000 COMPRESS THEN SERVE: SERVING THOUSANDS OF LORA ADAPTERS WITH LITTLE OVERHEAD

Anonymous authors

Paper under double-blind review

ABSTRACT

Fine-tuning large language models (LLMs) with low-rank adaptations (LoRAs) has become common practice, often yielding numerous copies of the same LLM differing only in their LoRA updates. This paradigm presents challenges for systems that serve real-time responses to queries that each involve a different LoRA. Prior works optimize the design of such systems but still require continuous loading and offloading of LoRAs, as it is infeasible to store thousands of LoRAs in GPU memory. To mitigate this issue, we investigate the efficacy of model compression when serving LoRAs. We propose a method for joint compression of LoRAs into a shared basis paired with LoRA-specific scaling matrices. We extend our algorithm to learn clusters of LoRAs that are more amenable to joint compression, allowing it to scale gracefully to large LoRA collections. Our experiments with up to 500 LoRAs demonstrate that compressed LoRAs preserve performance while offering major throughput gains in realistic serving scenarios with over a thousand LoRAs, maintaining 80% of the throughput of serving a single LoRA.

025 026

027

001

002 003 004

006 007 008

009 010

011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

028 The myriad uses for foundation models (FMs) have led to a proliferation of specialized models, each 029 fine-tuned to perform a downstream task. To avoid fine-tuning foundation models with billions of parameters, more parameter-efficient fine-tuning (PEFT) algorithms were proposed. An especially 031 successful PEFT method is low-rank adaptation (LoRA) (Hu et al., 2021), which learns low-rank 032 additive changes to neural network matrices. Because of the low-rank parameterization, these ma-033 trices (called adapter weights) contain orders-of-magnitude fewer parameters than the base model. 034 Still, LoRA can achieve performance on par with full fine-tuning (Hu et al., 2021).

LoRA's popularity has triggered a growing need to serve large collections of LoRA adapters at scale. For example, proprietary and open-source LLM providers offer fine-tuning services (OpenAI, 037 2024; TogetherAI, 2024; Predibase, 2024) with user bases likely in thousands or even hundreds of 038 thousands. As each user wants to use their own fine-tuned version of the LLM, simply serving a dedicated fine-tuned LLM per user becomes infeasible. To this end, S-LoRA (Sheng et al., 2023) considers a system where only the base LLM is placed on an inference server and individual LoRA 040 adapters are switched as needed at inference time. S-LoRA optimizes the system's inner workings 041 via custom CUDA kernels and memory management to increase the throughput when serving multi-042 ple LoRAs. Such multi-LoRA system design has also been adopted in vLLM (Kwon et al., 2023), a 043 state-of-the-art LLM serving engine. Despite the optimized system design, serving LoRAs still has 044 a fundamental limitation: when the number of adapters is large, they need to be constantly loaded 045 and offloaded from GPU memory to accommodate incoming requests, degrading throughput. 046

The problem of accommodating multiple LoRA adapters is also apparent when placing LLMs on 047 edge devices, where smaller LLMs are fine-tuned for various tasks and the adapters are swapped 048 depending on the task at hand (Gunter et al., 2024). In this case, the amount of adapters is smaller, 049 e.g., a few dozen (Gunter et al., 2024), but the memory constraints are also tighter due to the limited capacity of edge devices. 051

In this work, we consider the problem of compressing a collection of LoRAs. We have two key 052 objectives: (1) preserving the performance of the original LoRAs and (2) improving the throughput of serving many LoRAs. We formulate LoRA compression as a reconstruction problem, where the

054 goal is to approximate the original adapters via collections of matrices of a smaller total size. We 055 investigate an approach based on compressing LoRAs jointly by finding a shared basis and LoRA-056 specific scaling matrices, and propose a joint diagonalization-based algorithm (JD). To improve 057 reconstruction error for large numbers of LoRAs while keeping the number of parameters in check, 058 we propose a clustering approach where each cluster is compressed independently using the joint diagonalization algorithm. Our clustering algorithm is based on alternating between optimizing the cluster assignments and the per-cluster reconstruction error. 060

061 We showcase the benefits of joint compression in 062 Figure 1. We chose JD configurations that match the 063 constraints. When serving up to 64 unique LoRAs 064 we use JD without clustering and for 128 or more, we pick the number of clusters to match the perfor-065 mance of compressed and original LoRAs. In each 066 case, the GPU memory footprint of the compressed 067 and original LoRAs is matched for a fair compar-068 ison to the vLLM's multi-LoRA inference engine. 069 When serving over 1000 LoRAs, compression increases throughput by a factor of 1.6 and maintains 071 80% of the throughput of serving the base LLM (or a single LoRA merged into the LLM). Detailed results 073 are presented in Section 6.



Figure 1: Throughput gains when serving 1000s of compressed LoRAs with vLLM.

We summarize our main contributions below: 075

- We formulate the problem of compressing a collection of LoRAs and propose a joint compression scheme based on joint diagonalization.
- For large numbers of LoRAs, we scale the joint compression scheme by proposing a clustering algorithm where each cluster is jointly compressed to minimize the overall reconstruction error.
- We establish theoretical guarantees for the reconstruction error central to our compression formulation and verify the relation between reconstruction loss and performance empirically.
- We train a collection of 500 high-quality LoRAs for Mistral-7B-Instruct-v0.2 (Jiang et al., 2023a) on 500 natural instruction tasks (Wang et al., 2022) and demonstrate that our compression techniques preserve the performance of the original LoRAs. We will release the 500 LoRAs to facilitate future work on LoRA compression as well as the code for our method.
- We incorporate LoRA compression into a state-of-the-art LLM serving system and demonstrate that it is possible to serve over 1000 LoRAs across thousands of asynchronous requests with throughput comparable to serving a single LoRA.
- 880 090

074

076

077

079

081

082

083

084

085

087

RELATED WORK 2

091 092

095

Parameter-efficient fine-tuning (PEFT) has become prevalent for updating foundation models thanks to the need for efficiency in training and communication (Lialin et al., 2023). Many PEFT methods have been proposed, e.g. (Houlsby et al., 2019; Liu et al., 2022) and LoRA (Hu et al., 2021) became 094 the standard, partially due to the ease of switching between LoRAs in inference time.

104 Punica (Chen et al., 2023) introduces Segmented Gather Matrix-Vector Multiplication (SGMV) to 105 optimize multi-LoRA serving by parallelizing feature-weight multiplications in batches and grouping requests that utilize the same LoRA model. Our approach, by contrast, emphasizes parameter 106 reduction as a means to efficiently serve multiple LoRAs, providing an orthogonal strategy that 107 can be seamlessly integrated with Punica's methods to further enhance performance. In our vLLM

⁰⁹⁶ Several works improve LoRA (Liu et al., 2024; Wang et al., 2024), sometimes with algebraic methods like SVD (Meng et al., 2024; Zhang et al., 2023; Jiang et al., 2023b) or by leveraging its statis-098 tical properties (Zhu et al., 2024; Zeng & Lee, 2024). Relatively few, however, accelerate inference times. S-LoRA (Sheng et al., 2023) provides an efficient means of switching between LoRAs. (Wen & Chaudhuri, 2024) adapts training to reduce batch multiplications, accelerating inference. Our 100 method achieves a similar outcome (see Appendix D) without changing the LoRA formulation or 101 requiring that LoRAs be trained in a dedicated way; future improvements to LoRA will also benefit 102 from this aspect of our work (e.g., Meng et al. (2024)). 103

experiments, we leveraged the Punica kernel for multi-LoRA implementation, demonstrating the application of our method in conjunction with Punica's optimizations.

There are many efforts to compress models (Cheng et al., 2017; Gholami et al., 2022; Sharma et al., 2024; Li et al., 2018)—including some specifically for LoRAs—to accelerate inference. Predominantly, pruning and sparsification methods delete weights (Yadav et al., 2023a), and quantization methods reduce the weights' precision (Dettmers et al., 2024). Some works compress weights to reduce model size but typically require decompression and hence do not save GPU memory (Hershcovitch et al., 2024). Similarly to our work, while most methods increase speed at the cost of performance, a few note increased performance and generalization after compression (Yadav et al., 2023a; Nadjahi et al., 2023; Hershcovitch et al., 2024; Sharma et al., 2024).

Our work also relates to model merging (Choshen et al., 2022; Wortsman et al., 2022; Matena & Raffel, 2021) and mixture of experts methods (Muqeeth et al., 2024; Yadav et al., 2024). These methods reuse models trained by others (Choshen et al., 2023; Raffel, 2023), serving them together as one compressed model. Despite this similarity, these methods create a single general model that acts on any input, while our model allows for more performant per-task solutions.

123 124 125

126

3 RANK-BASED LORA COMPRESSION

LoRA updates are parameterized by pairs of matrices A, B, whose product BA updates the fixed weight matrices $W_0 \in \mathbb{R}^{d_B \times d_A}$ of a neural network foundation model. Given an input x to a layer, the output of the LoRA-updated model at this layer is $(W_0 + BA)x$.

In formulating our compression algorithms, we consider a collection of given LoRA adapters $\{(A_i, B_i)\}_{i=1}^n$ that we would like to serve. We let r_i refer to the rank of the LoRA adapter-pair (A_i, B_i) , i.e., $B_i \in \mathbb{R}^{d_B \times r_i}$, $A_i \in \mathbb{R}^{r_i \times d_A}$.

While our compression technique has access only to a collection of $\{(A_i, B_i)\}_{i=1}^n$ pairs, in our experiments we will assess the efficacy of compression by comparing how the compressed matrices perform relative to the uncompressed LoRAs on typical data. For this reason, although in this section we optimize a Frobenius norm reconstruction error relative to the product B_iA_i , in reality this is a proxy for the nonlinear and complex way that compression errors in the adapters impact transformer performance. Our experimental evaluation will thus focus on the performance of the compressed LoRAs against the uncompressed versions on real data in §6.

Our compression methods significantly reduce the overall number of parameters. Reducing parameters through compression theoretically accelerates storage and serving processes for a collection of LoRAs. This reduction, however, alters the computational dynamics during inference, so parameter reduction alone does not immediately imply faster throughput. In light of the complexities of GPU optimization, we experimentally assess the throughput under realistic conditions in §6.3.

146

147 3.1 JOINT DIAGONALIZATION148

For compression to scale to large numbers of LoRAs, the compressed number of parameters should not scale linearly with *n*. Hence seeking to compress each LoRA individually (e.g., via SVD as detailed in the experimental baselines) is inherently limited.

To address this, we suggest a Joint Diagonalization (JD) method, which optimizes a shared basis onto which we can project the set of n LoRAs. This will allow structure to be shared, implicitly grouping and/or merging the collection of LoRAs.

In this model, each LoRA product $B_i A_i$ is factorized into the form $U\Sigma_i V$, where U and V are shared across all LoRAs and Σ_i is specific to each LoRA. In this formulation, every Σ_i shares the same rank r. This allows U and V to be pre-loaded onto the GPU, with Σ_i loaded when necessary for each batch. The matrices Σ_i can be either diagonal or small square matrices, accelerating the forward pass compared to conventional multi-LoRA serving configurations.

161 Objective function. Motivated by the relationship of singular value decomposition to minimizing the Frobenius norm of the reconstruction error, we also propose to minimize the Frobenius norm of

the adapter matrix approximation error. Specifically, we use the following objective function:

$$\min_{\{\Sigma_i\}_{i=1}^n, U, V} \sum_{i=1}^n \|B_i A_i - U \Sigma_i V^\top\|_{\text{Fro}}^2.$$
(1)

Note this problem is *not* solved by a single matrix SVD, since U and V are shared among all terms but the Σ_i 's are not. Using the Frobenius norm has the added benefit of making the objective convex in each argument separately, suggesting the possibility of efficient optimization. This objective function is underdetermined, however, so we consider two constrained regimes below.

