
Maestro: Uncovering Low-Rank Structures via Trainable Decomposition

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Deep Neural Networks (DNNs) have been a large driver and enabler for AI break-
2 throughs in recent years. These models have been getting larger in their attempt to
3 become more accurate and tackle new upcoming use-cases, including AR/VR and
4 intelligent assistants. However, the training process of such large models is a costly
5 and time-consuming process, which typically yields a single model to fit all targets.
6 To mitigate this, various techniques have been proposed in the literature, including
7 pruning, sparsification or quantization of the model weights and updates. While
8 able to achieve high compression rates, they often incur computational overheads
9 or accuracy penalties. Alternatively, factorization methods have been leveraged
10 to incorporate low-rank compression in the training process. Similarly, such tech-
11 niques (e.g., SVD) frequently rely on the computationally expensive decomposition
12 of layers and are potentially sub-optimal for non-linear models, such as DNNs.

13 In this work, we take a further step in designing efficient low-rank models and
14 propose MAESTRO, a framework for trainable low-rank layers. Instead of regularly
15 applying a priori decompositions such as SVD, the low-rank structure is built
16 into the training process through a generalized variant of Ordered Dropout. This
17 method imposes an importance ordering via sampling on the decomposed DNN
18 structure. Our theoretical analysis demonstrates that our method recovers the SVD
19 decomposition of linear mapping on uniformly distributed data and PCA for linear
20 autoencoders. We further apply our technique on DNNs and empirically illustrate
21 that MAESTRO enables the extraction of lower footprint models that preserve
22 model performance while allowing for graceful accuracy-latency tradeoff for the
23 deployment to devices of different capabilities.

24 1 Introduction

25 Deep Learning has been experiencing an unprecedented uptake, with models achieving a
26 (super-)human level of performance in several tasks across modalities, giving birth to even more in-
27 telligent assistants and next-gen visual perception and generation systems. However, the price of this
28 performance is that models are getting significantly larger, with training and deployment becoming
29 increasingly costly. Therefore, techniques from Efficient ML become evermore relevant [27], and a
30 requirement for deployment in constrained devices, such as smartphones or IoT devices.

31 Typical techniques to compress the network involve *i) quantization*, i.e., reducing precision of the
32 model [52] or communicated updates [45, 2], *ii) pruning* the model during training, e.g., through
33 Lottery Ticket Hypothesis (LTH) [11], *iii) sparsification* of the network representation and updates,
34 i.e., dropping the subset of coordinates [49, 3] or *iv) low-rank approximation* [53, 9], i.e. keeping the
35 most relevant ranks of the decomposed network. Despite the benefits during deployment, that is a
36 lower footprint model, in many cases, the overhead during training time or the accuracy degradation

37 can be non-negligible. Moreover, many techniques can introduce multiple hyperparameters or the
 38 need to fine-tune to recover the lost accuracy.

39 In this work, we focus on training low-rank factorized models. Specifically, we pinpoint the challenges
 40 of techniques [53, 54] when decomposing the parameters of each layer in low-rank space and the need
 41 to find the optimal ranks for each one at training time. To solve this, we adopt and non-trivially extend
 42 the Ordered Dropout technique from [17] and apply it to find progressively the optimal decomposition
 43 for each layer of a network while training (Fig. 1). Critical differences to prior work include *i*) the
 44 non-uniformity of the search space (i.e. we allow for different ranks per layer), *ii*) the trainable aspect
 45 of the decomposition to reflect the data distribution, and *iii*) the gains to training and deployment time
 46 without sacrificing accuracy. Nevertheless, we also provide a latency-accuracy trade-off mechanism
 47 to deploy the network on even more constrained devices.

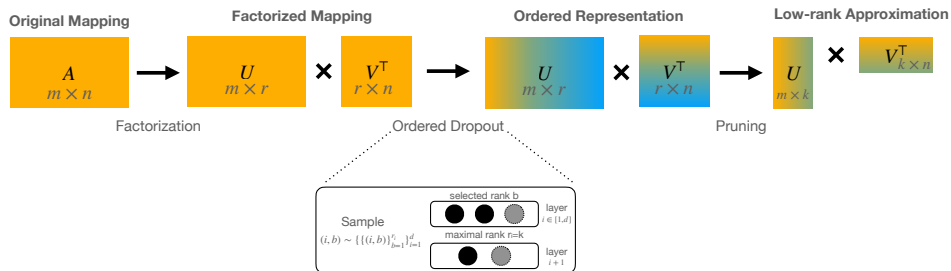


Figure 1: MAESTRO’s construction. To obtain low-rank approximation, the given linear map is decomposed and trained with ordered dropout to obtain an ordered representation that can be efficiently pruned.

48 Our contributions can be summarized as follows:

- 49 • We propose MAESTRO, a novel layer decomposition technique that enables learning low-rank layers
 50 in a progressive manner while training. We novelly fuse layer factorization and an extended variant
 51 of the ordered dropout, by embedding OD directly into the factorized weights. By decomposing
 52 layers and training on stochastically sampled low-rank models, we apply ordered importance
 53 decomposed representation of each layer. We combine this with a *hierarchical group-lasso*
 54 term [64] in the loss function to zero out redundant ranks and *progressively shrink* the rank space.
 55 This way, we enable computationally efficient training achieved by the proposed decomposition
 56 without relying on inexact and potentially computationally expensive decompositions such as SVD.
- 57 • MAESTRO is a theoretically motivated approach that embeds decomposition into training. First,
 58 we show that our new objective is able to recover *i*) the SVD of the target linear mapping for the
 59 particular case of uniform data distribution and *ii*) the Principal Component Analysis (PCA) of the
 60 data in the case of identity mapping.
- 61 • As MAESTRO’s decomposition is part of the training procedure, it also accounts for data distribution
 62 and the target function, contrary to SVD, which operates directly on learned weights. We show
 63 that this problem *already arises* for a simple linear model and empirically generalize our results in
 64 the case of DNNs, by applying our method to different types of layers (including fully-connected,
 65 convolutional, and attention) spanning across three datasets and modalities. We illustrate that our
 66 technique achieves better results than SVD-based baselines at a lower cost.

67 2 Related work

68 The topic of Efficient ML has received a lot of attention throughout the past decade as networks
 69 have been getting increasingly computationally expensive. Towards this end, we distinguish between
 70 training and deployment time, with the latter having a more significant impact and thus amortizes
 71 the potential overhead during training. Nevertheless, with the advent of Federated Learning [36],
 72 efficient training becomes increasingly relevant to remain tractable.

73 **Efficient inference.** For efficient deployment, there have been proposed various techniques that either
 74 optimize the architecture of the DNN in a hand-crafted [19] or automated manner (i.e. NAS) [50],
 75 they remove redundant computation by means of pruning parts of the network [12, 6, 11, 48, 30, 55,
 76 21, 55, 65, 15, 59, 33, 62] or utilise low-precision representation [52] of the neurons and activations.
 77 Closer to our method, there have been techniques leveraging low-rank approximation (e.g. SVD)
 78 for efficient inference [58, 43, 22, 56, 9]. Last, there is a category of techniques that dynamically

79 resize the network at runtime for compute, memory or energy efficiency, based on early-exiting [26]
 80 or dynamic-width [63] and leverage the accuracy-latency tradeoff.

81 **Efficient training.** On the other hand, techniques for efficient training become very relevant nowadays
 82 when scaling DNNs sizes [20] or deploying to embedded devices [32], and oftentimes offer additional
 83 gains at deployment time. Towards this goal, there have been employed methods where part of
 84 the network is masked [46] or dropped [1, 5] during training, with the goal of minimizing the
 85 training footprint. Similarly to early-exiting, there have been proposed multi-exit variants for efficient
 86 training [24, 34], and the same applies for width-based scaling [17, 8]. Last but not least, in the era of
 87 transformers and LLMs, where networks have scaled exponentially in size, PEFT-based techniques,
 88 such as adapter-based fine-tuning [18] (such as LoRA [20]), become increasingly important and make
 89 an important differentiator for tackling downstream tasks.

90 **Learning ordered representation.** Originally, Ordered Dropout (OD) was proposed as a mechanism
 91 for importance-based pruning for the easy extraction of sub-networks devised to allow for heteroge-
 92 neous federated training [17]. The earlier work that aims to learn ordered representation includes a
 93 similar technique to OD—Nested Dropout, which proposed a similar construction, applied to the
 94 representation layer in autoencoders [42] to enforce identifiability of the learned representation or
 95 the last layer of the feature extractor [16] to learn an ordered set of features for transfer learning.
 96 We leverage and non-trivially extend OD in our technique as a means to order ranks in terms of
 97 importance in a nested manner during training of a decomposed network that is progressively shrunk
 98 as redundant ranks converge to 0. Ranks selection is ensured through hierarchical group lasso penalty,
 99 as described in Sec. 3.3. Moreover, contrary to [17], which assumed a uniform width, our formulation
 100 allows for heterogeneous ranks per layer. Last, we leverage the ordered representation of ranks at
 101 inference time to further compress the model, allowing a graceful degradation of performance as a
 102 mechanism for the accuracy-latency trade-off.

103 3 MAESTRO

104 In this work, we focus on low-rank models as a technique to reduce the computational complexity and
 105 memory requirements of the neural network model. The main challenge that we face is the selection
 106 of the optimal rank or the trade-off between the efficiency and the rank for the given layer represented
 107 by linear mapping. Therefore, we devise an importance-based training technique, MAESTRO, which
 108 not only learns a mapping between features and responses, but also learns the decomposition of the
 109 trained network. This is achieved by factorizing all the layers in the network.

110 3.1 Formulation

111 **Low-rank approximation.** Our inspiration comes from the low-rank matrix approximation of a
 112 matrix $A \in \mathbb{R}^{m \times n}$. For simplicity, we assume that A has rank $r = \min\{m, n\}$ with $k \leq r$ distinct
 113 non-zero singular values $\tilde{\sigma}_1 > \tilde{\sigma}_2 > \dots > \tilde{\sigma}_k > 0$, with corresponding left and right singular vectors
 114 $\tilde{u}_1, \tilde{u}_2, \dots, \tilde{u}_k \in \mathbb{R}^m$ and $\tilde{v}_1, \tilde{v}_2, \dots, \tilde{v}_k \in \mathbb{R}^n$, respectively. For such a matrix, we can rewrite its
 115 best l -rank approximation as the following minimization problem

$$\min_{U \in \mathbb{R}^{m \times l}, V \in \mathbb{R}^{n \times l}} \left\| \sum_{i=1}^l u_i v_i^\top - A \right\|_F^2 \quad (1)$$

116 where c_i denotes the i -th row of matrix C and $\|\cdot\|_F$ denotes Frobenius norm. We note that Prob-
 117 lem (1) is non-convex and non-smooth. However, [60] showed that the randomly initialized gradient
 118 descent algorithm solves this problem in polynomial time. In this work, we consider the best rank
 119 approximation across all the ranks that leads us to the following objective

$$\min_{U \in \mathbb{R}^{m \times r}, V \in \mathbb{R}^{n \times r}} \frac{1}{r} \sum_{b=1}^r \|U_{:,b} V_{:,b}^\top - A\|_F^2, \quad (2)$$

120 where $C_{:,b}$ denotes the first b columns of matrix C . This objective, up to scaling, recovers SVD
 121 of A exactly, and for the case of distinct non-zero singular values, the solution is, up to scaling,
 122 unique [17]. This formulation, however, does not account for the data distribution, i.e., it cannot tailor
 123 the decomposition to capture specific structures that appear in the dataset.

124 **Data-dependent low-rank approximation.** Therefore, the next step of our construction is to
 125 extend this problem formulation with data that can further improve compression, reconstruction, and
 126 generalization, and incorporate domain knowledge. We assume that data comes from the distribution

127 $x \sim \mathcal{X}$ centered around zero, i.e., $\mathbf{E}_{x \sim \mathcal{X}}[x] = 0$ ¹, and the response is given by $y = Ax$. In this
 128 particular case, we can write the training loss as

$$\min_{U \in \mathbb{R}^{m \times r}, V \in \mathbb{R}^{n \times r}} \mathbf{E}_{x, y \sim \mathcal{X}} \left[\sum_{b=1}^r \frac{1}{r} \|U_{:b} V_{:b}^\top x - y\|^2 \right]. \quad (3)$$

129 It is important to note that the introduced problem formulation (3) is the same as the Ordered Dropout
 130 formulation of [17] for the neural network with a single hidden layer and no activations, and it can be
 131 solved using stochastic algorithms by sampling from the data distribution \mathcal{X} (subsampling) and rank
 132 distribution \mathcal{D} . However, there is an important distinction when we apply MAESTRO for deep neural
 133 networks. While FjORD applies uniform dropout across the width of the network for each layer, we
 134 propose to decompose each layer independently to uncover its – potentially different – optimal rank
 135 for deployment. We discuss details in the next paragraph.

