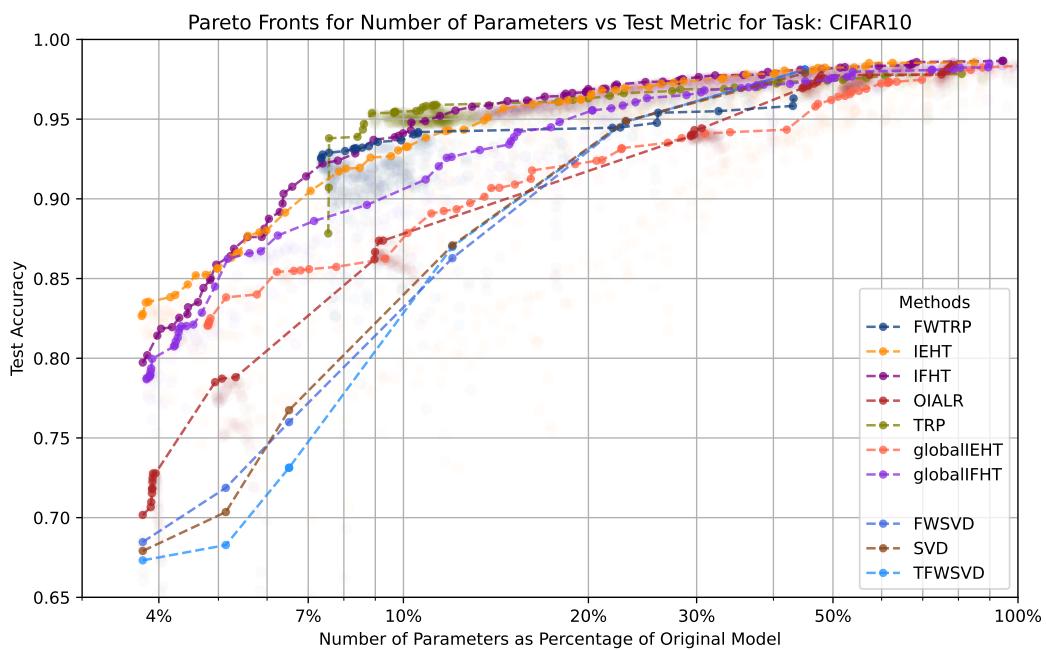


Figure 5: CIFAR-10 classification, illustrating also the globalIEHT and globalIFHT methods.