Full Σ_i approximation. The first method we call JD – Full. Without loss of generality, U and V can be constrained to be orthogonal, so long as Σ_i remains an unconstrained full matrix. JD – Full adopts this restriction to make the optimization better posed, but note it does not restrict the expressiveness of the objective equation 1. This setting yields the following optimization problem:

$$JD-Full_{r}(\{B_{i}A_{i}\}_{i=1}^{n}) = \underset{\substack{\{\Sigma_{i}\}_{i=1}^{n}\\U^{\top}U=V^{\top}V=I_{r}}}{\operatorname{argmin}} \sum_{i=1}^{n} \|B_{i}A_{i} - U\Sigma_{i}V^{\top}\|_{\operatorname{Fro}}^{2} \quad (\mathsf{JD}-\mathsf{Full})$$
(2)

An efficient alternating algorithm to solve this objective function can be found in Appendix A.

Diagonal Σ_i approximation. As an alternative, we can leave U, V unconstrained (other than to have r columns) and instead constrain the matrices Σ_i to be diagonal (but not necessarily positive). This formulation yields the following optimization problem:

$$JD-Diag_{r}(\{B_{i}A_{i}\}_{i=1}^{n}) = \underset{\{\Sigma_{i}\}_{i=1}^{n}, U, V}{\operatorname{argmin}} \sum_{i=1}^{n} \|B_{i}A_{i} - U\operatorname{diag}(\Sigma_{i})V^{\top}\|_{\operatorname{Fro}}^{2} \quad (\mathsf{JD} - \mathsf{Diag})$$
(3)

An efficient alternating least squares algorithm to optimize this objective can be found in Appendix A. This diagonal version has some per-LoRA parameter savings when compared to JD - Full, since the diagonal Σ_i only needs r parameters instead of r^2 .

3.2 CLUSTERING

As the number of LoRAs n grows and becomes more diverse, the rank r needed for Joint Diagonalization to achieve good performance will tend to increase. This increases the size and number of parameters of each Σ_i that needs to be stored, especially for JD-Full which will require $O(nr^2)$ storage for these matrices. If the necessary r is growing proportionally to n, then this storage will eventually become the bottleneck.

To resolve this limitation with very large n, we propose to group the n LoRAs into $|C_i|$ clusters C_i . Each cluster is given its own rank r JD compression, and the clusters are chosen such that the overall reconstruction error is minimized. Specifically, the overall objective is

$$\min_{\{\{C_j\}, U_j, V_j\}, \{\Sigma_i\}} \sum_j \sum_{i \in C_j} ||B_i A_i - U_j \Sigma_i V_j||_F^2$$

and we optimize this by alternating between cluster assignments and the JD of each cluster. The algorithm details are in Appendix A.3. Typically, the goal with large n is to have $|C_i|$ grow with n as r becomes fixed. Comparing k rank-r JD-Full clusters to a rank-kr JD-Full single cluster compression, observe that the clustered approach requires $O(dkr + nr^2)$ parameters, while the single-cluster approach requires $O(dkr + nk^2r^2)$ parameters due to the increased sizes of the $\Sigma_i s$. While these two approaches have the same rank, note that they may have different reconstruction abilities. Empirically, we find that multiple clusters significantly aid performance for $n \ge 100$.

216 4 THEORETICAL ANALYSIS217

In this section, we seek to better understand the role of the joint diagonalization method presented in §3.1 and how this understanding further motivates the clustering approach. We will focus on the full- Σ_i case with orthogonal U, V matrices. Note that, for the same r, the r-JD-Diag has at least as large reconstruction error as r-JD-Full since it imposes an additional constraint on the Σ_i .

Firstly, note that perfect reconstruction can be achieved if and only if r is large enough, since there exist U, V such that all the B_i, A_i are in the spans of U, V resp. if and only if $r \ge \tilde{r}$:

Proposition 1. Suppose that for all *i*, $rank(B_iA_i) = r_i$, and let

 $\tilde{r} = \max\left\{\operatorname{rank}([A_1,\ldots,A_n]),\operatorname{rank}([B_1^{\top}\ldots,B_n^{\top}])\right\}.$

Note $\max_i r_i \leq \tilde{r} \leq \sum_{i=1}^n r_i$. Then JD – Full (equation 2) with $r = \tilde{r}$ achieves lossless compression (perfect reconstruction), and using $r < \tilde{r}$ will give nonzero reconstruction error.

Due to training noise, \tilde{r} will equal $\sum_{i=1}^{n} r_i$ almost always. This implies that in most realistic settings, the joint diagonalization approach is a lossy reconstruction.

This reconstruction loss can be significant, as the following theorem shows (proved in Appendix B):

Theorem 1. Consider n LoRAs $(\{A_i, B_i\}_{i=1}^n)$ with $r, n \leq d^2$, and form the matrix 235

 $L = \left[\operatorname{vec}(B_1 A_1) \cdots \operatorname{vec}(B_n A_n) \right].$

Let σ_j be the singular values of L, sorted from largest to smallest, and let $\bar{\sigma}_j$ be the singular values of $\sum_{i=1}^{n} B_i A_i$. Then, using JD – Full (equation 2),

$$\sum_{j=1}^{r} \bar{\sigma}_{j}^{2} \leq \sum_{i=1}^{n} \|\Sigma_{i}\|_{\text{Fro}}^{2} = \sum_{i=1}^{n} \|U\Sigma_{i}V^{\top}\|_{\text{Fro}}^{2} \leq \sum_{j=1}^{\min(r^{2},n)} \sigma_{j}^{2}$$

implying the sum of squared Frobenius norms of the reconstructed LoRAs satisfies

$$= \frac{\sum_{i=1}^{n} \|U\Sigma_{i}V^{\top}\|_{\text{Fro}}^{2}}{\sum_{i=1}^{n} \|B_{i}A_{i}\|_{\text{Fro}}^{2}} \leq \frac{\sum_{j=1}^{\min(r^{2},n)}\sigma_{j}^{2}}{\sum_{j=1}^{n}\sigma_{j}^{2}} \leq 1, \text{ and } \frac{\sum_{i=1}^{n} \|U\Sigma_{i}V^{\top} - B_{i}A_{i}\|_{\text{Fro}}^{2}}{\sum_{i=1}^{n} \|B_{i}A_{i}\|_{\text{Fro}}^{2}} \geq 1 - \frac{\sum_{j=1}^{\min(r^{2},n)}\sigma_{j}^{2}}{\sum_{j=1}^{n}\sigma_{j}^{2}} \leq 1, \text{ and } \frac{\sum_{i=1}^{n} \|U\Sigma_{i}V^{\top} - B_{i}A_{i}\|_{\text{Fro}}^{2}}{\sum_{i=1}^{n} \|B_{i}A_{i}\|_{\text{Fro}}^{2}} \geq 1 - \frac{\sum_{j=1}^{\min(r^{2},n)}\sigma_{j}^{2}}{\sum_{j=1}^{n}\sigma_{j}^{2}} \leq 1, \text{ and } \frac{\sum_{i=1}^{n} \|U\Sigma_{i}V^{\top} - B_{i}A_{i}\|_{\text{Fro}}^{2}}{\sum_{i=1}^{n} \|B_{i}A_{i}\|_{\text{Fro}}^{2}} \geq 1 - \frac{\sum_{j=1}^{\min(r^{2},n)}\sigma_{j}^{2}}{\sum_{j=1}^{n} \sigma_{j}^{2}} \leq 1, \text{ and } \frac{\sum_{i=1}^{n} \|U\Sigma_{i}V^{\top} - B_{i}A_{i}\|_{\text{Fro}}^{2}}{\sum_{j=1}^{n} \|B_{i}A_{j}\|_{\text{Fro}}^{2}} \geq 1 - \frac{\sum_{j=1}^{\min(r^{2},n)}\sigma_{j}^{2}}{\sum_{j=1}^{n} \sigma_{j}^{2}} \leq 1, \text{ and } \frac{\sum_{i=1}^{n} \|B_{i}A_{i}\|_{\text{Fro}}^{2}}{\sum_{j=1}^{n} \|B_{i}A_{j}\|_{\text{Fro}}^{2}} \geq 1 - \frac{\sum_{j=1}^{n} \|B_{i}A_{j}\|_{\text{Fro}}^{2}}{\sum_{j=1}^{n} \|B_{i}A_{j}\|_{\text{Fro}}^{2}} \leq 1 - \frac{\sum_{j=1}^{n} \|B_{i}A_{j}\|_{\text{Fro}}^{2}}{\sum_{j=1}^{n} \|B_{i}A_{j}\|_{\text{Fro}}^{2}}} \leq 1 - \frac{\sum_{j=1}^{n} \|B_{i}A_{j}\|_{\text{Fro}}^{2}}}{\sum_{j=1}^{n} \|B_{i}A_{j}\|_{\text{Fro}}^{2}}} \leq 1 - \frac{\sum_{j=1}^{n} \|B_{i}A_{j}\|_{\text{Fro}}^{2}}}{\sum_{j=1}^{n} \|B_{i}A_{j}\|_{\text{Fro}}^{2}}} \leq 1 - \frac{\sum_{j=1}^{n} \|B_{i}A_{j}\|_{\text{Fro}}^{2}}}$$

In other words, if the singular values of L are not concentrated in the top r^2 entries, significant reconstruction error is unavoidable.

Remark 1 (Lower bound and merging). The lower bound $\sum_{j=1}^{r} \bar{\sigma}_{j}^{2}$ could be achieved by setting all the Σ_{i} equal, i.e., using a fully merged model instead of only merging the subspaces U, V and allowing Σ_{i} to vary with i.

Remark 2 (Upper bound and grouping). The upper bound is smallest when the LoRAs are relatively clustered, i.e., when groups of vectors $vec(B_iA_i)$ are similar. This situation raises the magnitude of the largest singular values of L, raising the upper bound in the proposition. As the LoRAs are $d \times d$ matrices that can be thought of as points in d^2 dimensional space, for typical values of d well into the hundreds, it is likely that unrelated LoRAs will be unclustered, i.e., they will have relatively low inner products with each other.

For the case of orthogonal LoRAs, the singular values of L are the norms of the LoRAs, and we immediately have the following corollary:¹

Corollary 1. Suppose (e.g., due to normalization) that the inputs to the joint diagonalization algorithm all have unit Frobenius norm, i.e., $||B_iA_i||_{\text{Fro}} = 1$. Moreover, assume that the LoRAs are all orthogonal in the sense $\text{tr}((B_iA_i)(B_jA_j)^{\top}) = 0$ for $i \neq j$. Then, using the JD – Full method equation 2, we have $1 \leq \sum_{i=1}^{n} ||\Sigma_i||_{\text{Fro}}^2 \leq \min(r^2, n)$, implying that the sum of squared Frobenius norms of the reconstructed LoRAs satisfies

267 268

259

225

226 227

228

229 230

231 232

233

236

$$1 - \frac{1}{n} \ge \frac{\sum_{i=1}^{n} \|U\Sigma_i V^{\top} - B_i A_i\|_{\text{Fro}}^2}{\sum_{i=1}^{n} \|B_i A_i\|_{\text{Fro}}^2} \le 1 - \min\left(\frac{r^2}{n}, 1\right).$$

¹A result for isotropic Gaussian LoRAs could be obtained via the quantiles of the Marchenko-Pastur Law.