136 **DNN low-rank approximation.** For Deep Neural Networks (DNNs), we seek to uncover the optimal
 137 ranks for a set of d linear mappings $W^1 \in \mathbb{R}^{m_1 \times n_1}, \dots, W^d \in \mathbb{R}^{m_d \times n_d}$, where W^i 's are model
 138 parameters and d is model depth, e.g., weights corresponding to linear layers², by decomposing
 139 them as $W^i = U^i (V^i)^\top$. We discuss how these are selected in the next section. To decompose the
 140 network, we aim to minimize the following objective:

$$\mathbf{E}_{x, y \sim \mathcal{X}} \left[\frac{1}{\sum_{i=1}^d r_i} \sum_{i=1}^d \sum_{b=1}^{r_i} l(h(U^1 (V^1)^\top, \dots, U_{:b}^i (V_{:b}^i)^\top, \dots, U^d (V^d)^\top, W^o, x), y) \right], \quad (4)$$

141 where $r_i = \min\{m_i, n_i\}$, l is a loss function, h is a DNN, and W^o are the other weights that we do
 142 not decompose. We note that our formulation aims to decompose each layer, while decompositions
 143 across layers do not directly interact. The motivation for this approach is to uncover low-rank
 144 structures within each layer that are not affected by inaccuracies from other layers due to multiple
 145 low-rank approximations.

146 3.2 Layer factorization

147 The following subsections discuss how model factorization is implemented for different model
 148 architectures.

149 **FC layers.** A 2-layer fully connected (FC) neural network can be expressed as $f(x) =$
 150 $\sigma(\sigma(xW_1)W_2)$, where W s are weight matrices of each FC layer, and $\sigma(\cdot)$ is any arbitrary acti-
 151 vation function, e.g., ReLU. The weight matrix W can be factorized as UV^\top .

152 **CNN layers.** For a convolution layer with dimension, $W \in \mathbb{R}^{m \times n \times k \times k}$ where m and n are the
 153 number of input and output channels, and k is the size of the convolution filters. Instead of directly
 154 factorizing the 4D weight of a convolution layer, we factorize the unrolled 2D matrix. Unrolling the
 155 4D tensor W leads to a 2D matrix with shape $W_{\text{unrolled}} \in \mathbb{R}^{mk^2 \times n}$, where each column represents the
 156 weight of a vectorized convolution filter. Factorization can then be conducted on the unrolled 2D
 157 matrix; see [53] for details.

158 **Transformers.** A Transformer layer consists of a stack of encoders and decoders [51]. The encoder
 159 and decoder contain three main building blocks: the multi-head attention layer, position-wise feed-
 160 forward networks (FFN), and positional encoding. We factorize all trainable weight matrices in the
 161 multi-head attention (HMA) and the FFN layers. The FFN layer factorization can directly adopt the
 162 strategy from the FC factorization. A p -head attention layer learns p attention mechanisms on the
 163 key, value, and query (K, V, Q) of each input token:

$$\text{MHA}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_p) W^O.$$

164 Each head performs the computation of:

$$\text{head}_i = \text{Attention}(QW_Q^{(i)}, KW_K^{(i)}, VW_V^{(i)}) = \text{softmax} \left(\frac{QW_Q^{(i)} W_K^{(i)\top} K^\top}{\sqrt{d/p}} \right) VW_V^{(i)}.$$

165 where d is the hidden dimension. The trainable weights $W_Q^{(i)}, W_K^{(i)}, W_V^{(i)}, i \in \{1, 2, \dots, p\}$ can
 166 be factorized by simply decomposing all learnable weights $W \cdot$ in an attention layer and obtaining
 167 $U \cdot V^\top$. [51].

¹We make this assumption for simplicity. It can be simply overcome by adding a bias term into the model.

²We can apply our decomposition on different types of layers, such as Linear, Convolutional and Transformers as shown in Sec. 3.2.

168 3.3 Training techniques

169 Having defined the decomposition of typical layers found in DNNs, we move to formulate the
 170 training procedure of our method, formally described in Algorithm 1. Training the model comprises
 171 an iterative process of propagating forward on the model by *sampling a rank* b_i per decomposed layer
 172 i up to maximal rank r_i (line 3). We calculate the loss, which integrates an additional *hierarchical*
 173 *group lasso* component (lines 4) and *backpropagate* on the sampled decomposed model (line 5). At
 174 the end of each epoch, we *progressively shrink* the network by updating the maximal rank r_i , based
 175 on an importance threshold ε_{ps} (line 11). We provide more details about each component below.

Algorithm 1: MAESTRO (Training Process)

Input: epochs E , dataset \mathcal{D} , model h parametrized by $U^1 \in \mathbb{R}^{m_1 \times r_1}$,
 $V^1 \in \mathbb{R}^{n_1 \times r_1}, \dots, U^d \in \mathbb{R}^{m_d \times r_d}, V^d \in \mathbb{R}^{n_d \times r_d}, W^o$, and hyperparameters $\lambda_{gl}, \varepsilon_{ps}$

```

1 for  $t \leftarrow 0$  to  $E - 1$  do // Epochs
2   for  $(x, y) \in \mathcal{D}$  do // Iterate over dataset
3     Sample  $(i, b) \sim \{\{(i, b)\}_{b=1}^{r_i}\}_{i=1}^d$ ;
4      $L = l(h(U^1(V^1)^\top, \dots, U_b^i(V_b^i)^\top, \dots, U^d(V^d)^\top), W^o, x, y) +$ 
        $+\lambda_{gl} \sum_{i=1}^d \sum_{b=1}^{r_i} (\|U_b^i\| + \|V_b^i\|)$  // compute loss
5     L.backward() // Update weights
6   end
7   for  $i \leftarrow 1$  to  $d$  do
8     for  $b \leftarrow 1$  to  $r_i$  do
9       // rank importance thresholding
10      if  $\|V_b^i\| \|U_b^i\| \leq \varepsilon_{ps}$  then
11         $r_i = b - 1$  // progressive shrinking
12        break
13      end
14    end
15  end
16 end
```

176 **Efficient training via sampling.** In Sec. 4, we show that for the linear case (3), the optimal solution
 177 corresponds to PCA over the linearly transformed dataset. This means that the obtained solution
 178 contains orthogonal directions. This property is beneficial because it directly implies that when we
 179 employ gradient-based optimization, not only is the gradient zero at the optimum, but the gradient
 180 with respect to each summand in Equation (3) is also zero. The same property is directly implied
 181 by overparametrization [35] or strong growth condition [44]. As a consequence, this enables us to
 182 sample only one summand at a time and obtain the same quality solution. When considering (4)
 183 as an extension to (3), it is unclear whether this property still holds, which would also imply that
 184 the set of stationary points of (3) is a subset of stationary points of the original objective without
 185 decomposition. However, in the experiments, we observed that sampling is sufficient to converge to a
 186 good-quality solution. If this only holds approximately, we one could leverage fine-tuning to recover
 187 the loss in performance.

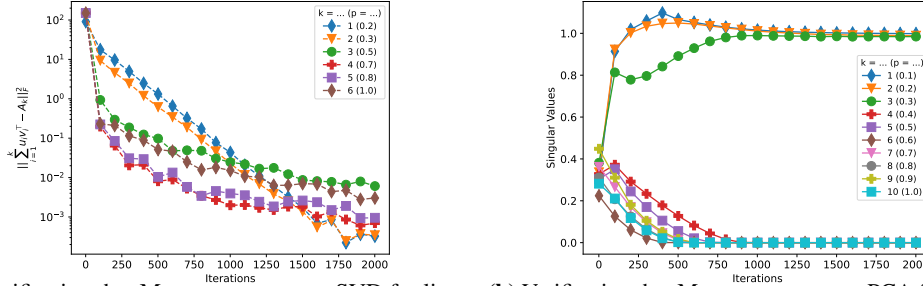
188 **Efficient rank extraction via hierarchical group-lasso.** By definition, (3) leads to an ordered set
 189 of ranks for each layer. This ordered structure enables efficient rank extraction and selection. To
 190 effectively eliminate unimportant ranks while retaining the important ones, thus leading to a more
 191 efficient model, we consider Hierarchical Group Lasso (HGL) [31] in the form

$$\lambda_{gl} \sum_{i=1}^d \sum_{b=1}^{r_i} (\|U_b^i\| + \|V_b^i\|), \quad (5)$$

192 where C_b denotes the matrix that contains all the columns of C except for the first $b - 1$ columns.

193 **Progressive shrinking.** HGL encourages that unimportant ranks become zero and can be effectively
 194 removed from the model. To account for this, for each layer we remove V_b^i and U_b^i (i.e., set $r_i = b - 1$)
 195 if $\|V_b^i\| \|U_b^i\| \leq \varepsilon_{ps}$, where ε_{ps} is a pre-selected threshold – and a hyperparameter of our method.

196 **Initialization.** Initialization is a key component of the training procedure [13, 37]. To adopt the
 197 best practices from standard non-factorized training, we follow a similar approach to [23, 53], where
 198 we first initialize the non-factorized model using standard initialization. For initializing factorized
 199 layers, we use the Singular Value Decomposition (SVD) of the non-factorized initialization – in a



(a) Verification that MAESTRO recovers SVD for linear mapping with uniform data. The plot displays the L2 distance between the best rank k and MAESTRO's approximation of mapping A . The target matrix was randomly generated 9×6 matrix with rank 3. p and k stand for relative and actual rank, respectively.

(b) Verification that MAESTRO recovers PCA for identity mapping. The plot displays the estimates of singular values. The data distribution has only 3 directions. It is expected that the top 3 ranks will converge to value one and the rest to zero. p and k stand for relative and actual rank, respectively.

Figure 2: Empirical showcase of theoretical properties of the MAESTRO's formulation.

200 full-rank form – to ensure that the resulting product matrix is the same as the original parameter
 201 decomposition. In addition, SVD is an optimal decomposition for the linear case with uniform data.
 202 However, in contrast with the adaptive baseline method [54] we only decompose once, rather than on
 203 every training iteration.

204 3.4 Train-once, deploy-everywhere

205 Up until now, we have described how our method works for training low-rank models, which yield
 206 computational, memory, network, and energy [57] bandwidth benefits during training. At deployment
 207 time, one can directly deploy the final model (rank r_i for each layer) on the device, which we
 208 acquire from performing a threshold sweep of ε_{ps} over the effective range of rank importance across
 209 layers. However, in case we want to run on even more constrained devices, such as mobile [4] or
 210 embedded [4] systems, the learned decomposition also gives us the flexibility to further compress
 211 the model in a straightforward manner, effectively trading off accuracy for a smaller model footprint.
 212 Inspired by [61], we propose to use greedy search. We begin with the current model and compare
 213 model performance across various low-rank models, each created by removing a certain percentage
 214 of ranks from each layer. We then eliminate the ranks that cause the least decrease in performance.
 215 This process is iterated until we reach the desired size or accuracy constraint. To make this approach
 216 efficient, we estimate the loss using a single mini-batch with a large batch size, for example, 2048.
 217 This also avoids issues with BatchNorm layers; see [61] for details.

218 In summary, MAESTRO comprises a technique for trainable low-rank approximation during training
 219 time that progressively compresses the model, reflecting the data distribution, and a method that
 220 enables a graceful trade-off between accuracy and latency for embedded deployment, by selecting
 221 the most important parts of the network. We validate these claims in Sec. 5.2 and 5.5, respectively.

222 4 Theoretical guarantees

223 In this section, we further investigate the theoretical properties of MAESTRO for the linear mappings,
 224 i.e., the setup of the problem formulation (3).

225 **Theorem 4.1** (Informal). *Let $A = \tilde{U}\tilde{\Sigma}\tilde{V}^\top$ be a SVD decomposition of A . Then, the minimization*
 226 *problem (3) is equivalent to PCA applied to the transformed dataset $x \rightarrow \tilde{\Sigma}\tilde{V}^\top x$, $x \sim \mathcal{X}$ projected*
 227 *on the column space of \tilde{U} .*

228 The formal statement can be found in Appendix D. Theorem 4.1 shows that MAESTRO can adapt to
 229 data distribution by directly operating on data $x \sim \mathcal{X}$ and also to the target mapping by projecting
 230 data to its right singular vectors scaled by singular values. In particular, we show that in the special
 231 case, when \mathcal{X} is the uniform distribution on the unit ball, (3), i.e., MAESTRO, exactly recovers
 232 truncated SVD of A , which is consistent with the prior results [17]. In the case A is the identity, it is
 233 straightforward to see that MAESTRO is equivalent to PCA. We can see that MAESTRO can efficiently
 234 extract low-rank solutions by filtering out directions corresponding to the null space of the target
 235 mapping A and directions with no data. We also numerically verify both of the special cases—PCA
 236 and SVD, by minimizing (3) using stochastic gradient descent (SGD) with \mathcal{D} being the uniform
 237 distribution. These preliminary experiments are provided in Fig. 2a and 2b.