This implies that for the common setting where $r^2 \ll n$, the reconstructed LoRAs will be significantly smaller than the original LoRAs and necessarily have significant reconstruction error.

The results in this section illustrate the tradeoffs of using joint diagonalization. If the LoRAs are similar or well-clustered, reconstruction error will be low. On the other hand, if the LoRAs are random and orthogonal, reconstruction error will be high.

Since the loss space of transformers is highly complex, increasing weight reconstruction error does 276 not necessarily imply degrading LLM performance. Interestingly, in Figure 3 below, we see that 277 while large reconstruction error rapidly decreases performance, moderate (but still relatively large, at 278 around 0.4) reconstruction error does not damage performance and may even slightly outperform the 279 zero-error setting. This observation motivates our focus on minimizing weight reconstruction error, 280 while also suggesting that our approach is capable of achieving something deeper than compression. 281 Specifically, the tendency of joint diagonalization is to find subspaces that are shared among many 282 LoRAs when r is large, and to *merge* subspaces when r is small. When r is particularly small, this 283 tendency towards averaging all or some of the LoRA signals directly connects to the concept of 284 merging LoRAs, whose empirical success (Shah et al., 2023; Huang et al., 2024) could explain the 285 success of our procedure despite the nonlinearity of transformers.

Experiments in Appendix G.9 explore this idea further, comparing reconstruction of real-world Lo-RAs to reconstruction of randomly sampled LoRAs. The reconstruction error is generally large, but significantly lower than the reconstruction error for random noise, indicating that there is a major shared component between the LoRAs that is being successfully retained.

That said, as the number of LoRAs grows, the shared component may not be significant enough to maintain sufficiently low reconstruction error with low rank r. This motivates the introduction of *clustering* in §3.2, since clustering seeks to find groups of LoRAs that are similar and better compressible by joint diagonalization. In particular, if the number of clusters $|C_j|$ grows with n, the reconstruction error may no longer degrade with n even when r is fixed.

In the extreme case where $|C_j| = n$, each LoRA is compressed independently. By the Eckart-Young Theorem, JD applied to a single LoRA reduces to an SVD, replacing each rank- r_i LoRA adapter B_iA_i with a reduced rank-r approximation, where typically $r < \frac{1}{n} \sum_{i=1}^{n} r_i$:

 $SVD_r(B_iA_i) = U_i\Sigma_iV_i^{\top}, \quad \forall i = 1,\dots,n.$ (4)

300 As $\Sigma_i V_i^{\top}$ can be saved as a single matrix, this approach has $rn(d_A + d_B)$ parameters. We refer to 301 this $|C_j| = n$ method as r - SVD and find that it underperforms our other methods, while outper-302 forming the baseline uncompressed LoRAs significantly. This result parallels Jiang et al. (2023b)'s 303 observation that lowering LoRA ranks is beneficial for multi-task learning and model merging.

5 TRAINING LORAS & EVALUATING TASK PERFORMANCE

5.1 TRAINING

We trained LoRA adapters on 500 natural instruction tasks (Wang et al., 2022) using Mistral-7B-Instruct-v0.2 (Jiang et al., 2023a) as the base model. All LoRA adapters were configured with a rank of 16, i.e., $\forall i, r_i = 16$.

We selected 10 diverse tasks (Table 2 in Appendix C) manually for consistent evaluation across experiments and randomly sampled an additional 490 tasks, resulting in a total of 500 tasks. These tasks were exclusively in English (both input and output), ensuring higher quality and thorough review (Wang et al., 2022). The tasks represent a realistic and varied set, not inherently clustered. Each task dataset was divided into training, validation, and test sets.

Hyperparameters, such as early stopping, were tuned using the validation sets. Evaluation on the test sets demonstrated that LoRA consistently outperformed the base model in terms of both Rouge scores and loss metrics, as shown in Table 1. Details are provided in Appendix C.

320

299

304

306 307

308

321 5.2 EVALUATION 322

We evaluated multiple metrics for the natural instruction tasks, including cross-entropy loss, Rouge-1, Rouge-L (Lin, 2004), exact match, and *agreement* between uncompressed and compressed LoRA.

| Metric | Base Model | LoRA |
|-------------|-------------------|------------------|
| Loss | 4.99 ± 3.11 | 0.43 ± 0.57 |
| Exact Match | 2.28 ± 7.89 | 66.66 ± 34.3 |
| Rouge-1 | 20.38 ± 18.90 | 76.74 ± 24.8 |
| Rouge-L | 19.66 ± 18.16 | 76.22 ± 25.2 |

Table 1: Comparison of metrics before and after LoRA training across 500 tasks.

331 332 333

351

352 353 354

355

356 357

358 359

360

330

324

Here, *agreement* measures the exact match in task-generations between the uncompressed LoRA model and the compressed LoRA model, rather than comparing to ground truth data. While detailed results and discussions for all metrics are provided in Appendix G, our primary focus in the main text is on Rouge-L. We find that all metrics correlate, but Rouge-L correlates most strongly with downstream utility. This finding aligns with prior work (Wang et al., 2022), which demonstrates that Rouge-L correlates well with classification accuracy.

While cross-entropy is used for optimization during training, identical generation outputs across models can yield different cross-entropy losses. Exact match is too rigid and does not account for the variability in task responses. Similarly, agreement does not capture the inexactness associated with most of our tasks, nor does it account for the performance gains or losses of the compressed LoRAs. Arguably, practitioners are primarily concerned with task performance in the settings for which the LoRA was designed, rather than exact generational agreement between models.

Joint diagonalization optimizes reconstruction error measured by the Frobenius norm, and our theoretical analysis in §4 bounds this reconstruction error. Empirically, reconstruction error and downstream Rouge-L performance correlate.

Instead of listing the absolute performance of different methods, we compute the performance dif-ference between the base model and the LoRA model for each task. We present the ratio

Performance relative to $LoRA := \frac{method-performance}{LoRA-performance}$

for the specific method in question, highlighting relative improvement with respect to the uncompressed LoRA and the base model.

- 6 EXPERIMENTS
- 6.1 TASK PERFORMANCE

For each method, we vary the number of n LoRAs that are compressed and the compression rank r. We run each experiment three times with different random seeds and report the mean and standard deviation. See Table 4 for results where we evaluate on the same ten manually selected tasks (Table 2) across settings. Every compressed collection of LoRAs contains these 10 tasks (i.e., indistribution tasks), and each collection contains the smaller collections as subsets.

We normalize each LoRA adapter to have a Frobenius norm of one prior to running joint diagonal ization. This normalization enhances performance and reduces the variance in reconstruction error.
 We restore the original norms of the LoRA adapters before reconstruction and testing.

Figure 2a illustrates the Total Parameter Saved Ratio versus the Number of Unique LoRAs served.
We only include methods that maintain over 99% of the original LoRA's performance (as measured by RougeL). Notably, our JD methods uniquely approach the compression efficacy of a single LoRA, and with clustering, this aggressive reduction in size also maintains performance in larger LoRA collections.

Figure 2b illustrates the Rouge-L scored of the compressed LoRAs divided by the Rouge-L score of
 the uncompressed LoRAs. It is interesting to note that JD variants often increase generalization and
 outperform the original LoRA. In Appendix G, we include multiple tables of results for additional
 metrics, relative as well as absolute.



(a) Total parameter saved ratio with the number of unique LoRAs served.

(b) Performance relative to LoRA with total parameter saved ratio

Figure 2: Performance after compression. In (a), we only include methods that maintain over 99% of the original LoRA's performance (as measured by RougeL). In (b), we compare the performance of compressed LoRAs relative to uncompressed ones, with higher values on both axes reflecting better performance. The Total Parameter Saved Ratio depicts the number of parameters saved for a system with a large number n of different LoRAs. It is computed as: $r_{total} := 1 - \frac{\text{num. parameters after compression}}{\text{num. parameters before compression}}$.

For efficiency, we limited the JD methods to ten iterations instead of pursuing full convergence. While the alternating algorithm quickly reaches an approximation of the minimizer, squeezing out the last few digits of precision takes many more iterations with limited to no performance gain. Appendix G.10 also evaluates an alternative eigenvalue iteration algorithm that more rapidly converges once U, V are close to a minimizer, with minimal performance differences.

6.2 PERFORMANCE AND RECONSTRUCTION ERROR

Figure 3 relates reconstruction error and performance. The y-axis measures the mean performance improvement of Rouge-L relative to uncompressed LoRA, and the x-axis quantifies the mean reconstruction error between the compressed reconstruction of the product BA and the original uncompressed product BA. Although the relationship between performance and reconstruction error is nonlinear, it demonstrates a generally decreasing, somewhat exponential trend. Notably, the minimal reconstruction error does not correlate with optimal performance, indicating that a degree of lossy reconstruction may be advantageous for enhancing generalization.

To select hyperparameters (compression rank and number of clusters) for the clustering experiments, we first assessed reconstruction error on a single LoRA module over a range of settings (see Appendix F). These preliminary experiments enabled efficient selection of cluster counts and rank values for compressing all LoRA modules.



387

388 389

390

391

392

393

394

396

397

398

399

400

401 402

403



Figure 4: Throughput ratio when serving varying numbers of LoRAs with vLLM.

512 1024

256

6.3 THROUGHPUT OF SERVING COMPRESSED LORAS

Results in the previous sections demonstrate how to select an appropriate joint compression setting guided by the reconstruction error, such that the performance of the original LoRAs is preserved.

Naturally, the rank and/or the number of clusters for the compression needs to increase as we compress larger LoRA collections to match LoRA performance.

In Figure 4 we study how throughput with various compression settings compares to the vLLM 435 multi-LoRA throughput with the matched GPU memory footprint. Specifically, for each number 436 of unique LoRAs served and each compression setting, we compute the corresponding number of 437 LoRAs to be placed on the GPU during serving and report the ratio of the two throughputs. For 438 example, when serving 64 unique LoRAs and using rank 64 JD-Full compression, we report the 439 ratio of throughputs of rank 64 JD-Full and vLLM multi-LoRA with 6 LoRAs allowed on the GPU 440 at a time (see Appendix E for details). As the number of unique LoRAs increases, vLLM multi-441 LoRA throughput degrades as it needs to schedule the requests and load and offload the adapters. 442 We note that vLLM multi-LoRA already employs many advanced system optimization techniques, such as efficient scheduling and non-blocking CPU-GPU communication when swapping LoRAs 443 (Sheng et al., 2023; Kwon et al., 2023), but system optimization alone is not sufficient to fully 444 mitigate throughput degradation when serving many LoRAs. 445

446 In Figure 4 we see that across all LoRA collection sizes our compression techniques improve the 447 throughput of vLLM multi-LoRA. Additionally, we highlight regions for each compression setting where compression is sufficiently moderate to achieve 99%+ of LoRA performance, according to 448 449 the results in Section 6.2. We also note that compression with a larger rank or too many clusters does not improve baseline throughput when serving a smaller number of LoRAs and should not be used 450 in such cases. For example, rank 16 JD-Full improves baseline throughput with 4 and 8 LoRAs, but 451 will underperform with more LoRAs, while 7 cluster rank 64 JD-Full does not improve throughput 452 with 64 or fewer LoRAs, but when serving 1000+ LoRAs it improves the throughput significantly 453 while maintaining the performance. To conclude, an appropriate joint compression setting improves 454 vLLM multi-LoRA throughput and preserves performance for LoRA collections of any size between 455 4 and 1024, as we showed in Figure 1. Specific compression settings for each LoRA collection size 456 are listed in Appendix E. 457

Finally, we note that vLLM extensively uses custom CUDA kernels. To accommodate our compression techniques, we minimally adjusted the vLLM code to generate additional kernels needed by the compressed LoRAs while we utilized Punica (Chen et al., 2023) kernel to further accelerate matrix multiplication. A pseudo code is given in E.4 to show how we utilize the batch multiplication kernel. There likely is room for improvement to optimize the newly added kernels.

462 463

Additional details In this experiment we considered a varying number of rank-16 LoRAs, using 464 a dataset of Shakespeare sonnets as inputs² arriving asynchronously. We measured throughput, i.e., 465 the number of requests served per second when generating ten tokens per request. The base model 466 was Mistral 7B Instruct as in the other experiments; we simulated random LoRAs and assigned 467 inputs to LoRAs at random. Experiments were conducted on H100 80GB GPU capped at 40% 468 memory consumption. This was done to reflect cost concerns in practical situations where a service provider might want to serve many LoRAs from cheaper hardware with lower memory than higher-469 end GPUs. This setting also takes into account the scenario where the LLM is large compared to the 470 size of GPU and yet a provider may want to serve many LoRAs efficiently using the same device. 471

- 472
- 473 474

475

6.4 **RECOMMENDATIONS**

JD-Full is generally preferred over JD-Diag, although for smaller numbers of LoRAs (less than 100), the performance difference is negligible. While JD-Full alone is effective up to 100 LoRAs, incorporating clustering at scales of 500 LoRAs significantly enhances performance.

We recommend the following procedure for hyperparameter selection. For 100 or fewer LoRAs, JD-Full can be utilized independently without substantial degradation, using a rank approximately equal to (number of LoRAs/2) + 7. Beyond 100 LoRAs, clustering becomes increasingly critical. A robust method for any number of LoRAs up to 500 involves employing JD-Full with clustering. Specifically, select a LoRA module from the middle of the network, apply a compression rank of 16,

²https://www.kaggle.com/datasets/shivamshinde123/william-shakespeares-sonnet/ data

and experiment with an exponentially increasing number of clusters. Compute the reconstruction
error for each setting on this module across all LoRAs—a computationally efficient process. Choose
the minimal number of clusters that achieves a reconstruction loss below 0.5, and then use these
settings to compress all LoRA modules. An example of this procedure applied to 500 LoRAs is
illustrated in Figure 5 in the Appendix.

We note that tuning hyperparameters as discussed above using reconstruction loss as a validation metric is especially convenient since it can be done efficiently on CPU without having to perform expensive LLM evaluation. As our experiments demonstrate, compression settings that achieve below 0.5 reconstruction loss reliably translate into preserving 99% or more of the LoRA performance, sometimes even outperforming the original LoRAs.

For inference, this procedure is executed as a preprocessing step before deploying our inference server. As new LoRAs are submitted, they are initially served uncompressed. A background cron job re-runs the compression algorithm on the CPU every six hours, and upon completion, updates the served LoRA parameters with the compressed versions.

501 7 DISCUSSION

502

This study introduces approaches to LoRA compression, addressing significant challenges facing
 foundation models and large language models. Our contributions include theoretical formulations,
 empirical validation, and practical implementations that enhance the understanding and application
 of LLMs in scalable environments.

507 The implications of our findings are manifold. Our theoretical guarantees for reconstruction error 508 not only increase confidence in the use of compressed models but also lay a groundwork for future 509 explorations in this area. Demonstrating that our compression techniques can preserve up to 100% 510 of the original LoRAs' performance highlights the effectiveness of our methods. Furthermore, in-511 tegrating LoRA compression into state-of-the-art LLM serving systems demonstrates the potential 512 for resource optimization, with throughput for thousands of LoRAs nearing that of a single LoRA.

The promising results of our study suggest several future research directions. First, further compression may be possible via quantization. Our joint-diagonalization compression and quantization are independent axes of approaching the problem and exploring a combined solution can be fruitful. Second, when scaling to hundreds of thousands of LoRAs, joint compression, while effective, will not be sufficient to fit all LoRAs onto the GPU, thus requiring a procedure to schedule the requests. Our clustering variant offers opportunities to develop an efficient scheduling mechanism that takes into account the cluster assignments of LoRAs corresponding to the incoming requests.

In conclusion, our research significantly advances the deployment of LLMs by providing robust,
 scalable, and efficient compression solutions. The ability of compressed LoRAs to maintain high
 performance while facilitating substantial resource savings opens new avenues for the broader application and adoption of LLMs across various industries. We encourage the community to build upon
 our findings and the shared LoRAs to further explore and enhance the utility of these technologies.

References

525 526

527

528

529 530

- Lequn Chen, Zihao Ye, Yongji Wu, Danyang Zhuo, Luis Ceze, and Arvind Krishnamurthy. Punica: Multi-tenant lora serving, 2023. URL https://arxiv.org/abs/2310.18547.
- Yu Cheng, Duo Wang, Pan Zhou, and Tao Zhang. A survey of model compression and acceleration for deep neural networks. *arXiv preprint arXiv:1710.09282*, 2017.
- Leshem Choshen, Elad Venezian, Noam Slonim, and Yoav Katz. Fusing finetuned models for better
 pretraining. *ArXiv*, abs/2204.03044, 2022.
- Leshem Choshen, Elad Venezian, Shachar Don-Yehiya, Noam Slonim, and Yoav Katz. Where to start? analyzing the potential value of intermediate models. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 1446–1470, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.90. URL https://aclanthology.org/2023.emnlp-main.90.

- 540 Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning 541 of quantized llms. Advances in Neural Information Processing Systems, 36, 2024. 542 Amir Gholami, Sehoon Kim, Zhen Dong, Zhewei Yao, Michael W Mahoney, and Kurt Keutzer. A 543 survey of quantization methods for efficient neural network inference. In Low-Power Computer 544 Vision, pp. 291–326. Chapman and Hall/CRC, 2022. 546 Tom Gunter, Zirui Wang, Chong Wang, Ruoming Pang, Andy Narayanan, Aonan Zhang, Bowen 547 Zhang, Chen Chen, Chung-Cheng Chiu, David Qiu, et al. Apple intelligence foundation language 548 models. arXiv preprint arXiv:2407.21075, 2024. 549 Moshik Hershcovitch, Leshem Choshen, Andrew Wood, Ilias Enmouri, Peter Chin, Swaminathan 550 Sundararaman, and Danny Harnik. Lossless and near-lossless compression for foundation models. 551 arXiv preprint arXiv:2404.15198, 2024. 552 553 Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, An-554 drea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for nlp. 555 In International conference on machine learning, pp. 2790–2799. PMLR, 2019. 556 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. arXiv preprint 558 arXiv:2106.09685, 2021. 559 Chengsong Huang, Qian Liu, Bill Yuchen Lin, Tianyu Pang, Chao Du, and Min Lin. Lorahub: 561 Efficient cross-task generalization via dynamic lora composition, 2024. 562 Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, 563 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 564 Mistral 7b. arXiv preprint arXiv:2310.06825, 2023a. 565 566 Weisen Jiang, Baijiong Lin, Han Shi, Yu Zhang, and James T Kwok. Byom: Building your own 567 multi-task model for free. arXiv preprint arXiv:2310.01886, 2023b. 568 Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph 569 Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model 570 serving with pagedattention. In Proceedings of the 29th Symposium on Operating Systems Prin-571 *ciples*, pp. 611–626, 2023. 572 573 Chunyuan Li, Heerad Farkhoor, Rosanne Liu, and Jason Yosinski. Measuring the intrinsic dimension 574 of objective landscapes, 2018. 575 Vladislav Lialin, Vijeta Deshpande, and Anna Rumshisky. Scaling down to scale up: A guide to 576 parameter-efficient fine-tuning. arXiv preprint arXiv:2303.15647, 2023. 577 578 Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In Text Summarization 579 Branches Out, pp. 74-81, Barcelona, Spain, July 2004. Association for Computational Linguis-580 tics. URL https://aclanthology.org/W04-1013. 581 Haokun Liu, Derek Tam, Mohammed Muqeeth, Jay Mohta, Tenghao Huang, Mohit Bansal, and 582 Colin A Raffel. Few-shot parameter-efficient fine-tuning is better and cheaper than in-context 583 learning. Advances in Neural Information Processing Systems, 35:1950–1965, 2022. 584 585 Shih-Yang Liu, Chien-Yi Wang, Hongxu Yin, Pavlo Molchanov, Yu-Chiang Frank Wang, Kwang-586 Ting Cheng, and Min-Hung Chen. Dora: Weight-decomposed low-rank adaptation, 2024. Michael Matena and Colin Raffel. Merging models with fisher-weighted averaging. arXiv preprint 588 arXiv:2111.09832, 2021. 589 Fanxu Meng, Zhaohui Wang, and Muhan Zhang. Pissa: Principal singular values and singular 591 vectors adaptation of large language models. arXiv preprint arXiv:2404.02948, 2024. 592
- 593 Mohammed Muqeeth, Haokun Liu, Yufan Liu, and Colin Raffel. Learning to route among specialized experts for zero-shot generalization. *arXiv preprint arXiv:2402.05859*, 2024.

- 594 Kimia Nadjahi, Kristjan Greenewald, Rickard Brüel Gabrielsson, and Justin Solomon. Slicing mu-595 tual information generalization bounds for neural networks. In ICML 2023 Workshop Neural 596 Compression: From Information Theory to Applications, 2023. URL https://openreview. 597 net/forum?id=cbLcwK3SZi. 598 Openai fine-tuning api. https://platform.openai.com/docs/guides/ OpenAL. fine-tuning, 2024. 600 601 Predibase. Multi-lora inference server that scales to 1000s of fine-tuned llms. https:// 602 loraexchange.ai, 2024. 603 604 Colin Raffel. Building machine learning models like open source software. Communications of the 605 ACM, 66(2):38-40, 2023. 606 Viraj Shah, Nataniel Ruiz, Forrester Cole, Erika Lu, Svetlana Lazebnik, Yuanzhen Li, and Varun 607 Jampani. Ziplora: Any subject in any style by effectively merging loras. arXiv preprint 608 arXiv:2311.13600, 2023. 609 610 Pratyusha Sharma, Jordan T. Ash, and Dipendra Misra. The truth is in there: Improving reasoning 611 in language models with layer-selective rank reduction. In The Twelfth International Confer-612 ence on Learning Representations, 2024. URL https://openreview.net/forum?id= 613 ozX92bu8VA. 614 Ying Sheng, Shiyi Cao, Dacheng Li, Coleman Hooper, Nicholas Lee, Shuo Yang, Christopher Chou, 615 Banghua Zhu, Lianmin Zheng, Kurt Keutzer, Joseph E. Gonzalez, and Ion Stoica. S-lora: Serving 616 thousands of concurrent lora adapters, 2023. 617 618 Together fine-tuning. https://www.together.ai/products# TogetherAI. 619 fine-tuning, 2024. 620 Sheng Wang, Boyang Xue, Jiacheng Ye, Jiyue Jiang, Liheng Chen, Lingpeng Kong, and Chuan Wu. 621 Prolora: Partial rotation empowers more parameter-efficient lora. ArXiv, abs/2402.16902, 2024. 622 URL https://api.semanticscholar.org/CorpusID:268032580. 623 624 Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei, An-625 jana Arunkumar, Arjun Ashok, Arut Selvan Dhanasekaran, Atharva Naik, David Stap, et al. 626 Super-naturalinstructions: Generalization via declarative instructions on 1600+ nlp tasks. arXiv 627 preprint arXiv:2204.07705, 2022. 628 Yeming Wen and Swarat Chaudhuri. Batched low-rank adaptation of foundation models, 2024. 629 630
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi,
 Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick
 von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger,
 Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Huggingface's transformers: Stateof-the-art natural language processing, 2020.
- Mitchell Wortsman, Gabriel Ilharco, Samir Yitzhak Gadre, Rebecca Roelofs, Raphael Gontijo Lopes, Ari S. Morcos, Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon Kornblith, and
 Ludwig Schmidt. Model soups: averaging weights of multiple fine-tuned models improves ac curacy without increasing inference time. In *International Conference on Machine Learning*,
 2022.
- Prateek Yadav, Leshem Choshen, Colin Raffel, and Mohit Bansal. Compet: Compression for communicating parameter efficient updates via sparsification and quantization. *arXiv preprint arXiv:2311.13171*, 2023a.
- Prateek Yadav, Derek Tam, Leshem Choshen, Colin Raffel, and Mohit Bansal. TIES-merging: Resolving interference when merging models. In *Thirty-seventh Conference on Neural In*formation Processing Systems, 2023b. URL https://openreview.net/forum?id= xtaX3WyCjl.