238 We showed that MAESTRO could recover SVD in a particular case of the linear model and the uniform
 239 data distribution on the unit ball. We note that in this case, SVD is optimal, and we cannot acquire
 240 better decomposition. Therefore, it is desired that MAESTRO is equivalent to SVD in this scenario. In
 241 the more general setting, we argue that MAESTRO decomposition should be preferable to SVD due to
 242 the following reasons:

- 243 • MAESTRO formulation is directly built into the training and tailored to obtain the best low-rank
 244 decomposition, while SVD relies on linearity assumption.
- 245 • SVD does not account for data, and even in the linear NN case, the learned singular vectors might
 246 exhibit wrong ordering. We demonstrate this issue using a simple example where we take matrix
 247 A with rank 3. We construct the dataset \mathcal{X} in such a way that the third singular vector is the most
 248 important, the second one is the second, and the first is the third most important direction. Clearly,
 249 SVD does not look at data. Therefore, it cannot capture this phenomenon. We showcase that
 250 MAESTRO learns the correct order; see Fig. 5 of the Appendix.
- 251 • Pre-factorizing models allow us to apply hierarchical group-lasso penalty [64] for decomposed
 252 weights to directly regularize the rank of different layers.
- 253 • SVD is computationally expensive and can only run rarely, while MAESTRO is directly built into
 254 the training and, therefore, does not require extra computations. In addition, MAESTRO supports
 255 rank sampling so training can be made computationally efficient.

256 5 Experiments

257 We start this section by describing the setup of our experiments, including the models, datasets and
 258 baselines with which we compare MAESTRO. We then compare MAESTRO against the baselines on
 259 accuracy and MAC and discuss the results. Subsequently, we analyze the behaviour of our system
 260 in-depth and provide additional insights on the performance of our technique, along with an ablation
 261 study and sensitivity analysis to specific hyperparameters. Finally, we showcase the performance
 262 of models upon deployment and how we can derive a smaller footprint model with some accuracy
 263 trade-off, without the need to fine-tune.

264 5.1 Experimental setup

265 **Models & datasets.** The datasets and models considered in our experiments span across four datasets,
 266 concisely presented along with the associated models on Tab. 1. We have implemented our solution
 267 with PyTorch [38](v1.13.0) trained our models on NVidia A100 (40G) GPUs. Details for the learning
 268 tasks and hyperparameters used are presented in the Appendix.

269 **Baselines.** We have selected various baselines from the literature that we believe are closest to
 270 aspects of our system. On the *pruning* front, we compare with the IMP [40] and RareGems [48]
 271 techniques, themselves based on the LTH [11]. On the *quantization* front, we compare with
 272 XNOR-Net [41]. With respect to *low-rank* meth-
 273 ods, we compare with Spectral Initialisation [23], Pufferfish [53] and Cuttlefish [54].

Table 1: Datasets and models for evaluation. The network footprints depict the vanilla variants of the models.

Dataset	Model	# GMACs	# Params (M)	Task
MNIST	LeNet	$2e^{-4}$	0.04	Image classification
CIFAR10	ResNet-18	0.56	11.18	Image classification
CIFAR10	VGG-19	0.40	20.00	Image classification
TinyImageNet	ResNet-50	5.19	53.9	Image classification
Multi30k	6-layer Transformer	1.37	48.98	Translation (en-ge)

277 5.2 Performance comparison

278 We start off by comparing MAESTRO with various baselines from the literature across different
 279 datasets and types of models³. Results are depicted in Tab. 2 and 3, while additional performance
 280 points of MAESTRO for different model footprints are presented in the Appendix F.2 and F.3.

281 **Comparisons with low-rank methods.** The low-rank methods we are comparing against are
 282 Pufferfish [53] and Cuttlefish [54]. These methods try to reduce training and inference runtime while
 283 preserving model accuracy by leveraging low-rank approximations. For ResNet-18, we achieve
 284 $94.19 \pm 0.07\%$ for 4.08M parameters and $93.97 \pm 0.25\%$ for 2.19M parameters compared to the
 285 94.17% of Pufferfish at 3.3M parameters. For VGG-19, we achieve +0.41pp (percentage points)
 286 higher accuracy compared to Pufferfish and -0.29pp to Cuttlefish at 44.8% and 93.2% of the sizes,

³The operating points we select for MAESTRO are the closest lower to the respective baseline in terms of footprint. Where the result is not present in the Tab. 2, we provide the λ_{gp} value so that it can be referenced from the Appendix, Tab. 11, 12.

Table 2: Maestro vs. baselines on CIFAR10.

Variant	Model	Acc. (%)	GMACs	Params. (M)
Non-factorized	ResNet-18	93.86 \pm 0.20	0.56	11.17
Pufferfish	ResNet-18	94.17	0.22	3.336
Cuttlefish	ResNet-18	93.47	0.3	3.108
IMP	ResNet-18	92.12	-	0.154
RareGems	ResNet-18	92.83	-	0.076
XNOR-Net	ResNet-18	90.06	-	0.349 [†]
MAESTRO [†] ($\lambda_{gp} = 16e^{-6}$)	ResNet-18	94.19\pm0.07	0.39 \pm 0.00	4.08 \pm 0.02
MAESTRO [†] ($\lambda_{gp} = 64e^{-6}$)	ResNet-18	93.86 \pm 0.11	0.15 \pm 0.00	1.23 \pm 0.00
Non-factorized	VGG-19	92.94 \pm 0.17	0.40	20.56
Pufferfish	VGG-19	92.69	0.29	8.37
Cuttlefish	VGG-19	93.39	0.15	2.36
RareGems	VGG-19	86.28	-	5.04
IMP	VGG-19	92.86	-	5.04
XNOR-Net	VGG-19	88.94	-	0.64 [†]
Spectral Init.*	VGG-19	83.27	-	\approx 0.4
MAESTRO [†] ($\lambda_{gp} = 32e^{-6}$)	VGG-19	93.10 \pm 0.10	0.13 \pm 0.00	2.20 \pm 0.03
MAESTRO [†] ($\lambda_{gp} = 512e^{-6}$)	VGG-19	88.53 \pm 0.13	0.03\pm0.00	0.35\pm0.00

*Results from original work; †: XNOR-Net employs binary weights and activations; although the overall #trainable parameters remain the same as the vanilla network, each model weight is quantized from 32-bit to 1-bit. Therefore, we report a compression rate of 3.125% ($1/32$).

287 respectively. Finally, comparing with the spectral initialization [23] for VGG-19, we achieve +5.26pp
 288 higher accuracy for 87.5% of parameter size. Detailed results are shown in Tab. 2. This performance
 289 benefits also apply in the case of Transformers (Tab. 3), where MAESTRO performs 6% better in terms
 290 of perplexity at 25% of the cost (MACs) and 51.7% of the size (parameters) compared to Pufferfish.

291 **Comparisons with pruning methods.** The next family of baselines is related to the LTH [11].
 292 Specifically, we compare against IMP [40] and witness from Tab. 2 that MAESTRO can achieve
 293 +1.25pp ($\lambda_{gp} = 128e^{-6}$) and +0.24pp ($\lambda_{gp} = 32e^{-6}$) higher accuracy for ResNet-18 and VGG-19
 294 respectively. Although we cannot scale to the size that RareGems [48] for ResNet-18, the sparsity
 295 that they achieve is unstructured, which most modern hardware cannot take advantage of. In contrast,
 296 our technique performs ordered structured sparsity, compatibly with most computation targets. On
 297 the other hand, for VGG-19, we achieve +6.82pp higher accuracy at 43.6% of the footprint.

298 **Comparisons with quantized models.** We also compare against XNOR-Net [41], which binarizes
 299 the network to achieve efficient inference. Training continues to happen in full precision, and
 300 inference performance is really dependent on the operation implementation of the target hardware.
 301 Nonetheless, assuming a compression rate of 3.125%, for the same size of network, we achieve
 302 +1.08pp ($\lambda_{gp} = 512e^{-6}$) and +2.18pp ($\lambda_{gp} = 256e^{-6}$) higher accuracy on ResNet-18 and VGG-19.

303 5.3 Training behaviour of MAESTRO

304 Having shown the relative performance of our framework to selected baselines, we now move to
 305 investigate how our method behaves, with respect to its convergence and low-rank approximations.

306 **Model and rank convergence.** In Fig. 3, we present the training dynamics for MAESTRO. Fig. 3a
 307 illustrates the evolution of total rank throughout the training steps. We observe that the ranks are
 308 pruned incrementally. This aligns with the observations made during Pufferfish [53] training, where
 309 the authors suggest warm-start training with full precision to enhance the final model performance.
 310 In our situation, we do not need to integrate this heuristic because MAESTRO automatically prunes
 311 rank. Fig. 3b reveals the ranks across layers after training. We notice an intriguing phenomenon: the
 312 ranks are nested for increasing λ_{gl} . This could imply apart from a natural order of ranks within each
 313 layer, a global order. We briefly examine this captivating occurrence in the following section, and we
 314 plan to investigate it more thoroughly in future work, as we believe this might contribute to a superior
 315 rank selection and sampling process. Lastly, Fig. 3c depicts the progression of training loss. We find
 316 that our hypothesis, that sampling does not adversely impact training, is also supported empirically.

317 5.4 Ablation study

318 In this section, we examine the impact of each component on the performance of MAESTRO.
 319 Specifically, we run variants of our method *i)* without the *hierarchical group lasso regularization*
 320 (*HGL*), *ii)* without progressive *shrinking* (*PS*). Additionally, we integrate *iii)* an *extra full low-rank*
 321 *pass* ($b = r_i$) into the training at each step to assess whether extra sampling would be beneficial.
 322 The results are displayed in Tab. 4. As anticipated, our findings confirm that neither inclusion of
 323 hierarchical group lasso with a tuned λ_{gl} nor progressive shrinking impair the final performance,

Table 3: Maestro vs. baselines on Multi30k.

Variant	Model	Perplexity	GMACs	Params. (M)
Non-factorized	Transformer	9.85 \pm 0.10	1.370	53.90
Pufferfish*	Transformer	7.34 \pm 0.12	0.996	26.70
MAESTRO [†]	Transformer	6.90\pm0.07	0.248\pm0.0032	13.80\pm0.113

*Results from original work; † tuned λ_{gp} from $\{2^i/100; i \in 0, \dots, 9\}$

Table 4: Ablation study for ResNet18 on CIFAR10

Variant	Acc. (%)	GMACs	Params. (M)
MAESTRO	94.19 \pm 0.39	0.39 \pm 0.0008	4.08 \pm 0.020
w/out GL	94.04 \pm 0.10	0.56 \pm 0.0000	11.2 \pm 0.000
w/out PS	94.12 \pm 0.36	0.39 \pm 0.0010	4.09 \pm 0.027
w/ full-training	94.05 \pm 0.32	0.39 \pm 0.0004	4.09 \pm 0.032

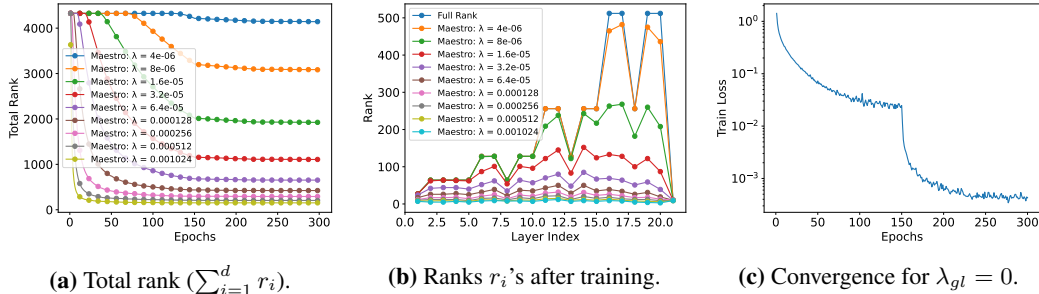


Figure 3: Training dynamics of MAESTRO for ResNet18 on CIFAR10.

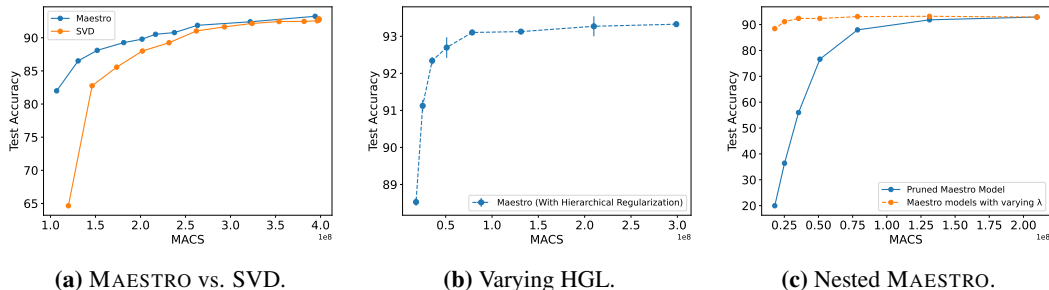


Figure 4: Accuracy-latency trade-off of MAESTRO under different settings for VGG19 on CIFAR10.