- Prateek Yadav, Colin Raffel, Mohammed Muqeeth, Lucas Caccia, Haokun Liu, Tianlong Chen, Mohit Bansal, Leshem Choshen, and Alessandro Sordoni. A survey on model moerging: Recycling and routing among specialized experts for collaborative learning. *arXiv preprint arXiv:2408.07057*, 2024.
- 46524653 Yuchen Zeng and Kangwook Lee. The expressive power of low-rank adaptation, 2024.
- Qingru Zhang, Minshuo Chen, Alexander Bukharin, Pengcheng He, Yu Cheng, Weizhu Chen, and Tuo Zhao. Adaptive budget allocation for parameter-efficient fine-tuning. In *The Eleventh International Conference on Learning Representations*, 2023.
- Jiacheng Zhu, Kristjan Greenewald, Kimia Nadjahi, Haitz Sáez de Ocáriz Borde, Rickard Brüel
 Gabrielsson, Leshem Choshen, Marzyeh Ghassemi, Mikhail Yurochkin, and Justin Solomon.
 Asymmetry in low-rank adapters of foundation models. *arXiv preprint arXiv:2402.16842*, 2024.

A JOINT DIAGONALIZATION ALGORITHMS

A.1 ALTERNATING METHODS

Our goal is to derive algorithms that optimize equation 1. Common to both methods, we expand the objective functional:

$$\sum_{i} \|B_{i}A_{i} - U\Sigma_{i}V^{\top}\|_{\text{Fro}}^{2} = \sum_{i} \operatorname{tr}((B_{i}A_{i} - U\Sigma_{i}V^{\top})(B_{i}A_{i} - U\Sigma_{i}V^{\top})^{\top}) \text{ by definition}$$
$$= \sum_{i} \left[\operatorname{tr}(B_{i}A_{i}A_{i}^{\top}B_{i}^{\top}) - 2\operatorname{tr}(B_{i}A_{i}V\Sigma_{i}^{\top}U^{\top}) + \operatorname{tr}(U\Sigma_{i}V^{\top}V\Sigma_{i}^{\top}U^{\top})\right]$$
$$= \operatorname{const.} - 2\sum_{i} \operatorname{tr}(B_{i}A_{i}V\Sigma_{i}^{\top}U^{\top}) + \sum_{i} \|U\Sigma_{i}V^{\top}\|_{\text{Fro}}^{2}.$$
(5)

Using this expansion, we now consider the two settings discussed in $\S3.1$.

Case 1: Non-diagonal Σ_i , orthogonal U, V. Setting the derivative of equation 5 with respect to Σ_i to zero, we find

$$\Sigma_i = \Sigma_i^*(U, V) = U^\top B_i A_i V.$$
(6)

We simplify our objective function after plugging in this expression:

$$\sum_{i} \|B_{i}A_{i} - U\Sigma_{i}V^{\top}\|_{\text{Fro}}^{2} + \text{const.} = \sum_{i} \left[\|\Sigma_{i}\|_{\text{Fro}}^{2} - 2\text{tr}(B_{i}A_{i}V\Sigma_{i}^{\top}U^{\top})\right] \text{ from equation 5}$$
$$= \sum_{i} \left[\text{tr}(U^{\top}B_{i}A_{i}VV^{\top}A_{i}^{\top}B_{i}^{\top}U) - 2\text{tr}(B_{i}A_{i}VV^{\top}A_{i}^{\top}B_{i}^{\top}UU^{\top})\right] \text{ from equation 6}$$
$$= -\sum_{i} \text{tr}(B_{i}A_{i}VV^{\top}A_{i}^{\top}B_{i}^{\top}UU^{\top}).$$

Substituting equation 6, we find

$$U_{opt}, V_{opt} = \arg\max_{\substack{U^{\top}U=I\\VV^{\top}=I}} \sum_{i=1}^{n} \|U^{\top}B_{i}A_{i}V\|_{\text{Fro}}^{2} = \arg\max_{\substack{U^{\top}U=I\\VV^{\top}=I}} \sum_{i=1}^{n} \|\Sigma_{i}^{*}(U,V)\|_{\text{Fro}}^{2}.$$
 (7)

Note that

$$\sum_{i=1}^{n} \|U^{\top} B_{i} A_{i} V\|_{\text{Fro}}^{2} = \operatorname{tr} \left(\left(\sum_{i=1}^{n} B_{i} A_{i} V V^{\top} A_{i}^{\top} B_{i}^{\top} \right) U U^{\top} \right)$$
$$= \operatorname{tr} \left(\left(\sum_{i=1}^{n} B_{i}^{\top} A_{i}^{\top} U U^{\top} A_{i} B_{i} \right) V V^{\top} \right)$$

by the identity $||A||_{\text{Fro}}^2 = \text{tr}(A^{\top}A)$. Hence, we optimize equation 7 by alternating between U and V:

• U iteration: Define $M \coloneqq \sum_i B_i A_i V V^{\top} A_i^{\top} B_i^{\top}$. Parenthesizing this expression properly requires only O((m + n)r) storage/computation time. With this definition, we maximize $\operatorname{tr}(MUU^{\top})$ over U satisfying $U^{\top}U = I$. Since M is positive semidefinite, the optimum is to take U to be the r eigenvectors of M with largest eigenvalue, equivalent to an SVD problem.

to take U to be the r eigenvectors of M with largest eigenvalue, equivalent to an SVD problem. • V **iteration:** Define $N \coloneqq \sum_i A_i^\top B_i^\top UU^\top B_i A_i$. Similarly to the previous step, we take V to contain the r eigenvectors of N with largest eigenvalue, again solvable using an SVD.

This method decreases the objective in each step.

Case 2: Diagonal Σ_i . If constrain Σ_i to be diagonal, we interpret our objective function equation 1 as a "triple least squares" problem. We compute gradients:

759
760
$$\nabla_U \sum_i \|B_i A_i - U\Sigma_i V^\top\|_{\text{Fro}}^2 = 2 \sum_i (U\Sigma_i V^\top - B_i A_i) V\Sigma_i^\top$$

$$\nabla_V \sum_i^{\circ} \|B_i A_i - U\Sigma_i V^{\top}\|_{\text{Fro}}^2 = 2\sum_i^{\circ} (V\Sigma_i^{\top} U^{\top} - A_i^{\top} B_i^{\top}) U\Sigma_i$$

$$\nabla_{\Sigma_i} \sum_i \|B_i A_i - U\Sigma_i V^\top\|_{\text{Fro}}^2 = 2U^\top (U\Sigma_i V^\top - B_i A_i) V$$

These expressions suggest efficient $r \times r$ linear systems to solve for U, V:

$$U = \left(\sum_{i} B_{i}A_{i}V\Sigma_{i}^{\top}\right) \left(\sum_{i} \Sigma_{i}V^{\top}V\Sigma_{i}^{\top}\right)^{-1}$$

$$V = \left(\sum_{i} A_{i}^{\top}B_{i}^{\top}U\Sigma_{i}\right) \left(\sum_{i} \Sigma_{i}^{\top}U^{\top}U\Sigma_{i}\right)^{-1}.$$

$$V = \left(\sum_{i} A_{i}^{\top}B_{i}^{\top}U\Sigma_{i}\right) \left(\sum_{i} \Sigma_{i}^{\top}U^{\top}U\Sigma_{i}\right)^{-1}.$$

For Σ_i , we extract the diagonal from our gradient above:

776
777
778
779
diag
$$(U^{\top}U\Sigma_iV^{\top}V)_j = (U^{\top}U\Sigma_iV^{\top}V)_{jj}$$

 $= \sum_m (U^{\top}U)_{jm}\Sigma_{imm}(V^{\top}V)_{mj}$

779
780
$$= (U^{\top}U \circ V^{\top}V) \operatorname{diag}(\Sigma_{i})$$
781
$$= (U^{\top}D \land U) = \sum_{i=1}^{m} (U^{\top}D) = (A, U)$$

$$\operatorname{diag}(U^{\top}B_iA_iV)_j = \sum_m (U^{\top}B_i)_{jm}(A_iV)_{mj}$$

783
784
$$= \sum_{m}^{m} (U^{\top} B_i)_{jm} (V^{\top} A_i^{\top})_{jm}$$

 $= (U^{\top} B_i \circ V^{\top} A_i^{\top}) \mathbf{1}$ \implies diag $(\Sigma_i) = (U^{\top}U \circ V^{\top}V)^{-1}(U^{\top}B_i \circ V^{\top}A_i^{\top})\mathbf{1}$

Here \circ denotes the Hadamard product.

Combining these expressions, we use a simple coordinate descent algorithm cycling between the following three steps:

- 1. Solve for U
- 2. Solve for V
 - 3. Solve for the Σ_i 's

4. Optionally, normalize so $\sum_{i} \|\Sigma_i\|_{\text{Fro}}^2 = 1$

A.2 ADDITIONAL EIGENVALUE ITERATION ALGORITHM

For the first case in A.1, we introduce an alternative algorithm that eschews the use of SVD. This alternative is optimized for GPU execution, enabling tractable runs to convergence.

To derive this algorithm, we employ Lagrange multipliers to formulate the derived objective from equation 7:

$$U_{opt}, V_{opt} = \arg \max_{\substack{U^{\top}U=I\\VV^{\top}=I}} \sum_{i=1}^{n} \|U^{\top}B_{i}A_{i}V\|_{\text{Fro}}^{2},$$
(8)

yielding the expression

$$\Lambda = -\frac{1}{2} \| U^{\top} B_i A_i V \|_{\text{Fro}}^2 - \frac{1}{2} \operatorname{tr}(X^{\top} (I - U^{\top} U)) - \frac{1}{2} \operatorname{tr}(Y^{\top} (I - V^{\top} V)).$$
(9)

Taking the derivatives gives

$$\nabla_U \Lambda = -\sum_i B_i (A_i V) (V^\top A_i^\top) (B_i^\top U) + UX$$
(10)

$$\nabla_V \Lambda = -\sum_i A_i^\top (B_i^\top U) (U^\top B_i) (A_i V) + VY$$
(11)

Setting these derivatives to zero shows

$$\sum_{i} B_i(A_i V) (V^{\top} A_i^{\top}) (B_i^{\top} U) = UX$$
(12)

$$\sum_{i} A_i^{\top} (B_i^{\top} U) (U^{\top} B_i) (A_i V) = VY.$$
(13)

Here, one can show that the Lagrange multiplier matrices X and Y are diagonal and nonnegative, since the problem reduces to an eigenvalue problem when either U or V is fixed; this is essentially the argument behind the alternating algorithm in Appendix A. Hence, taking inspiration from classical eigenvalue iteration, we use the following updates to improve our estimates of U and V:

$$U_0^{(k+1)} \leftarrow \sum_i B_i(A_i V^{(k)})((V^{(k)})^\top A_i^\top)(B_i^\top U^{(k)})$$
(14)

$$V_0^{(k+1)} \leftarrow \sum_i A_i^{\top} (B_i^{\top} U^{(k)}) ((U^{(k)})^{\top} B_i) (A_i V^{(k)})$$
(15)

$$U^{(k+1)} \leftarrow \text{orthogonalize}(U_0^{(k+1)}) \tag{16}$$

$$V^{(k+1)} \leftarrow \text{orthogonalize}(V_0^{(k+1)}) \tag{17}$$

Here, the function orthogonalize orthogonalizes the columns of a matrix, e.g. by using the Qpart of the reduced-size QR factorization. Although we lack a formal convergence proof, in practice we find that this method reliably reaches a local optimum of our problem.

By executing matrix operations in the specified sequence, these computations can be rapidly performed on GPUs. Note the expressions above are parenthesized to avoid constructing a large matrix product as an intermediate computation.

A.3 CLUSTERING ALGORITHM

Initialization: We run joint diagonalization with a single U, V then perform k-means with $|C_j|$ clusters on the space of Σ_i 's. This gives us our first clusters and we can use random initialization U_j, V_j for each cluster but the Σ_i can be maintained as initialization.

Step 1: Using the alternating JD algorithms from earlier in this section, we optimize the problem $\min_{U_j, V_j, \Sigma_i} \sum_{i \in C_i} ||B_i A_i - U_j \Sigma_i V_j^\top||_F^2$ for each *j* independently.

Step 2: New cluster assignment for $i : \min_j \min_{\Sigma_i} ||B_i A_i - U_j \Sigma_i V_j^{\top}||_F^2$. If any assignment changes we go to Step 1, else we have converged.

B PROOF OF THEOREM 1

Proof. For the lower bound, note that by Jensen's inequality,

$$\sum_{i=1}^{n} \| U^{\top} B_i A_i V \|_{\text{Fro}}^2 \ge \left\| U^{\top} \sum_{i=1}^{n} B_i A_i V \right\|_{\text{Fro}}^2$$

for any U, V. Hence,

$$\sup_{U,V \in \text{St}(k,d)} \sum_{i=1}^{n} \|U^{\top} B_{i} A_{i} V\|_{\text{Fro}}^{2} \ge \sup_{U,V \in \text{St}(k,d)} \left\|U^{\top} \sum_{i=1}^{n} B_{i} A_{i} V\right\|_{\text{Fro}}^{2}.$$
 (18)

By the definition of singular value decomposition, the right hand side of equation 18 is maximized with U, V being the top r singular vectors of $\sum_{i=1}^{n} B_i A_i$, yielding $\left\| U^{\top} \sum_{i=1}^{n} B_i A_i V \right\|_{\text{Fro}}^2 = \sum_{i=1}^{r} \bar{\sigma}_i^2$. Recalling that $\Sigma_i = U^{\top} B_i A_i V$ yields the lower bound.

For the upper bound, recall that $\Sigma_i = U^{\top} B_i A_i V$. Rearranging,

$$\operatorname{vec}(\Sigma_i) = (V^\top \otimes U^\top)\operatorname{vec}(B_i A_i)$$

Define

868

869 870 871

872 873

875

876 877 878

884

885 886

893

894 895

896

897

898

899

900

901

902

903

904 905

906

$$\overline{\Sigma} \coloneqq [\operatorname{vec}(\Sigma_1), \dots, \operatorname{vec}(\Sigma_n)].$$

By our previous simplification,

$$\bar{\Sigma} = (V^\top \otimes U^\top)L.$$

Now

$$\sum_{i=1}^{n} \|\Sigma_i\|_{\mathrm{Fro}}^2 = \|\bar{\Sigma}\|_{\mathrm{Fro}}^2 = \mathrm{tr}\left(((V \otimes U)(V \otimes U)^{\top})(LL^{\top})\right)$$

Since U, V are orthogonal and size $d \times r$, the top r^2 eigenvalues of the symmetric matrix $(V \otimes U)(V \otimes U)^{\top}$ will be equal to 1, and the rest will equal 0. The eigenvalues of the symmetric matrix LL^{\top} will be equal to the squared singular values of L. We can then apply the Von Neumann trace inequality to obtain the upper bound.

The last statement follows from the Pythagorean theorem and the fact that the Σ_i is a projection of $B_i A_i$ to the U, V subspace.

Note that we have only used the fact that the matrix $(V \otimes U)$ has singular values equal to 1; we have not used the fact that it has Kronecker product structure. On the other hand, each vector $vec(B_iA_i)$ is a sum of r_i Kronecker products and cannot be expressed as a Kronecker product. As a result, while the upper bound in the Von Neumann trace inequality is achieved if the eigenvectors of the two matrices align, the Kronecker product structure is a severe constraint and the upper bound we have provided is generous.

C TRAINING LORAS

We trained LoRA adapters on 500 natural instruction tasks (Wang et al., 2022) using Mistral-7B-Instruct-v0.2 (Jiang et al., 2023a) as the base model. All LoRA adapters were configured with a rank of 16, i.e., $\forall i, r_i = 16$. We selected 10 diverse tasks manually for consistent evaluation across experiments and randomly sampled an additional 490 tasks, resulting in a total of 500 tasks. These tasks were exclusively in English (both input and output), ensuring higher quality and thorough review (Wang et al., 2022). Each task dataset was divided into training, validation, and test sets (80-10-10). Hyperparameters, such as early stopping, were tuned using the validation loss. Evaluation on the test sets demonstrated that LoRA consistently outperformed the base model in terms of both Rouge scores and loss metrics (see Table 1).

Table 2: Main Evaluation Tasks

| Task Number | Name | Туре | Domain |
|-------------|--|---------------------|---------------------|
| task280 | stereoset_classification_stereotype_type | classification | stereoset |
| task190 | snli_classification | snli | image captions |
| task391 | causal_relationship | commonsense | cause and effect |
| task290 | tellmewhy_question_answerability | answerability | story |
| task1391 | winogrande_easy_answer_generation | commonsense | social and physical |
| task1342 | amazon_us_reviews_title | title generation | amazon reviews |
| task442 | com_qa_paraphrase_question_generation | question generation | wikipedia |
| task620 | ohsumed_medical_subject_headings_answer_generation | keyword tagging | scientific |
| task1598 | nyc_long_text_generation | data to text | restaurants |
| task039 | qasc_find_overlapping_words | overlap extraction | natural science |

In Table 3 we include all 500 tasks that were used.

| 9 | 1 | 8 | |
|---|---|---|--|
| 9 | 1 | 9 | |

Table 3: List of Tasks

| Task ID | Description | Task ID | Description | Task ID | Description |
|-----------------------|--|----------------------|---|----------------------|---|
| tack 280 | etarteret description starantime tune | task100 | anti algorification | tack 201 | caucal relationship |
| task200 | tellmewhy.question answerability | task1391 | winogrande easy answer generation | task1342 | amazon us reviews title |
| task442 | com qa paraphrase question generation | task620 | ohsumed medical subject headings answer generation | task1598 | nyc long text generation |
| task039 | qase find overlapping words | task769 | qed summarization | task1448 | disease entity extraction ncbi dataset |
| task247 | dream answer generation | task513 | argument stance classification | task875 | emotion classification |
| task1551 | every the element from kth element | task583 | udeps eng coarse pos tagging | task1334 | head ga answer generation |
| task270 | csrg counterfactual context generation | task1487 | organism substance extraction anem dataset | task679 | hope edi english text classification |
| task456 | matres intention classification | task385 | socialiqa incorrect answer generation | task1607 | ethos text classification |
| task278 | stereoset antistereotype sentence generation | task022 | cosmosqa passage inappropriate binary | task210 | logic2text structured text generation |
| task137 task1378 | quare correct answer generation | task1194 | kth largest element | task029 task1529 | scitail1 1 classification |
| task453 | swag.answer generation | task102 | commongen sentence generation | task460 | qasper answer generation |
| task1204 | atomic classification hinderedby | task1384 | deal or no dialog classification | task1572 | samsum summary |
| task699 | mmmlu high school biology answer generation | task1631 | openpi answer generation | task1722 | civil comments threat classification |
| task580 task1283 | socialiga answer generation | task605 | longest common subsequence in lists | task1152 task723 | bard analogical reasoning causation |
| task084 | babilit supporting fact identification | task201 | mnli neutral classification | task956 | leetcode strong password check |
| task167 | strategyqa question generation | task1192 | food flavor profile | task300 | storycloze order generation |
| task1714 | convai3 sentence generation | task388 | torque token classification | task516 | senteval conjoints inversion |
| task127 task322 | scan action command all generation | task362 task607 | spolin yesand response classification | task1158 task1566 | bard analogical reasoning manipulating items |
| task076 | splash correcting SOL mistake | task1451 | drug dose extraction | task1135 | xcsr en commonsense mc classification |
| task341 | winomt gender anti classification | task267 | concatenate and reverse elements from i to j | task1720 | civil comments toxicity classification |
| task1452 | location entity extraction btc corpus | task131 | scan action command long generation | task685 | mmmlu clinical knowledge answer generation |
| task727 | mmmlu prehistory answer generation | task1590 | diplomacy text generation | task1731 | quartz question answering |
| task047 task1326 | answering science questions | task929 task615 | products reviews classification movies answer generation | task1592 task1216 | vahoo answers topics classification |
| task689 | mmalu college mathematics answer generation | task1156 | bard analogical reasoning tools | task1657 | gooaq question generation |
| task833 | poem sentiment classification | task1206 | atomic classification isbefore | task1151 | swap max min |
| task244 | count elements in set union | task1562 | zest text modification | task043 | essential terms answering incomplete questions |
| task044 | essential terms identifying essential words | task722 | mmmlu random topic answer generation | task183 | rhyme generation |
| task503 task616 | colored assumption | task155 task1724 | count nouns and verbs | task353 task288 | casino negotiation elicit preference classification |
| task092 | check prime classification | task707 | mmmlu high school microeconomics answer generation | task577 | curiosity dialogs classification |
| task742 | Ihoestq frequency answer generation | task706 | mmmlu high school mathematics answer generation | task1401 | obqa sentence generation |
| task1393 | superglue copa text completion | task1198 | atomic classification owant | task966 | ruletaker fact checking from context |
| task219 | rocstories title answer generation | task1211 | atomic classification hassubevent | task050 | multire answerability |
| task494 task068 | abductiventi incorrect answer generation | task15/9 task566 | quare incorrect answer generation | task170 task333 | break decompose questions |
| task593 | sciq explanation generation | task667 | mmmlu business ethics answer generation | task130 | scan action command long generation |
| task161 | count words containing letter | task507 | position of numerical elements in list | task1502 | hatexplain classification |
| task505 | count numerical elements in list | task633 | dbpedia 14 answer generation | task1645 | medical question pair classification |
| task1486 | cell extraction anem dataset | task1146 | country capital | task1380 | quarel correct option generation |
| task1088 task1294 | array of products wiki as answer verification | task033 task080 | winogrande answer generation | task085 task489 | unnatural addsub arithmetic mwsc question generation |
| task1721 | civil comments obscenity classification | task1713 | convai3 sentence generation | task721 | mmmlu medical genetics answer generation |
| task1403 | check validity date mmddyyyy | task746 | yelp restaurant review classification | task728 | mmmlu professional accounting answer generation |
| task889 | goemotions classification | task1583 | bless meronym classification | task1665 | trianglecopa question generation |
| task708 | mmmlu high school physics answer generation | task1419 | mathqa gain | task963 | librispeech asr next word prediction |
| task454 task579 | swag incorrect answer generation | task308 task753 | Jeopardy answer generation all | task828 task1404 | copa cause effect commonsense date conversion |
| task1201 | atomic classification xintent | task901 | freebase ga category question generation | task1567 | propara question generation |
| task1319 | country by barcode prefix | task858 | inquisitive span detection | task1200 | atomic classification xeffect |
| task492 | mwse incorrect answer generation | task675 | google wellformed query sentence generation | task094 | conala calculate mean |
| task1506 | celebrity minimal dob span | task694 | mmmlu econometrics answer generation | task614 | glucose cause event detection |
| task1390 task457 | wschxed coreference | task1355 task1565 | sent comp summarization | task/14 task834 | mmmlu human sexuality answer generation mathdataset classification |
| task642 | esnli classification | task732 | mmmlu public relations answer generation | task1605 | ethos text classification |
| task326 | jigsaw obscene classification | task1292 | yelp review full text categorization | task716 | mmmlu jurisprudence answer generation |
| task1479 | organization entity extraction btc corpus | task1147 | country currency | task153 | tomqa find location hard clean |
| task1495 | adverse drug event classification | task1196 | atomic classification oeffect | task1489 | sarcasmdetection tweet classification |
| task 294 task 1197 | atomic classification oreact | task157 task754 | count vowers and consonants syamp common-division question answering | task147 task1599 | ars argument similarity gay marriage smcalflow classification |
| task1420 | mathqa general | task1285 | kpa keypoint matching | task587 | amazonfood polarity correction classification |
| task1338 | peixian equity sentiment classifier | task116 | com2sense commonsense reasoning | task713 | mmmlu human aging answer generation |
| task431 | senteval object count | task067 | abductivenli answer generation | task934 | turk simplification |
| task01/ | amazonreview category text generation | task090 | mmmiu elementary mathematics answer generation | task840 | ambigga text generation |
| task1398 | obga question generation | task1518 | limit answer generation | task628 | xlwic different meaning sentence generation |
| task1286 | openbookga question answering | task1596 | event2mind text generation 2 | task298 | storycloze correct end classification |
| task645 | summarization | task903 | deceptive opinion spam classification | task594 | sciq question generation |
| task413 | mickey en sentence perturbation generation | task719 | mmmlu management answer generation | task672 | nummersense |
| task 1418 | bless semantic relation classification | task4/5 | yelp polarity classification | task357 | casino negotiation small talk classification |
| task750 | adua multiple choice answering | task1320 | country domain tld | task034 | winogrande question modification object |
| task692 | mmnlucomputer security answer generation | task1406 | kth smallest element | task119 | semeval 2019 task10 geometric mathematical answer |
| task211 | logic2text classification | task363 | sst2 polarity classification | task1087 | two number sum |
| task083 | babi t1 answer generation | task1385 | anli r1 entailment | task1308 | amazonreview category classification |
| task 1656 | gooaq answer generation | task892 task1207 | gap reverse coreterence resolution | task499 task564 | extract and add numbers from list |
| task 1409 | ijesaw identity attack classification | task120/ | collatz conjecture | task304 | ascouse classification |
| task1520 | qa srl answer generation | task703 | mmmlu high school geography answer generation | task318 | stereoset gender classification |
| task366 | synthetic return primes | task335 | hateeval aggressive classification en | task600 | longest common substring in two strings |
| task477 | cls english dvd classification | task138 | detoxifying-lms classification fluency | task291 | semeval 2020 task4 commonsense validation |
| task074 | squal L question generation | task1389 | hellaswag completion | task192 | hotpotqa sentence generation |
| task666 | crows-pairs stereotype classification | task1582 | xquau en question generation bless hypernym generation | task290 task1728 | sorycroze correct end classification web nlg data to text |
| task701 | mmmlu high school computer science answer generation | task275 | enhanced wsc paraphrase generation | task107 | splash question to SQL |
| task079 | conala concat strings | task1157 | bard analogical reasoning rooms for containers | task1167 | penn treebank coarse pos tagging |
| task403 | creak commonsense inference | task359 | casino negotiation vouch fair classification | task517 | emo classify emotion of dialogue |
| task351 | winomt gender identifiability anti classification | task964 | hbrispeech asr text auto completion | task904 | hate speech offensive classification |
| task 1509 | evalution antonyms | tasko/J | max element lists | task0.00 | arc easy answer generation |