324 but they do significantly enhance the efficiency of MAESTRO. Moreover, sampling more ranks at
 325 each training step does not improve the final performance, and, in fact, it hampers training efficiency,
 326 making it approximately twice as computationally demanding.

327 **5.5 Accuracy-latency trade-off at training and deployment time**

328 In Fig. 4, we illustrate various approaches to balance latency (proxied through MACs operations)
 329 and accuracy in model training and deployment. Fig. 4a demonstrates how MAESTRO ($\lambda_{gl} = 0$) can
 330 be pruned effectively for deployment using the greedy search method discussed in Section 3.4. We
 331 contrast this with the greedy pruning of a non-factorized model that has been factorized using SVD.
 332 We reveal that this straightforward baseline does not measure up to the learned decomposition of
 333 MAESTRO and results in a significant performance decrease. Next, Fig. 4b portrays the final accuracy
 334 and the number of model parameters for varying hierarchical group lasso penalties. This leads to the
 335 optimal latency-accuracy balance for both training and inference. However, it's crucial to point out
 336 that each model was trained individually, while greedy pruning only necessitates a single training
 337 cycle. Lastly, we delve into the observation of nested ranks across increasing λ_{gl} . Fig. 4c displays
 338 the performance of MAESTRO ($\lambda_{gl} = 0$) across different ranks selected by smaller models MAESTRO
 339 ($\lambda_{gl} > 0$). Intriguingly, we observe that MAESTRO ($\lambda_{gl} = 0$) performs very well—for instance, we
 340 can decrease its operations in half (and parameters by 10 \times) and still maintain an accuracy of 87.7%
 341 without fine-tuning, just by reusing rank structure from independent runs. As aforementioned, we
 342 intend to further explore this in the future.

343 **6 Conclusion and future work**

344 In this work, we have presented MAESTRO, a method for trainable low-rank approximation of DNNs
 345 that leverages progressive shrinking by applying a generalized variant of Ordered Dropout to the
 346 factorized weights. We have shown the theoretical guarantees of our work in the case of linear
 347 models and empirically demonstrated its performance across different types of models, datasets, and
 348 modalities. Our evaluation has demonstrated that MAESTRO outperforms competitive compression
 349 methods at a lower cost. In the future, we plan to expand our technique to encompass more advanced
 350 sampling techniques and apply it to different distributed learning scenarios, such as Federated
 351 Learning, where data are natively non-independent or identically distributed (non-IID).

References

- 352
- 353 [1] Samiul Alam, Luyang Liu, Ming Yan, and Mi Zhang. Fedrolex: Model-heterogeneous federated learning
354 with rolling sub-model extraction. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun
355 Cho, editors, *Advances in Neural Information Processing Systems*, 2022.
- 356 [2] Dan Alistarh, Demjan Grubic, Jerry Li, Ryota Tomioka, and Milan Vojnovic. Qsgd: Communication-
357 efficient sgd via gradient quantization and encoding. *Advances in neural information processing systems*,
358 30, 2017.
- 359 [3] Dan Alistarh, Torsten Hoefler, Mikael Johansson, Sarit Khirirat, Nikola Konstantinov, and Cédric Renggli.
360 The convergence of sparsified gradient methods. *arXiv preprint arXiv:1809.10505*, 2018.
- 361 [4] Mario Almeida, Stefanos Laskaridis, Abhinav Mehrotra, Lukasz Dudziak, Ilias Leontiadis, and Nicholas D
362 Lane. Smart at what cost? characterising mobile deep neural networks in the wild. In *Proceedings of the*
363 *21st ACM Internet Measurement Conference*, pages 658–672, 2021.
- 364 [5] Sebastian Caldas, Jakub Konečný, Brendan McMahan, and Ameet Talwalkar. Expanding the reach of
365 federated learning by reducing client resource requirements, 2019.
- 366 [6] Miguel A Carreira-Perpinán and Yerlan Idelbayev. “learning-compression” algorithms for neural net
367 pruning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages
368 8532–8541, 2018.
- 369 [7] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical
370 image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255.
371 Ieee, 2009.
- 372 [8] Enmao Diao, Jie Ding, and Vahid Tarokh. Hetero{fl}: Computation and communication efficient federated
373 learning for heterogeneous clients. In *International Conference on Learning Representations*, 2021.
- 374 [9] Łukasz Dudziak, Mohamed S Abdelfattah, Ravichander Vipperla, Stefanos Laskaridis, and Nicholas D
375 Lane. Shrinkml: End-to-end asr model compression using reinforcement learning. *INTERSPEECH*, 2019.
- 376 [10] D. Elliott, S. Frank, K. Sima’an, and L. Specia. Multi30k: Multilingual english-german image descriptions.
377 pages 70–74, 2016.
- 378 [11] Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural
379 networks. In *International Conference on Learning Representations*, 2019.
- 380 [12] Song Han, Huizi Mao, and William J Dally. Deep compression: Compressing deep neural networks with
381 pruning, trained quantization and huffman coding. *arXiv preprint arXiv:1510.00149*, 2015.
- 382 [13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing
383 human-level performance on imagenet classification. In *Proceedings of the IEEE international conference*
384 *on computer vision*, pages 1026–1034, 2015.
- 385 [14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition.
386 In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- 387 [15] Yihui He, Xiangyu Zhang, and Jian Sun. Channel pruning for accelerating very deep neural networks. In
388 *Proceedings of the IEEE international conference on computer vision*, pages 1389–1397, 2017.
- 389 [16] Samuel Horváth, Aaron Klein, Peter Richtárik, and Cédric Archambeau. Hyperparameter transfer learning
390 with adaptive complexity. In *International Conference on Artificial Intelligence and Statistics*, pages
391 1378–1386. PMLR, 2021.
- 392 [17] Samuel Horváth, Stefanos Laskaridis, Mario Almeida, Ilias Leontiadis, Stylianos Venieris, and Nicholas
393 Lane. FjORD: Fair and accurate federated learning under heterogeneous targets with ordered dropout.
394 *Advances in Neural Information Processing Systems*, 34:12876–12889, 2021.
- 395 [18] Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Ges-
396 mundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for nlp. In *International*
397 *Conference on Machine Learning*, pages 2790–2799. PMLR, 2019.
- 398 [19] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco
399 Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision
400 applications. *arXiv preprint arXiv:1704.04861*, 2017.

- 401 [20] Edward Hu, Yelong Shen, Phil Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Lu Wang, and Weizhu Chen. Lora:
402 Low-rank adaptation of large language models, 2021.
- 403 [21] Hengyuan Hu, Rui Peng, Yu-Wing Tai, and Chi-Keung Tang. Network trimming: A data-driven neuron
404 pruning approach towards efficient deep architectures. *arXiv preprint arXiv:1607.03250*, 2016.
- 405 [22] Max Jaderberg, Andrea Vedaldi, and Andrew Zisserman. Speeding up convolutional neural networks with
406 low rank expansions. *arXiv preprint arXiv:1405.3866*, 2014.
- 407 [23] Mikhail Khodak, Neil Tenenholz, Lester Mackey, and Nicolo Fusi. Initialization and regularization of
408 factorized neural layers. *arXiv preprint arXiv:2105.01029*, 2021.
- 409 [24] Minjae Kim, Sangyoon Yu, Suhyun Kim, and Soo-Mook Moon. DepthFL : Depthwise federated learning
410 for heterogeneous clients. In *The Eleventh International Conference on Learning Representations*, 2023.
- 411 [25] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- 412 [26] Stefanos Laskaridis, Alexandros Kouris, and Nicholas D Lane. Adaptive inference through early-exit
413 networks: Design, challenges and directions. In *Proceedings of the 5th International Workshop on
414 Embedded and Mobile Deep Learning*, pages 1–6, 2021.
- 415 [27] Stefanos Laskaridis, Stylianos I Venieris, Alexandros Kouris, Rui Li, and Nicholas D Lane. The future of
416 consumer edge-ai computing. *arXiv preprint arXiv:2210.10514*, 2022.
- 417 [28] Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. *CS 231N*, 7(7):3, 2015.
- 418 [29] Yann LeCun, Corinna Cortes, and CJ Burges. Mnist handwritten digit database. *ATT Labs [Online]*.
419 Available: <http://yann.lecun.com/exdb/mnist>, 2, 2010.
- 420 [30] Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning filters for efficient
421 convnets. *arXiv preprint arXiv:1608.08710*, 2016.
- 422 [31] Michael Lim and Trevor Hastie. Learning interactions via hierarchical group-lasso regularization. *Journal
423 of Computational and Graphical Statistics*, 24(3):627–654, 2015.
- 424 [32] Ji Lin, Ligeng Zhu, Wei-Ming Chen, Wei-Chen Wang, Chuang Gan, and Song Han. On-device training
425 under 256kb memory. In *Annual Conference on Neural Information Processing Systems (NeurIPS)*, 2022.
- 426 [33] Zhuang Liu, Mingjie Sun, Tinghui Zhou, Gao Huang, and Trevor Darrell. Rethinking the value of network
427 pruning. *arXiv preprint arXiv:1810.05270*, 2018.
- 428 [34] Zicheng Liu, Da Li, Javier Fernandez-Marques, Stefanos Laskaridis, Yan Gao, Łukasz Dudziak, Stan Z Li,
429 Shell Xu Hu, and Timothy Hospedales. Federated learning for inference at anytime and anywhere. *arXiv
430 preprint arXiv:2212.04084*, 2022.
- 431 [35] Siyuan Ma, Raef Bassily, and Mikhail Belkin. The power of interpolation: Understanding the effectiveness
432 of sgd in modern over-parametrized learning. In *International Conference on Machine Learning*, pages
433 3325–3334. PMLR, 2018.
- 434 [36] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas.
435 Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and
436 statistics*, pages 1273–1282. PMLR, 2017.
- 437 [37] Dmytro Mishkin and Jiri Matas. All you need is a good init. *arXiv preprint arXiv:1511.06422*, 2015.
- 438 [38] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming
439 Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017.
- 440 [39] David Patterson, Joseph Gonzalez, Urs Hölzle, Quoc Le, Chen Liang, Lluís-Miquel Munguia, Daniel
441 Rothchild, David R So, Maud Texier, and Jeff Dean. The carbon footprint of machine learning training
442 will plateau, then shrink. *Computer*, 55(7):18–28, 2022.
- 443 [40] Mansheej Paul, Feng Chen, Brett W. Larsen, Jonathan Frankle, Surya Ganguli, and Gintare Karolina
444 Dziugaite. Unmasking the lottery ticket hypothesis: What’s encoded in a winning ticket’s mask? In *The
445 Eleventh International Conference on Learning Representations*, 2023.
- 446 [41] Mohammad Rastegari, Vicente Ordonez, Joseph Redmon, and Ali Farhadi. Xnor-net: Imagenet classi-
447 fication using binary convolutional neural networks. In *Computer Vision–ECCV 2016: 14th European
448 Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV*, pages 525–542.
449 Springer, 2016.