| 966 967 | We use Huggingface (Wolf et al., 2020) in our implementation. For the base model, we use quanti- zation with configuration: |
|------------|--|
| 968 | |
| 969 | BitsAndBytesConfig(|
| 970 | load_in_4bit=True, |

| 510 | roud_rn_nore | iiuc, |
|-----|---------------|---------------------|
| 971 | bnb_4bit_use_ | _double_quant=True, |
| | bnb_4bit_quar | nt_type="nf4", |

- 972 bnb_4bit_compute_dtype=torch.bfloat16, 973) 974 975 and LoRA configuration: 976 LoraConfig(977 r=16, 978 lora_alpha=32, 979 target_modules=["q_proj", "k_proj", "v_proj"], 980 lora_dropout=0.05, 981 bias="none", 982 task_type="CAUSAL_LM", 983 init_lora_weights=init_lora_weights, 984)
- 985 986 987

D AVOIDING BATCHED MATRIX MULTIPLICATION (BMM)

Fast LoRA (Wen & Chaudhuri, 2024) aims to alleviate the batched matrix multiplication (BMM)
bottleneck when serving many LoRAs. They propose an adapter parameterization that replaces addition with elementwise multiplication, avoiding BMM and improving LoRA throughput at lower
ranks. Our JD LoRA formulation also circumvents or heavily reduces the impact of BMM as discussed below, and both individual and joint compression methods can be applied to Fast LoRAs.

In the envisioned deployment scenario, a service provider hosts a large collection of LoRAs. Upon receiving a request, each user specifies both the input data and the desired LoRA identifier. The provider then processes the base model augmented with the specified LoRA for each user's data. As a provider is batching a collection of requests for GPU parallelization, they can expect to frequently have more than one unique LoRA identifier per batch.

Traditionally, a specific LoRA is integrated into the base model by transforming $W_0 \rightarrow W_0 + B_i A_i$. Serving multiple LoRAs conventionally would necessitate maintaining and executing a separate copy of the base model for each LoRA, bringing substantial computational overhead. Alternatively, the computation for $W_0 x$ and $B_i A_i x$ can be performed independently and subsequently merged. This strategy necessitates only a single instance of $W_0 x$ computation and storage of LoRA-specific parameters rather than the entire base model.

Consider the batch processing of **BAx**, where boldface indicates that B_i , A_i are stacked into tensors of dimensions $(b \times m \times r)$ and $(b \times r \times n)$ respectively, with batched data x shaped $(b \times l \times n)$:

1008 1009

1014

1015

1016

$$\mathbf{Ax} \leftrightarrow (b \times r \times n) \times (b \times l \times n) \rightarrow (b \times l \times r) \text{ bmm}$$
$$\mathbf{B}(\mathbf{Ax}) \leftrightarrow (b \times m \times r) \times (b \times l \times r) \rightarrow (b \times l \times m) \text{ bmm}.$$

Here, "bmm" denotes batched matrix multiplication, a known bottleneck in both throughput and latency. Consider the corresponding operations for our joint compression scheme, $U\Sigma V^{\top}x$:

1012 1013 $V^{\top} \mathbf{x} \leftrightarrow (\tilde{r} \times n) \times (b \times l \times n) \rightarrow (b \times l \times \tilde{r})$ broadcasted

$$\mathbf{\Sigma}(V^{\top}\mathbf{x}) \leftrightarrow (b \times \tilde{r}) \times (b \times l \times \tilde{r}) \rightarrow (b \times l \times \tilde{r})$$
 broadcasted

$$U(\mathbf{\Sigma}V^{\top}\mathbf{x}) \leftrightarrow (m \times \tilde{r}) \times (b \times l \times \tilde{r}) \rightarrow (b \times l \times m)$$
 broadcasted

In our optimized setup, batched matrix multiplications can be completely circumvented if the Σ_i matrices are diagonal. If not, given that $\tilde{r} \ll m, n$, any required batched matrix multiplication remains computationally inexpensive.

1020

E GPU MEMORY USAGE COMPUTATION FOR JD COMPRESSION

1021 1022

The GPU memory consumption is primarily influenced by the number of parameters that need to be
 stored and processed during inference. In this section, we introduce the detail of how we compute
 the GPU consumption of our method, and how we find the number of vLLM multi-LoRA that share
 the same GPU utilization.

1026 • D: Hidden dimension size (e.g., D = 4098). 1027 • r: Rank of the shared basis matrices for compression (e.g., r = 16, 32, 64). 1028 • N: Maximum number of LoRA modules being served simultaneously (max_lora_num). 1029 • c: Number of clusters in our clustering method (e.g., c = 7, 10, 25). 1030 In Figure 1, we use different JD-compression settings for serving different number of unique LoRAs. 1031 Specifically: 1032 • Serving 4 unique LoRAs: 1033 Ours: rank 16 JD-Full. 1034 vLLM multiLoRA baseline: max-gpu-lora = 2. 1035 Serving 8 unique LoRAs: 1036 Ours: rank 16 JD-Full. 1037 vLLM multiLoRA baseline: max-gpu-lora = 2. • Serving 16 unique LoRAs: 1039 Ours: rank 32 JD-Full. 1040 vLLM multiLoRA baseline: max-gpu-lora = 3. 1041 • Serving 32 unique LoRAs: Ours: rank 64 JD-Full. 1043 vLLM multiLoRA baseline: max-gpu-lora = 5. Serving 64 unique LoRAs: 1044 Ours: rank 64 JD-Full. 1045 vLLM multiLoRA baseline: max-gpu-lora = 6. 1046 Serving 128 unique LoRAs: 1047 Ours: 7 clusters, rank 16 JD-Full. 1048 vLLM multiLoRA baseline: max-gpu-lora = 8. 1049 Serving 256 unique LoRAs: 1050 Ours: 10 clusters, rank 16 JD-Full. 1051 vLLM multiLoRA baseline: max-gpu-lora = 10. 1052 • Serving 512 unique LoRAs: 1053 Ours: 25 clusters, rank 16 JD-Full. 1054 vLLM multiLoRA baseline: max-gpu-lora = 26. Serving 1024 unique LoRAs: 1055 Ours: 7 clusters, rank 64 JD-Full. 1056 vLLM multiLoRA baseline: max-gpu-lora = 60. 1057 1058 BASELINE GPU MEMORY USAGE E.1 1059 The baseline for our comparison is the standard LoRA method with a rank of 16. The total parameter 1061 count for the baseline is given by: 1062 1063 $Params_{baseline} = D \times 2 \times 16.$ 1064 This accounts for the parameters in the LoRA-adapted layers, where the factor of 2 represents the weights and biases. 1067 1068 E.2 GPU MEMORY USAGE FOR JD FULL METHOD 1069 1070 For the Joint Decomposition (JD) Full method without clustering, the total parameter count is: 1071 1072 Params_{ID Full} = $D \times 2 \times r + N \times r^2$. 1074 • $D \times 2 \times r$: Parameters for the base model adapted with rank-r LoRA. 1075 • $N \times r^2$: Additional parameters introduced by each of the N LoRA modules, each of size $r \times r$. The GPU memory usage ratio relative to the baseline is: 1077 1078 $\label{eq:GPU} \text{GPU Usage Ratio}_{\text{JD}_\text{Full}} = \frac{\text{Params}_{\text{JD}_\text{Full}}}{\text{Params}_{\text{baseline}}} = \frac{D \times 2 \times r + N \times r^2}{D \times 2 \times 16}.$ 1079

1080 E.3 GPU MEMORY USAGE FOR CLUSTERING METHOD

When employing clustering, the parameter count changes due to the addition of cluster-specificparameters:

1084 1085 1086

1091 1092

1093

1098

1099

1108

 $Params_{Clustering} = D \times 2 \times r \times c + N \times (r^2 + 1).$

• $D \times 2 \times r \times c$: Parameters for the base model adapted with rank-r LoRA across c clusters.