- 450 [42] Oren Rippel, Michael Gelbart, and Ryan Adams. Learning Ordered Representations with Nested Dropout.
451 In *International Conference on Machine Learning (ICML)*, pages 1746–1754, 2014.
- 452 [43] Tara N Sainath, Brian Kingsbury, Vikas Sindhwani, Ebru Arisoy, and Bhuvana Ramabhadran. Low-rank
453 matrix factorization for deep neural network training with high-dimensional output targets. In *2013 IEEE
454 international conference on acoustics, speech and signal processing*, pages 6655–6659. IEEE, 2013.
- 455 [44] Mark Schmidt and Nicolas Le Roux. Fast convergence of stochastic gradient descent under a strong growth
456 condition. *arXiv preprint arXiv:1308.6370*, 2013.
- 457 [45] Frank Seide, Hao Fu, Jasha Droppo, Gang Li, and Dong Yu. 1-bit stochastic gradient descent and its
458 application to data-parallel distributed training of speech dnns. In *Fifteenth annual conference of the
459 international speech communication association*, 2014.
- 460 [46] Hakim Sidahmed, Zheng Xu, Ankush Garg, Yuan Cao, and Mingqing Chen. Efficient and private federated
461 learning with partially trainable networks. *arXiv preprint arXiv:2110.03450*, 2021.
- 462 [47] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recogni-
463 tion. In *International Conference on Learning Representations*, 2015.
- 464 [48] Kartik Sreenivasan, Jy yong Sohn, Liu Yang, Matthew Grinde, Alliot Nagle, Hongyi Wang, Eric Xing,
465 Kangwook Lee, and Dimitris Papailiopoulos. Rare gems: Finding lottery tickets at initialization. In Alice H.
466 Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information
467 Processing Systems*, 2022.
- 468 [49] Ananda Theertha Suresh, X Yu Felix, Sanjiv Kumar, and H Brendan McMahan. Distributed mean estimation
469 with limited communication. In *International Conference on Machine Learning*, pages 3329–3337. PMLR,
470 2017.
- 471 [50] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In
472 *International conference on machine learning*, pages 6105–6114. PMLR, 2019.
- 473 [51] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
474 Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing
475 systems*, pages 5998–6008, 2017.
- 476 [52] Erwei Wang, James J Davis, Ruizhe Zhao, Ho-Cheung Ng, Xinyu Niu, Wayne Luk, Peter YK Cheung,
477 and George A Constantinides. Deep Neural Network Approximation for Custom Hardware: Where we’ve
478 been, where we’re going. *ACM Computing Surveys (CSUR)*, 52(2):1–39, 2019.
- 479 [53] Hongyi Wang, Saurabh Agarwal, and Dimitris Papailiopoulos. Pufferfish: communication-efficient models
480 at no extra cost. *Proceedings of Machine Learning and Systems*, 3:365–386, 2021.
- 481 [54] Hongyi Wang, Saurabh Agarwal, Yoshiki Tanaka, Eric P Xing, Dimitris Papailiopoulos, et al. Cuttlefish:
482 Low-rank model training without all the tuning. *arXiv preprint arXiv:2305.02538*, 2023.
- 483 [55] Wei Wen, Chunpeng Wu, Yandan Wang, Yiran Chen, and Hai Li. Learning structured sparsity in deep
484 neural networks. *Advances in neural information processing systems*, 29, 2016.
- 485 [56] Simon Wiesler, Alexander Richard, Ralf Schlüter, and Hermann Ney. Mean-normalized stochastic gradient
486 for large-scale deep learning. In *2014 IEEE International Conference on Acoustics, Speech and Signal
487 Processing (ICASSP)*, pages 180–184. IEEE, 2014.
- 488 [57] Carole-Jean Wu, Ramya Raghavendra, Udit Gupta, Bilge Acun, Newsha Ardalani, Kiwan Maeng, Gloria
489 Chang, Fiona Aga, Jinshi Huang, Charles Bai, et al. Sustainable ai: Environmental implications, challenges
490 and opportunities. *Proceedings of Machine Learning and Systems*, 4:795–813, 2022.
- 491 [58] Jian Xue, Jinyu Li, and Yifan Gong. Restructuring of deep neural network acoustic models with singular
492 value decomposition. In *Interspeech*, pages 2365–2369, 2013.
- 493 [59] Tien-Ju Yang, Yu-Hsin Chen, and Vivienne Sze. Designing energy-efficient convolutional neural networks
494 using energy-aware pruning. In *Proceedings of the IEEE conference on computer vision and pattern
495 recognition*, pages 5687–5695, 2017.
- 496 [60] Tian Ye and Simon S Du. Global convergence of gradient descent for asymmetric low-rank matrix
497 factorization. *Advances in Neural Information Processing Systems*, 34:1429–1439, 2021.
- 498 [61] Jiahui Yu and Thomas Huang. Autoslim: Towards one-shot architecture search for channel numbers. *arXiv
499 preprint arXiv:1903.11728*, 2019.

- 500 [62] Jiahui Yu and Thomas S Huang. Universally slimmable networks and improved training techniques. In
501 *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1803–1811, 2019.
- 502 [63] Jiahui Yu, Linjie Yang, Ning Xu, Jianchao Yang, and Thomas Huang. Slimmable neural networks. In
503 *International Conference on Learning Representations*, 2019.
- 504 [64] Ming Yuan and Yi Lin. Model selection and estimation in regression with grouped variables. *Journal of*
505 *the Royal Statistical Society: Series B (Statistical Methodology)*, 68(1):49–67, 2006.
- 506 [65] Michael Zhu and Suyog Gupta. To prune, or not to prune: exploring the efficacy of pruning for model
507 compression. *arXiv preprint arXiv:1710.01878*, 2017.

508 Appendix

509 Contents of the Appendix

510	A Broader impact	14
511	B Limitations	14
512	C Extended Background	15
513	D Theoretical Properties of Low-Rank Layers	15
514	E Experimental setup	17
515	E.1 Datasets	17
516	E.2 Models	17
517	E.3 Hyperparameter selection	18
518	E.4 Deciding against decomposition	19
519	F Extended evaluation	19
520	F.1 MAESTRO recovers correct ordering	19
521	F.2 Training behaviour of MAESTRO	19
522	F.3 Model size-accuracy trade-off at training and deployment time	20

523 A Broader impact

524 The goal of our work is to make the training and deployment of DNNs more efficient, affecting the
525 total computation, memory and bandwidth of systems, as well as the energy they require to run the
526 respective tasks. DNN model training requires significant amounts of energy, whether in a data center
527 or at the edge [57, 39]. However, such techniques should not be used as an excuse to make data
528 centers less green, but rather as a complementary measure to further reduce the carbon footprint of
529 Deep Learning.

530 Additionally, as our technique involves a training-aware methodology for progressively selecting
531 ranks, it depends on the quality of data used in training. Deploying the model in the wild for various
532 downstream tasks may result in behavior different from the intended outcomes. Therefore, it should
533 be thoroughly tested before deployment to ensure it adheres to the required Service Level Objectives
534 (SLOs), especially in performance-critical use cases, such as self-driving vehicles or UAV navigation.

535 B Limitations

536 In this work, we have proposed a method for trainable low-rank approximation of DNNs that provides
537 performance benefits for both training and inference times. While we suggest that this could have
538 repercussions on the energy consumption of these tasks, we have not yet evaluated this hypothesis
539 experimentally across different devices, be they data center-grade or at the edge.

540 Additionally, we have applied our technique to CNN and Transformer models spanning across vision
541 and NLP tasks. While we anticipate generalization to any type of network, it remains to be seen
542 whether our techniques can also be applied to alternative types of layers, such as recurrent ones, and
543 the benefits they may bring.

544 Although we have provided a thorough investigation of the behaviour of our proposed system,
545 the only way we can control the end footprint of the model during training is via the λ_{gl} and ε_{ps}
546 hyperparameters. However, there is no guarantee about the final footprint of the model. If we are
547 willing to sacrifice accuracy, then the technique illustrated in Sec. 3.4 and evaluated in Sec. 5.5 is a
548 start. More robust ways of globally ranking per-layer importances are left as future work.

549 Lastly, our sampling method during training is uniform up to the maximum rank during progressive
550 shrinking. Although this method has proven effective, alternative sampling methods could potentially

551 accelerate rank exploration, thereby hastening the shrinking and convergence of the network during
 552 training.

553

554 C Extended Background

555 **Ordered Dropout.** Ordered Dropout is a technique of importance-based, nested and ordered pruning
 556 that works along the indices of a layer’s parameters (neurons, filters, etc.) Introduced by [17],
 557 the authors describe a training technique where a layer’s width is discretised in $|P|$ values, where
 558 $P = \{s_1, s_2, \dots, s_{|P|}\}$, and at each training step, they sample $p \sim U_P$ to get a specific subnetwork,
 559 extracted by selecting the first $\lceil p * K_l - 1 \rceil$ neurons per layer and dropping the rest. In contrast to
 560 our work, sampling is happening directly on model parameters (rather than ranks) and is uniform
 561 across layers (i.e. a single p-value is set). Nested-ness refers to the fact that larger p-value models
 562 include the parameters of lower p-values and importance-based pruning means that via stochastic
 563 sampling, the right-most (in terms of index) parameters train on progressively less data due to the
 564 probability of sampling and nestedness (i.e. all data pass from the parameters of minimal subnetwork,
 565 less pass the higher the p-value).

566 D Theoretical Properties of Low-Rank Layers

567 In this section, we show that for the case of linear mappings, i.e., the problem formulation discussed in
 568 (3), MAESTRO acts as PCA applied to the original dataset \mathcal{X} projected onto the space weighted by the
 569 corresponding singular values. Before proceeding with the theorem, we first recall the assumptions
 570 and notations introduced in the main paper.

571 We denote $C_{:b}$ as the first b columns of matrix C , $C_{:a,:b}$ denotes the first a rows, and b columns of a
 572 matrix C , $a+1 :$ denotes the all the columns/rows from index $a+1$, $:$ denotes the all the columns/rows,
 573 and for vectors, we use a single subscript. As discussed in the main paper, we reformulate the original
 574 least squares problems to the following decomposition problem

$$575 \min_{U \in \mathbb{R}^{m \times r}, V \in \mathbb{R}^{n \times r}} \mathbf{E}_{x, y \sim \mathcal{X}} \left[\mathbf{E}_{b \sim \mathcal{D}} \left[\|U_{:b} V_{:b}^\top x - y\|^2 \right] \right], \quad (6)$$

575 where \mathcal{D} is a distribution that samples $b \in \{1, 2, \dots, r\}$ with probability $p_b > 0$ and we assume that
 576 y is linked with x through linear map A , i.e., $y = Ax$.

577 **Theorem D.1.** *Let $A = \tilde{U} \tilde{\Sigma} \tilde{V}^\top$ be a SVD decomposition of A . Then, the minimization problem (6)
 578 is equivalent to PCA applied to the transformed dataset $x \rightarrow \tilde{\Sigma} \tilde{V}^\top x$, $x \sim \mathcal{X}$ projected on the column
 579 space of \tilde{U} . Concretely, we can first solve*

$$580 \min_{U \in \mathbb{R}^{m \times r}, V \in \mathbb{R}^{n \times r}} \mathbf{E}_{z \sim \mathcal{X}} \left[\mathbf{E}_{b \sim \mathcal{D}} \left[\left\| (U_{:b} V_{:b}^\top - I) \tilde{\Sigma} \tilde{V}^\top x \right\|^2 \right] \right], \quad (7)$$

580 and then we can obtain the solutions of (6) using $U^* = \tilde{U}^\top \bar{U}$, $V^* = \tilde{V}^\top \bar{V}$, where \bar{U}, \bar{V} belong to
 581 the set of optimal solutions of problem (7).

582 In the particular case, where \mathcal{X} is a uniform distribution on the unit ball, (6) recovers the best rank
 583 approximation of A across all ranks, i.e., up to the scale of U and V recovers truncated SVD. In the
 584 case, A is identity, (6) leads to standard PCA decomposition.

585 *Proof.* From the assumptions that $y = Ax$ and $A = \tilde{U} \tilde{\Sigma} \tilde{V}^\top$, we can rewrite (6) as

$$\min_{U \in \mathbb{R}^{m \times r}, V \in \mathbb{R}^{n \times r}} \mathbf{E}_{x \sim \mathcal{X}} \left[\mathbf{E}_{b \sim \mathcal{D}} \left[\left\| (U_{:b} V_{:b}^\top - \tilde{U} \tilde{\Sigma} \tilde{V}^\top) x \right\|^2 \right] \right].$$

586 Since \tilde{U} is orthogonal, we have $\|z\| = \|\tilde{U}^\top z\|$. Therefore, the above problem is equivalent to

$$\min_{U \in \mathbb{R}^{m \times r}, V \in \mathbb{R}^{n \times r}} \mathbf{E}_{x \sim \mathcal{X}} \left[\mathbf{E}_{b \sim \mathcal{D}} \left[\left\| (\tilde{U}^\top U_{:b} V_{:b}^\top - \tilde{\Sigma} \tilde{V}^\top) x \right\|^2 \right] \right],$$

587 which is also equivalent to

$$\min_{U \in \mathbb{R}^{m \times r}, V \in \mathbb{R}^{n \times r}} \mathbf{E}_{x \sim \mathcal{X}} \left[\mathbf{E}_{b \sim \mathcal{D}} \left[\left\| \left(U_{:b} V_{:b}^\top - \tilde{\Sigma} \tilde{V}^\top \right) x \right\|^2 \right] \right]$$

588 after reparametrization. The next step involves injecting identity in the form $\tilde{V} \tilde{V}^\top$ as that leads to the
589 equivalent reformulation

$$\min_{U \in \mathbb{R}^{m \times r}, V \in \mathbb{R}^{n \times r}} \mathbf{E}_{x \sim \mathcal{X}} \left[\mathbf{E}_{b \sim \mathcal{D}} \left[\left\| \left(U_{:b} V_{:b}^\top \tilde{V} - \tilde{\Sigma} \right) \tilde{V}^\top x \right\|^2 \right] \right].$$

590 As for the previous case, we can reparametrise the problem to obtain

$$\min_{U \in \mathbb{R}^{m \times r}, V \in \mathbb{R}^{n \times r}} \mathbf{E}_{x \sim \mathcal{X}} \left[\mathbf{E}_{b \sim \mathcal{D}} \left[\left\| \left(U_{:b} V_{:b}^\top - \tilde{\Sigma} \right) \tilde{V}^\top x \right\|^2 \right] \right].$$

591 Let $k = \text{rank}(\tilde{\Sigma}) = \text{rank}(A) \leq r$ and $z = \tilde{V}^\top x$. Furthermore, let $g = \tilde{\Sigma} z$ for any $z \in \mathbb{R}^n$, then
592 $g_{k+1:} = \vec{0}$. This, combined with the nested structure of the optimization problem, implies that the
593 optimal solution for U has to be of the form $u_{i,k+1:} = \vec{0}$ for all interesting (non-zero mapping)
594 directions, i.e., there exists $x \in \mathcal{X}$ such that $v_i^\top \tilde{V}^\top x \neq 0$. These are the only interesting solutions
595 since the case where for all $x \in \mathcal{X} : v_i^\top \tilde{V}^\top x = 0$ yields zero mapping on \mathcal{X} , which is not of interest
596 and could be dropped, e.g., using group lasso penalty discussed in the main part. Therefore, to solve
597 the original problem, we could first solve the following problem

$$\min_{U \in \mathbb{R}^{k \times r}, V \in \mathbb{R}^{n \times r}} \mathbf{E}_{z \sim \mathcal{X}} \left[\mathbf{E}_{b \sim \mathcal{D}} \left[\left\| \left(U_{:k,:b} V_{:b}^\top - \tilde{\Sigma}_{:k,:} \right) z \right\|^2 \right] \right]$$

598 and then reconstruct the corresponding solution of the original problem by appending zeros to the
599 resulting matrix U . By a similar argument, we can argue that for all non-zero mapping directions, it
600 has to be the case that $v_{i,k+1:} = \vec{0}$. Therefore, solving the original minimization reduces to

$$\min_{U \in \mathbb{R}^{k \times r}, V \in \mathbb{R}^{k \times r}} \mathbf{E}_{z \sim \mathcal{X}} \left[\mathbf{E}_{b \sim \mathcal{D}} \left[\left\| \left(U_{:b} V_{:b}^\top - \tilde{\Sigma}_{:k,:k} \right) z_{:k} \right\|^2 \right] \right].$$

601 This can be further simplified using reparametrization $V^\top \rightarrow V^\top \tilde{\Sigma}_{:k,:k}^{-1}$, which leads to

$$\min_{U \in \mathbb{R}^{k \times r}, V \in \mathbb{R}^{k \times r}} \mathbf{E}_{z \sim \mathcal{X}} \left[\mathbf{E}_{b \sim \mathcal{D}} \left[\left\| \left(U_{:b} V_{:b}^\top - I_k \right) \tilde{\Sigma}_{:k,:k} z_{:k} \right\|^2 \right] \right], \quad (8)$$

602 where I_k is $k \times k$ identity. If \mathcal{X} is centred around zero, then $\tilde{\Sigma}_{:k,:k} z_{:k}$ is also centred around zero,
603 and the above problem is up to scaling equivalent to PCA of $\tilde{\Sigma}_{:k,:k} z_{:k}$ as shown by Rippel et al. [42].
604 Since $\tilde{\Sigma}$ is a diagonal matrix with only $k \times k$ non-zero upper left sub-matrix, therefore, PCA on
605 $\tilde{\Sigma}_{:k,:k} z_{:k}$ is equivalent to PCA on $\tilde{\Sigma} z$ by appending zeros to the obtained principal component vectors.
606 Thus, we can write an equivalent formulation

$$\min_{U \in \mathbb{R}^{m \times r}, V \in \mathbb{R}^{n \times r}} \mathbf{E}_{z \sim \mathcal{X}} \left[\mathbf{E}_{b \sim \mathcal{D}} \left[\left\| \left(U_{:b} V_{:b}^\top - I \right) \tilde{\Sigma} \tilde{V}^\top x \right\|^2 \right] \right].$$

607 Furthermore, let \tilde{U}, \tilde{V} belong to the set of optimal solutions of problem (7). Then $U^* = \tilde{U}^\top \tilde{U}, V^* =$
608 $\tilde{V}^\top \tilde{V}$ belong to the set of optimal solutions of problem (6). This can be proved by reversing our
609 construction and ignoring scaling since (7) is scaling invariant.

610 For the case \mathcal{X} is a uniform distribution on the unit ball, we have $\tilde{\Sigma}_{:k,:k} z_{:k}$ is a k -dimensional ellipsoid
611 with principal axes being standard basis vectors $\{e_i\}_{i=1}^k$, where the length of the axes is given by
612 ordered singular values, i.e., the first basis vector corresponds to the largest singular vector. Therefore,
613 its principal component vectors correspond to the basis vectors. Following our construction, one can
614 see that the solution to the original problems leads to truncated SVD up to the scaling factor.

615 For the case A is an identity, we have $k = r = m = n$, $\tilde{\Sigma}$ is an identity, and $\tilde{U} = \tilde{V}$. Under this
616 setting, the principal component vectors obtained from (8) corresponds to principal component vectors
617 of \mathcal{X} in basis given by columns of \tilde{U} . Similarly, as in the previous case, reversing the transformations
618 to return back to the original problem, we conclude that the optimal solution of the original problem
619 corresponds to principal component vectors of \mathcal{X} since we reverse the transformation by \tilde{U}^\top . \square

620 **E Experimental setup**

621 **E.1 Datasets**

622 **MNIST.** The MNIST dataset [29] is a database of 28×28 greyscale handwritten digits, with a training
 623 set of 60k examples and a test set of 10k samples.

624 **CIFAR-10.** The CIFAR10 dataset [25] is a computer vision dataset that consists of 32×32 RGB
 625 images classified into 10 labels. It is split into 50k training images and 10k test images which are
 626 balanced across labels.

627 **WMT16.** The WMT dataset from statmt is machine translation dataset, spanning news commentaries
 628 and parliament proceedings, that aims to investigate the applicability of machine translation techniques
 629 when translating between language pairs. Specifically, we focus on the task of German-English
 630 language translation of image descriptions, commonly referred to as **Multi30k** [10]. We only utilise
 631 the text modality for the translation task. Data is taken straight from `torchtext`.

632 **TinyImagenet.** The TinyImagenet dataset [28] is a image classification challenge similar to
 633 ILSVRC [7]. The task it to classify an 64×64 RGB image among 200 classes, with each class
 634 having 500 training samples. The test set contains 10,000 images.

635 **E.2 Models**

636 **LeNet.** LeNet is a simple convolutional network, introduced by LeCun at al. for recognizing
 637 handwritten digits [29]. It consists of a sequence of two convolutional layers, followed by three
 638 fully-connected layers. However, we are using a ReLU instead of the initially proposed sigmoid
 639 activation. The detailed architecture of the network is depicted in Tab. 5

640 **ResNet.** ResNet [14] is a deep neural network whose prominent feature is the existence of skip (or
 641 residual) connections, that is connections that perform identity mappings merged with the target layer
 642 it joins with through summation. Multiple residual blocks are stacked to form the network. The result
 643 is an easier to optimise network that offers enhanced accuracy. We use ResNet-18 in our experiments,
 644 the architecture of which is depicted in Tab. 6, except for TinyImageNet, where we use ResNet-50.

Table 5: Detailed architecture of the LeNet-5 architecture used in our experiments. Each convolution and linear layer is followed by a ReLU activation that is omitted from the table. The shapes for convolution layers follows (m, n, k, k) .

Parameter	Shape	Layer hyper-parameter
layer1.conv1.weight	$1 \times 6 \times 5 \times 5$	stride:1;padding:1
pooling.max	N/A	kernel size:2;stride:1;dilation:1
layer2.conv2.weight	$6 \times 16 \times 5 \times 5$	stride:1;padding:0;dilation:1
pooling.max	N/A	kernel size:2;stride:2
layer3.fc1.weight	256×120	N/A
layer4.fc2.weight	120×84	N/A
layer5.fc3.weight	84×10	N/A

Table 6: The hybrid ResNet architecture for the CIFAR-10 and TinyImageNet datasets used in the experiments.

Layer Name	ResNet-18	ResNet-50
conv1	$3 \times 3, 64, \text{stride } 1, \text{padding } 1$	$7 \times 7, 64, \text{stride } 2, \text{padding } 1$
conv2_x	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 2$	$3 \times 3 \text{ maxpool, stride } 2$ $\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$
conv3_x	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 2$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$
conv4_x	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 2$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 6$
conv5_x	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 2$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$
	Avg Pool, 10-dim FC, SoftMax	Avg Pool, 20-dim FC, SoftMax

645 **VGG.** VGG [47] is a also a convolutional network that leverages smaller 3×3 convolutions that
 646 enables deeper architecture than before. For our experiments we are using VGG-19, the architecture
 647 of which is depicted in Tab. 7.

648 **Transformers.** The transformer architecture [51] has been lately revolutionising deep learning.
 649 Based on the notion of self-attention, for each input token, it produces a weighted combination of
 650 other relevant tokens weighed by the attention weight. Each attention unit has three weight matrices,
 651 namely W_Q, W_K, W_V , for query, key and value weights respectively producing the equivalent
 652 vectors. Attention is defined as the scaled dot product between key and query. For our translation
 653 task, we use the architecture depicted in Tab. 9.

Table 7: Detailed architecture of the VGG-19 architecture used in our experiments. There is a BatchNorm layer followed by a ReLU activation (omitted in the table) after each convolution layer. The shapes for convolution layers follows (m, n, k, k) .

Parameter	Shape	Layer hyper-parameter
layer1.conv1.weight	$3 \times 64 \times 3 \times 3$	stride:1;padding:1
layer2.conv2.weight	$64 \times 64 \times 3 \times 3$	stride:1;padding:1
pooling.max	N/A	kernel size:2;stride:2
layer3.conv3.weight	$64 \times 128 \times 3 \times 3$	stride:1;padding:1
layer4.conv4.weight	$128 \times 128 \times 3 \times 3$	stride:1;padding:1
pooling.max	N/A	kernel size:2;stride:2
layer5.conv5.weight	$128 \times 256 \times 3 \times 3$	stride:1;padding:1
layer6.conv6.weight	$256 \times 256 \times 3 \times 3$	stride:1;padding:1
layer7.conv7.weight	$256 \times 256 \times 3 \times 3$	stride:1;padding:1
layer8.conv8.weight	$256 \times 256 \times 3 \times 3$	stride:1;padding:1
pooling.max	N/A	kernel size:2;stride:2
layer9.conv9.weight	$256 \times 512 \times 3 \times 3$	stride:1;padding:1
layer10.conv10.weight	$512 \times 512 \times 3 \times 3$	stride:1;padding:1
layer11.conv11.weight	$512 \times 512 \times 3 \times 3$	stride:1;padding:1
layer12.conv12.weight	$512 \times 512 \times 3 \times 3$	stride:1;padding:1
pooling.max	N/A	kernel size:2;stride:2
layer13.conv13.weight	$512 \times 512 \times 3 \times 3$	stride:1;padding:1
layer14.conv14.weight	$512 \times 512 \times 3 \times 3$	stride:1;padding:1
layer15.conv15.weight	$512 \times 512 \times 3 \times 3$	stride:1;padding:1
layer16.conv16.weight	$512 \times 512 \times 3 \times 3$	stride:1;padding:1
pooling.avg	N/A	kernel size:2
classifier.weight	512×10	N/A
classifier.bias	10	N/A

Table 8: Detailed information of the encoder layer in the Transformer architecture in our experiment

Parameter	Shape	Hyper-param.
embedding	9521×512	padding index: 1
positional encoding	N/A	N/A
dropout	N/A	$p = 0.1$
encoder.self-attention.wq(W^Q)	512×512	N/A
encoder.self-attention.wk(W^K)	512×512	N/A
encoder.self-attention.wv(W^V)	512×512	N/A
encoder.self-attention.wo(W^O)	512×512	N/A
encoder.self-attention.dropout	N/A	$p = 0.1$
encoder.self-attention.layernorm	512	$\epsilon = 10^{-6}$
encoder.ffn.layer1	512×2048	N/A
encoder.ffn.layer2	2048×512	N/A
encoder.layernorm	512	$\epsilon = 10^{-6}$
dropout	N/A	$p = 0.1$

Table 9: Detailed information of the decoder layer in the Transformer architecture in our experiment

Parameter	Shape	Hyper-param.
embedding	9521×512	padding index: 1
positional encoding	N/A	N/A
dropout	N/A	$p = 0.1$
decoder.self-attention.wq(W^Q)	512×512	N/A
decoder.self-attention.wk(W^K)	512×512	N/A
decoder.self-attention.wv(W^V)	512×512	N/A
decoder.self-attention.wo(W^O)	512×512	N/A
decoder.self-attention.dropout	N/A	$p = 0.1$
decoder.self-attention.layernorm	512	$\epsilon = 10^{-6}$
decoder.enc-attention.wq(W^Q)	512×512	N/A
decoder.enc-attention.wk(W^K)	512×512	N/A
decoder.enc-attention.wv(W^V)	512×512	N/A
decoder.enc-attention.wo(W^O)	512×512	N/A
decoder.enc-attention.dropout	N/A	$p = 0.1$
decoder.enc-attention.layernorm	512	$\epsilon = 10^{-6}$
decoder.ffn.layer1	512×2048	N/A
decoder.ffn.layer2	2048×512	N/A
decoder.layernorm	512	$\epsilon = 10^{-6}$
dropout	N/A	$p = 0.1$

654 **E.3 Hyperparameter selection**

655 **LeNet.** We use a standard configuration that is commonly employed for training LeNet models — a
656 step size of 0.01, a momentum of 0.9, and no weight decay. We train for a total of 20 epochs.

657 **VGG and ResNet-18.** Similarly, we use a standard configuration that is commonly employed for
 658 training VGG and ResNet-18 models — a step size of 0.01, a momentum of 0.9, weight decay of
 659 $1e^{-4}$, and a learning schedule with step size reductions by a factor of 10 at epochs 150 and 250. We
 660 train for a total of 300 epochs.

661 **ResNet-50.** Similarly, we use a standard configuration that is commonly employed for training
 662 ResNet-50 models — a step size of 0.01, a momentum of 0.9, weight decay of $1e^{-4}$, and a learning
 663 schedule with step size reductions by a factor of 10 at epochs 30 and 60. We train for a total of 90
 664 epochs.

665 **Transformers.** For the Transformer model, we use the Adam optimizer with an initial learning rate
 666 at 0.001, $\beta s = (0.9, 0.98)$, $\varepsilon = 10^{-8}$ batch size at 256. We also conduct gradient norm clipping with
 667 norm bound at 0.25. The entire training takes 400 epochs. For the vanilla warm-up training, we use
 668 warm-up epoch $E_{wu} = 10$. We enable label smoothing, weight sharing for the source and target
 669 word embedding, and weight sharing between target word embedding and the last dense layer. The
 670 learning rate schedule follows directly from the one proposed [51].

671 E.4 Deciding against decomposition

672 During inference, if the rank of a given layer is so large that keeping it as a non-decomposed layer is
 673 more efficient, we opt not to decompose that particular layer.

674 F Extended evaluation

675 F.1 MAESTRO recovers correct ordering

676 In the main text, we pointed out that SVD fails to consider data. Indeed, even in the case of linear NN,
 677 the acquired singular vectors may exhibit incorrect ordering. To illustrate this problem, we provide a
 678 simple example in which we use a matrix A with a rank of 3. We organize the dataset \mathcal{X} such that the
 679 third singular vector has the highest importance, followed by the second and then the first singular
 680 vector in decreasing order of significance. It is clear that SVD doesn't consider the data, and as a
 681 result, it cannot comprehend this behavior. Below (in Fig. 5), we demonstrate how MAESTRO is able
 682 to correctly discern the order.

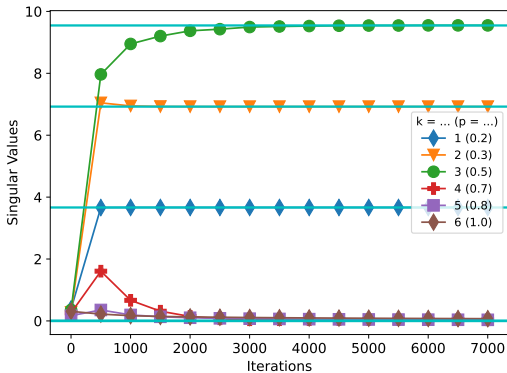


Figure 5: Verification that MAESTRO recovers correct order of importance. Target mapping is of rank 3, and the dataset is constructed in such a way that the singular vectors have reversed the order of importance. p and k stand for relative and actual rank, respectively.

683 F.2 Training behaviour of MAESTRO

684 For completeness, we also include an extended version of Fig. 3 from the main paper, where we
 685 presented the training dynamics for MAESTRO. Fig. 6, 7 and 8 present similar plots, but across both
 686 MNIST and CIFAR-10. Specifically, Fig. 6 illustrates the evolution of total rank throughout the
 687 training steps. We observe that the ranks are pruned incrementally. This aligns with the observations
 688 made during Pufferfish [53] training, where the authors suggested warm-start training with full

689 precision to enhance the final model performance. In our case, the necessity to implement this
 690 heuristic is avoided, as MAESTRO prunes rank automatically. Fig. 7 demonstrates the ranks across
 691 layers post-training. An intriguing trend is observed: the ranks are nested for increasing λ_{gl} ,
 692 suggesting a potential inherent ordering of ranks not only within each layer but also possibly a global
 693 one. We provide a preliminary exploration of this fascinating pattern in the subsequent section and
 694 intend to probe it more deeply in future studies. We believe this may enhance the rank selection
 695 and sampling process. Finally, Fig. 8 portrays the evolution of the training loss. Our premise that
 696 sampling does not negatively affect training is validated by empirical performance.

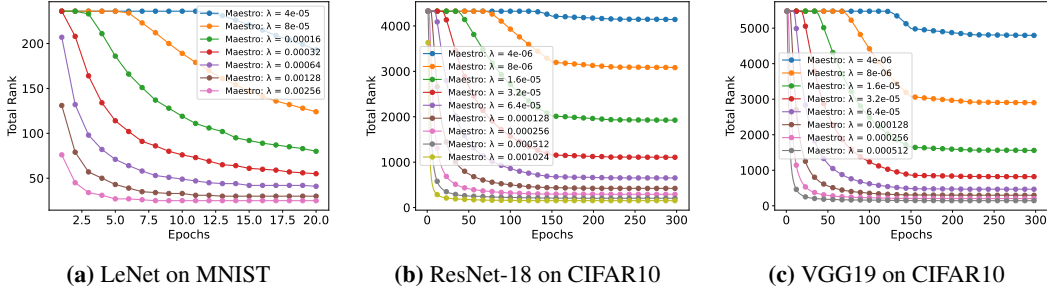


Figure 6: Total rank ($\sum_{i=1}^d r_i$) progression during training.

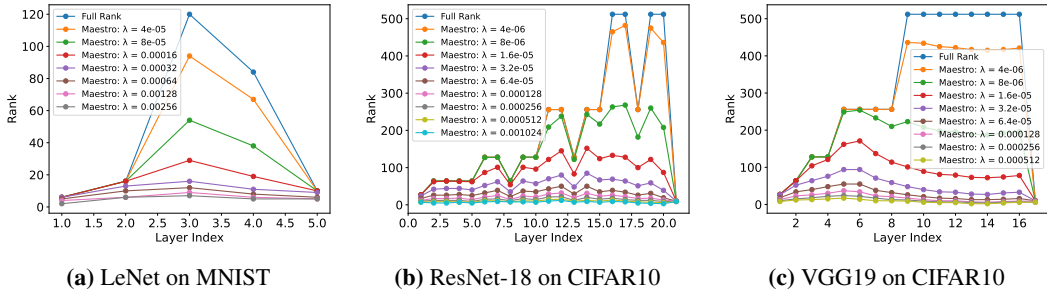


Figure 7: Ranks r_i 's across different layers after training.

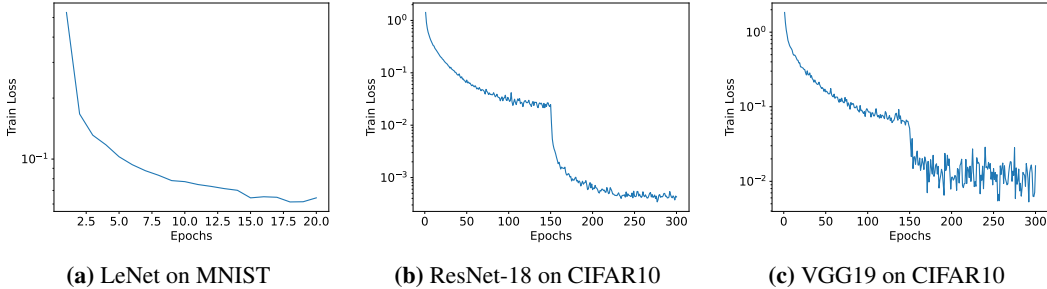


Figure 8: Convergence of MAESTRO with $\lambda_{gl} = 0$.

697 F.3 Model size-accuracy trade-off at training and deployment time

698 In addition to the original illustrations, we present an extended interpretation of Fig. 4, where we
 699 depict diverse strategies to maintain a balance between model size and accuracy in the process
 700 of model training and deployment. In Fig. 9, we demonstrate the effective pruning of MAESTRO
 701 ($\lambda_{gl} = 0$) for deployment, utilizing the greedy search methodology discussed in Section 3.4. This is
 702 juxtaposed with the greedy pruning of a model not originally factorized but later factorized through
 703 SVD. Our findings reveal that this straightforward baseline does not match the performance of
 704 MAESTRO's learned decomposition, leading to a considerable performance drop.

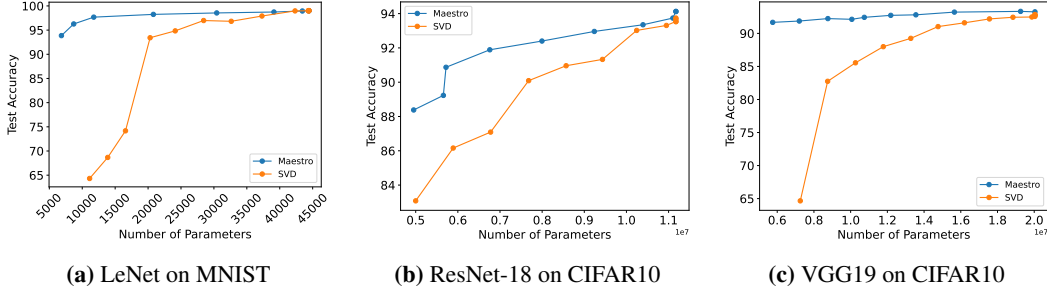


Figure 9: Accuracy-latency trade-off comparing MAESTRO ($\lambda_{gl}=0$) and SVD.

705 Subsequently, Fig. 10 displays the end accuracy and the count of model parameters corresponding to
 706 various hierarchical group lasso penalties. This results in an optimal compromise between latency
 707 and accuracy for both the training and inference stages. It’s worth noting, though, that each model
 708 was trained separately, in contrast to greedy pruning, which demands just a single training round.
 709 Additionally, we scrutinize the training expense for each model illustrated in Fig. 10, the results of
 710 which are exhibited in Tables 10, 11, 12, 13 and 14, where we display and the final accuracy of the
 711 model, MACs and the number of parameters for inference, and relative total training cost in terms of
 712 the number of model parameters and MACs compared to the non-factorized model. Interestingly,
 713 smaller models are not only advantageous in terms of inference efficiency, but they can also be trained
 714 at a small portion of the cost required by full-rank models. On the downside, the smallest models
 715 cause a non-negligible reduction in performance.

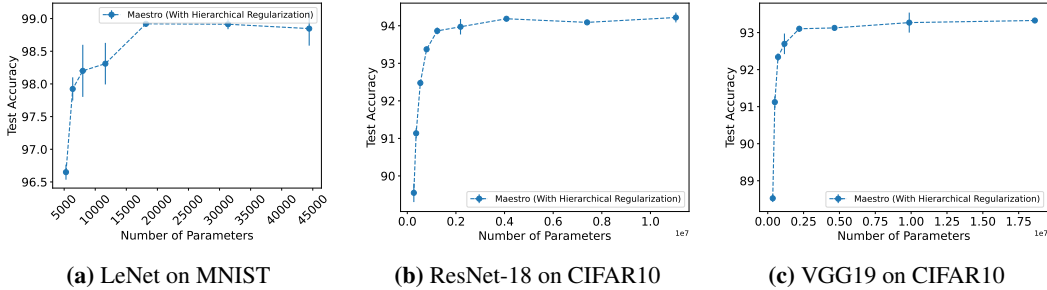


Figure 10: Impact of hierarchical group lasso on the accuracy-memory trade-off. Exact values are provided in Tables 10, 11 and 12, respectively.

Table 10: LeNet performance on MNIST for different regularization parameters. The last column in the table displays the relative total training cost in terms of the number of Multiply-Accumulate operations (MACs) and model parameters, compared to the non-factorized model.

Variant	Acc. (%)	MACs (Inf.)	Params. (Inf.)	Rel. MACs / Params. (Train.)
Non-Factorized	98.99 ± 0.09	281640 ± 0 (1.00 \times)	44426 ± 0 (1.00 \times)	$1.00 \times / 1.00 \times$
MAESTRO ($\lambda_{gp} = 0.$)	99.06 ± 0.09	281640 ± 0 (1.00 \times)	44426 ± 0 (1.00 \times)	$1.14 \times / 1.49 \times$
MAESTRO ($\lambda_{gp} = 8e^{-5}$)	98.91 ± 0.09	268577 ± 389 (0.95 \times)	31363 ± 0 (0.71 \times)	$1.08 \times / 1.14 \times$
MAESTRO ($\lambda_{gp} = 16e^{-5}$)	98.92 ± 0.05	255369 ± 217 (0.91 \times)	44426 ± 217 (0.41 \times)	$1.06 \times / 0.80 \times$
MAESTRO ($\lambda_{gp} = 32e^{-5}$)	98.31 ± 0.39	237084 ± 6268 (0.84 \times)	18155 ± 271 (0.26 \times)	$0.93 \times / 0.53 \times$
MAESTRO ($\lambda_{gp} = 64e^{-5}$)	98.20 ± 0.49	178165 ± 19098 (0.63 \times)	7996 ± 662 (0.18 \times)	$0.77 \times / 0.33 \times$
MAESTRO ($\lambda_{gp} = 128e^{-5}$)	97.92 ± 0.22	131789 ± 8965 (0.47 \times)	6375 ± 77 (0.14 \times)	$0.54 \times / 0.21 \times$
MAESTRO ($\lambda_{gp} = 256e^{-5}$)	96.65 ± 0.14	99969 ± 6252 (0.35 \times)	5293 ± 214 (0.12 \times)	$0.39 \times / 0.14 \times$

716 Lastly, we delve deeper into the observation of nested ranks with increasing λ_{gl} . Fig. 11 outlines
 717 the performance of MAESTRO ($\lambda_{gl} = 0$) across various ranks chosen by smaller models MAESTRO
 718 ($\lambda_{gl} > 0$). We observe that MAESTRO ($\lambda_{gl} = 0$) delivers impressive results—for example, we can
 719 reduce its parameters by 10x for VGG while preserving an accuracy of 87.7% without any fine-tuning
 720 simply by leveraging rank structure from separate runs. For LeNet, a reduction in model size by a
 721 factor of three is achievable without sacrificing accuracy. Last, for ResNet-18 the reduction is 1.7 \times .
 722 As highlighted earlier, we aim to delve deeper into this subject in future studies.

Table 11: ResNet-18 performance on CIFAR10 for different regularization parameters. The last column in the table displays the relative total training cost in terms of the number of Multiply-Accumulate operations (MACs) and model parameters, compared to the non-factorized model.

Variants	Acc. (%)	GMACs (Inf.)	Params. (M) (Inf.)	Rel. MACs / Params. (Train.)
Non-Factorized	93.86 \pm 0.20	0.56 \pm 0 (1.00 \times)	11.2 \pm 0 (1.00 \times)	1.00 \times / 1.00 \times
MAESTRO ($\lambda_{gp} = 0.$)	94.04 \pm 0.10	0.56 \pm 0 (1.00 \times)	11.2 \pm 0 (1.00 \times)	1.10 \times / 1.13 \times
MAESTRO ($\lambda_{gp} = 4e^{-6}$)	94.22 \pm 0.16	0.55 \pm 0.0047 (1.00 \times)	11.1 \pm 0.030 (0.99 \times)	1.09 \times / 1.10 \times
MAESTRO ($\lambda_{gp} = 8e^{-6}$)	94.09 \pm 0.01	0.49 \pm 0.0002 (0.89 \times)	7.41 \pm 0.004 (0.66 \times)	1.00 \times / 0.85 \times
MAESTRO ($\lambda_{gp} = 16e^{-6}$)	94.19 \pm 0.07	0.39 \pm 0.0008 (0.70 \times)	4.08 \pm 0.020 (0.37 \times)	0.83 \times / 0.58 \times
MAESTRO ($\lambda_{gp} = 32e^{-6}$)	93.97 \pm 0.25	0.25 \pm 0.0013 (0.45 \times)	2.19 \pm 0.007 (0.20 \times)	0.60 \times / 0.36 \times
MAESTRO ($\lambda_{gp} = 64e^{-6}$)	93.86 \pm 0.11	0.15 \pm 0.0006 (0.27 \times)	1.23 \pm 0.004 (0.11 \times)	0.39 \times / 0.22 \times
MAESTRO ($\lambda_{gp} = 128e^{-6}$)	93.37 \pm 0.07	0.094 \pm 0.0006 (0.17 \times)	0.79 \pm 0.009 (0.07 \times)	0.25 \times / 0.13 \times
MAESTRO ($\lambda_{gp} = 256e^{-6}$)	92.48 \pm 0.04	0.064 \pm 0.0002 (0.12 \times)	0.54 \pm 0.006 (0.05 \times)	0.16 \times / 0.08 \times
MAESTRO ($\lambda_{gp} = 512e^{-6}$)	91.14 \pm 0.16	0.044 \pm 0.0004 (0.08 \times)	0.37 \pm 0.007 (0.03 \times)	0.11 \times / 0.05 \times
MAESTRO ($\lambda_{gp} = 1024e^{-6}$)	89.55 \pm 0.30	0.032 \pm 0.0002 (0.06 \times)	0.27 \pm 0.007 (0.02 \times)	0.07 \times / 0.03 \times

Table 12: VGG19 performance on CIFAR10 for different regularization parameters. The last column in the table displays the relative total training cost in terms of the number of Multiply-Accumulate operations (MACs) and model parameters, compared to the non-factorized model.

Variants	Acc. (%)	GMACs (Inf.)	Params. (M) (Inf.)	Rel. MACs / Params. (Train.)
Non-Factorized	92.94 \pm 0.17	0.40 \pm 0 (1.00 \times)	20 \pm 0 (1.00 \times)	1.00 \times / 1.00 \times
MAESTRO ($\lambda_{gp} = 0.$)	93.06 \pm 0.17	0.40 \pm 0 (1.00 \times)	20 \pm 0 (1.00 \times)	1.10 \times / 1.12 \times
MAESTRO ($\lambda_{gp} = 4e^{-6}$)	93.33 \pm 0.08	0.39 \pm 0.0017 (0.97 \times)	18.8 \pm 0 (0.94 \times)	1.06 \times / 1.04 \times
MAESTRO ($\lambda_{gp} = 8e^{-6}$)	93.27 \pm 0.33	0.30 \pm 0.0017 (0.76 \times)	9.91 \pm 0.008 (0.49 \times)	0.90 \times / 0.73 \times
MAESTRO ($\lambda_{gp} = 16e^{-6}$)	93.13 \pm 0.07	0.21 \pm 0.0014 (0.53 \times)	4.66 \pm 0.052 (0.23 \times)	0.69 \times / 0.46 \times
MAESTRO ($\lambda_{gp} = 32e^{-6}$)	93.10 \pm 0.10	0.13 \pm 0.0009 (0.33 \times)	2.20 \pm 0.025 (0.11 \times)	0.47 \times / 0.27 \times
MAESTRO ($\lambda_{gp} = 64e^{-6}$)	92.70 \pm 0.34	0.08 \pm 0.0005 (0.20 \times)	1.17 \pm 0.010 (0.06 \times)	0.30 \times / 0.16 \times
MAESTRO ($\lambda_{gp} = 128e^{-6}$)	92.34 \pm 0.12	0.05 \pm 0.0005 (0.13 \times)	0.72 \pm 0.002 (0.04 \times)	0.19 \times / 0.09 \times
MAESTRO ($\lambda_{gp} = 256e^{-6}$)	91.12 \pm 0.19	0.04 \pm 0.0007 (0.09 \times)	0.50 \pm 0.023 (0.02 \times)	0.12 \times / 0.05 \times
MAESTRO ($\lambda_{gp} = 512e^{-6}$)	88.53 \pm 0.13	0.03 \pm 0.0003 (0.06 \times)	0.35 \pm 0.003 (0.02 \times)	0.08 \times / 0.03 \times

Table 13: Transformer performance on Multi30k for different regularization parameters. The last column in the table displays the relative total training cost in terms of the number of Multiply-Accumulate operations (MACs) and model parameters, compared to the non-factorized model.

Variants	Acc. (%)	Ppl.	GMACs (Inf.)	Params. (M) (Inf.)	Rel. MACs / Params. (Train.)
Non-Factorized	65.33 \pm 1.13	9.85 \pm 0.10	1.370 \pm 0.0000 (1.00 \times)	53.9 \pm 0.000 (1.00 \times)	1.00 \times / 1.00 \times
MAESTRO ($\lambda_{gp} = 0.32$)	61.30 \pm 0.26	12.99 \pm 0.31	1.125 \pm 0.0030 (0.82 \times)	45.1 \pm 0.101 (0.84 \times)	1.03 \times / 1.14 \times
MAESTRO ($\lambda_{gp} = 0.64$)	63.78 \pm 0.14	9.37 \pm 0.32	0.957 \pm 0.0112 (0.70 \times)	39.1 \pm 0.413 (0.73 \times)	0.95 \times / 1.05 \times
MAESTRO ($\lambda_{gp} = 1.28$)	66.14 \pm 0.08	7.02 \pm 0.17	0.570 \pm 0.0088 (0.42 \times)	25.3 \pm 0.315 (0.47 \times)	0.75 \times / 0.86 \times
MAESTRO ($\lambda_{gp} = 2.56$)	66.08 \pm 0.09	6.90 \pm 0.07	0.248 \pm 0.0032 (0.18 \times)	13.8 \pm 0.113 (0.26 \times)	0.47 \times / 0.58 \times
MAESTRO ($\lambda_{gp} = 5.12$)	57.70 \pm 0.13	13.97 \pm 0.43	0.123 \pm 0.0002 (0.09 \times)	9.3 \pm 0.001 (0.17 \times)	0.28 \times / 0.39 \times

Table 14: ResNet50 performance on Tiny-Imagenet-200 for different regularization parameters. The last column in the table displays the relative total training cost in terms of the number of Multiply-Accumulate operations (MACs) and model parameters, compared to the non-factorized model.

Variants	Acc. (%)	GMACs (Inf.)	Params. (M) (Inf.)	Rel. MACs / Params. (Train.)
Non-Factorized	61.74 \pm 0.27	5.19 \pm 0.0000 (1.00 \times)	23.9 \pm 0.000 (1.00 \times)	1.22 \times / 1.22 \times
MAESTRO ($\lambda_{gp} = 0.$)	61.05 \pm 0.09	5.19 \pm 0.0000 (1.00 \times)	23.9 \pm 0.000 (1.00 \times)	1.21 \times / 1.20 \times
MAESTRO ($\lambda_{gp} = 4e^{-5}$)	60.13 \pm 0.34	4.72 \pm 0.0013 (0.91 \times)	18.8 \pm 0.017 (0.79 \times)	0.81 \times / 0.69 \times
MAESTRO ($\lambda_{gp} = 8e^{-5}$)	59.20 \pm 0.40	3.01 \pm 0.0064 (0.58 \times)	9.64 \pm 0.023 (0.40 \times)	0.00 \times / 0.00 \times
MAESTRO ($\lambda_{gp} = 16e^{-5}$)	58.35 \pm 0.40	1.49 \pm 0.0142 (0.29 \times)	4.48 \pm 0.022 (0.19 \times)	0.61 \times / 0.54 \times
MAESTRO ($\lambda_{gp} = 32e^{-5}$)	56.52 \pm 0.08	0.72 \pm 0.0022 (0.14 \times)	2.25 \pm 0.013 (0.09 \times)	0.51 \times / 0.47 \times

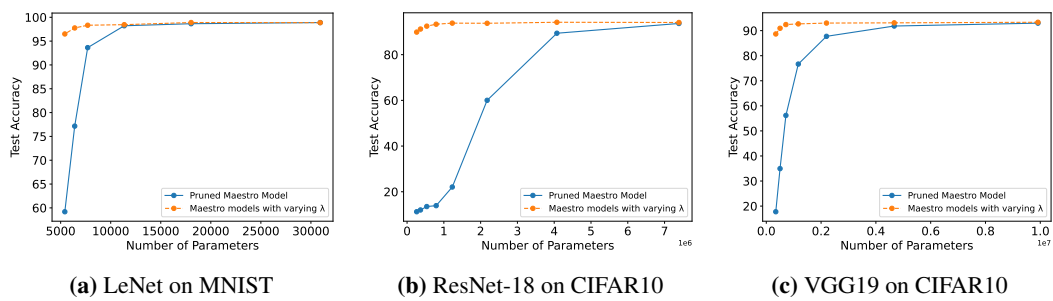


Figure 11: MAESTRO with progressive pruning to showcase nested rank importance structure. The original model corresponds to an evaluation in Fig. 10, and pruned models are based on MAESTRO with $\lambda_{gl} = 0$, and they are pruned using the same ranks as selected by MAESTRO with $\lambda_{gl} > 0$.