• $N \times (r^2 + 1)$: Additional parameters for each LoRA module and cluster assignments.

1089 The GPU memory usage ratio is:

$$\text{GPU Usage Ratio}_{\text{Clustering}} = \frac{\text{Params}_{\text{Clustering}}}{\text{Params}_{\text{baseline}}} = \frac{D \times 2 \times r \times c + N \times (r^2 + 1)}{D \times 2 \times 16}.$$

1094 1095 E.4 PUNICA

> In our vLLM experiments, we specifically utilized the Punica kernel for implementing multi-LoRA, applying our approach in conjunction with Punica's capabilities. Our custom function, add_lora_slice_with_sigma, implements the following key steps:

Initialize Buffers: Creates temporary storage for intermediate calculations if not already provided.

- 2. Apply Matrix A: Transforms x using matrix A, storing the result in buffer.
- 3. Apply Matrix Sigma: Further transforms buffer using Sigma, storing the result in buffer_sigma.
- 4. Apply Matrix B and Update y: Finally, transforms buffer_sigma using B, applies scaling, and updates a slice of y in place.
- Below is the pseudocode for add_lora_slice_with_sigma, illustrating the integration:

Listing 1: Pseudocode for 'add_lora_slice_with_sigma'

```
1109
       Function add_lora_slice_with_sigma(y, x, wa_t_all, wb_t_all, wsigma_t_all
1110
           , indices, layer_idx, scale, y_offset, y_slice_size, buffer=None):
1111 2
           # Initialize buffers if not provided
1112 3
           if buffer is None:
               buffer = create_tensor(shape=(x.size(0), R), dtype=float32)
1113 4
               buffer_sigma = create_tensor(shape=(buffer.size(0), R), dtype=
1114 <sup>5</sup>
                   float.32)
1115
           # Step 1: Apply matrix A
1116 7
           dispatch_bgmv_low_level(buffer, x, wa_t_all, indices, layer_idx,
1117
               scale=1.0)
           # Step 2: Apply matrix Sigma
1118 8
           dispatch_bgmv_low_level(buffer_sigma, buffer, wsigma_t_all, indices,
1119 <sup>9</sup>
               layer_idx, scale=1.0)
1120
           # Step 3: Apply matrix B and update y slice
1121<sub>11</sub>
           dispatch_bgmv_low_level(y, buffer_sigma, wb_t_all, indices, layer_idx
1122
               , scale, y_offset, y_slice_size)
       End Function
112312
```

1124 1125

F SELECTING NUMBER OF CLUSTERS

1126 1127

To identify optimal hyperparameters for the clusters compression method, we analyzed the relationship between reconstruction error and the parameter saved ratio for a single LoRA module, as shown in Figure 5. By comparing the results across different numbers of Low-Rank Adaptation (LoRA) configurations (100 and 500, depicted in subfigures 5a and 5b), we were able to observe the tradeoff between model size reduction and reconstruction accuracy. Based on these findings, we selected the rank and number-of-clusters hyperparameters that effectively balance these two objectives. The chosen settings were then used to conduct full-scale experiments.





1141(a) Recon. Error vs Parameter Saved Ratio for 100(b) Recon.1142LoRAsLoRAs1143

(b) Recon. Error vs Parameter Saved Ratio for 500 LoRAs

Figure 5: Comparison of reconstruction error against the parameter saved ratio for different numbers of LoRA configurations for a single LoRA module. The left subplot shows results for 100 LoRAs, while the right subplot displays results for 500 LoRAs. These plots illustrate the trade-off between reconstruction accuracy and compression efficiency, providing insights into optimal parameter settings for compression.

1149

1150 G ADDITIONAL RESULTS

1152

1156

1158

This section elaborates on the results that underpin the figures presented in the main text and showcases a consistent correlation across various evaluation metrics. Additionally, we assess the significance of achieving convergence and the performance of compression on new unseen LoRA models.

1157 G.1 RELATIVE ROUGE-L PERFORMANCE AND COMPRESSION RATE

Table 4 presents comprehensive results from the experiments underlying Figure 2a for each evaluation task. Additionally, we incorporate results using the Ties-merging benchmark (Yadav et al., 2023b), which consolidates all LoRA-adapters into a single adapter of identical configuration and parameter count; this integration significantly compromises performance.

1163

1165

1164 G.2 Absolute Rouge-L Performance and Compression Rate

1166Table 5 provides the full results behind Table 4, but with Rouge-L scores instead of relative perfor-1167mance compared to LoRA.

1168

1173

1169 G.3 RELATIVE ROUGE-1 PERFORMANCE AND COMPRESSION RATE

Table 6 provides full results for relative performance of Rouge-1, which shows the same trends as the results for relative performance of Rouge-L (Table 4).

1174 G.4 Absolute Rouge-1 Performance and Compression Rate

Table 7 provides full results for absolute performance of Rouge-1, which shows the same trends as the results for absolute performance of Rouge-L (Table 5).

- 1178 1179
- 1180 G.5 RELATIVE EXACT-MATCH PERFORMANCE AND COMPRESSION RATE

Table 8 provides full results for relative performance of exact-match, which shows the same trends as the results for relative performance of Rouge-L (Table 4).

1184

- 1185 G.6 LOSS AND COMPRESSION RATE
- 1187Table 9 provides full results for test loss (cross-entropy), which shows the same trends as the results
for relative performance of Rouge-L (Table 4).

| TIES SVD 10 diagonal (D) | SVD 2 500 500 SVD 2 SVD 4 SVD 4 SVD 4 SVD 4 SVD 16 16 D 32 D 64 D 128 D 256 D 16 F 32 F 64 F 128 F | $\begin{array}{c} task039 \\ 0.26 \pm 0.00 \\ 1.00 \pm 0.00 \\ 0.59 \pm 0.00 \\ 0.55 \pm 0.00 \\ 0.57 \pm 0.00 \\ 0.37 \pm 0.00 \\ 0.37 \pm 0.00 \\ 1.00 \pm 0.00 \\$ | $\begin{array}{c} task 190 \\ 0.02 \pm 0.00 \\ 1.00 \pm 0.00 \\ 0.57 \pm 0.02 \\ 0.41 \pm 0.00 \\ 0.26 \pm 0.00 \\ 0.26 \pm 0.00 \\ 1.07 \pm 0.02 \\ 1.04 \pm 0.01 \\ 1.02 \pm 0.01 \\ 1.00 \pm 0.00 \\ 1.01 \pm 0.01 \\ 1.03 \pm 0.01 \\ 1.01 \pm 0.01 \\ 1.00 \pm 0.00 \\ 0.00 $ | $\begin{array}{c} task280\\ 0.19 \pm 0.00\\ 1.00 \pm 0.00\\ 0.45 \pm 0.04\\ 0.18 \pm 0.05\\ 0.20 \pm 0.05\\ 0.01 \pm 0.00\\ 1.00 \pm 0.00\\ 1.$ | $\begin{array}{c} task290\\ 0.42 \pm 0.00\\ 1.00 \pm 0.00\\ 0.10 \pm 0.01\\ 0.03 \pm 0.01\\ 0.01 \pm 0.02\\ 0.00 \pm 0.00\\ 1.00 \pm 0.00\\ 0.99 \pm 0.00\\ \end{array}$ | $\begin{array}{ c c c c c c c c c c c c c c c c c c c$ | $\begin{array}{c} task442 \\ 0.47 \pm 0.00 \\ 1.00 \pm 0.00 \\ 0.31 \pm 0.00 \\ 0.29 \pm 0.00 \\ 0.29 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} task620\\ 0.11\pm0.00\\ 1.00\pm0.00\\ 0.65\pm0.00\\ 0.65\pm0.00\\ 0.64\pm0.00\\ 0.57\pm0.00\\ 1.00\pm0.01\\ 0.99\pm0.02\\ 1.01\pm0.00\\ 1.00\pm0.01\\ 0.90\pm0.02\\ 1.01\pm0.00\\ 0.$ | $\begin{array}{c} task 1342 \\ 0.23 \pm 0.00 \\ 1.00 \pm 0.00 \\ 0.57 \pm 0.00 \\ 0.57 \pm 0.02 \\ 0.37 \pm 0.00 \\ 1.00 \pm 0.10 \\ 0.99 \pm 0.08 \\ 1.00 \pm 0.01 \end{array}$ | $\begin{array}{c} task 1391 \\ 0.19 \pm 0.00 \\ 1.00 \pm 0.00 \\ 0.82 \pm 0.01 \\ 0.32 \pm 0.04 \\ 0.01 \pm 0.00 \\ 0.01 \pm 0.00 \\ 1.00 \pm 0.01 \\ 1.09 \pm 0.01 \\ 1.01 \pm 0.01 \end{array}$ | $\begin{array}{c} task 1598 \\ 0.77 \pm 0.00 \\ 1.00 \pm 0.00 \\ 0.85 \pm 0.00 \\ 0.84 \pm 0.00 \\ 0.82 \pm 0.00 \\ 0.43 \pm 0.00 \\ 1.00 \pm 0.01 \\ 1.00 \pm 0.01 \\ 1.01 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.28 \pm 0.21 \\ 1.00 \pm 0.00 \\ 0.62 \pm 0.23 \\ 0.48 \pm 0.28 \\ 0.44 \pm 0.30 \\ 0.31 \pm 0.26 \\ \hline 1.00 \pm 0.04 \\ 1.00 \pm 0.01 \\ \hline \end{array}$ | 1.00/1.00 0.00/0.00 1.00 / 1.00 1.00 / 1.00 1.00 / 1.00 1.00 / 1.00 1.00 / 1.00 0.88 / 0.88 0.75 / 0.75 0.50 / 0.50 |
|---|--|---|---|---|---|---|---|--|--|--|--|--|--|
| TIES SVD 10 diagonal (D) | base lora 10 50 500 SVD 2 SVD 4 SVD 4 SVD 16 16 D 32 D 64 D 128 D 256 D 16 F 32 F 64 F 128 F | $\begin{array}{c} 0.26\pm0.00\\ 1.00\pm0.00\\ 0.59\pm0.00\\ 0.57\pm0.00\\ 0.37\pm0.00\\ 0.37\pm0.00\\ 1.00\pm0.00\\ 1.00\pm0.00\\ 1.00\pm0.00\\ 1.00\pm0.00\\ 1.00\pm0.00\\ 1.00\pm0.00\\ 1.00\pm0.00\\ 1.00\pm0.00\\ 1.02\pm0.00\\ 1.02$ | $\begin{array}{c} 0.02 \pm 0.00 \\ 1.00 \pm 0.00 \\ 0.57 \pm 0.02 \\ 0.41 \pm 0.00 \\ 0.40 \pm 0.00 \\ 0.26 \pm 0.00 \\ 1.07 \pm 0.02 \\ 1.04 \pm 0.01 \\ 1.02 \pm 0.01 \\ 1.02 \pm 0.01 \\ 1.05 \pm 0.01 \\ 1.03 \pm 0.01 \\ 1.03 \pm 0.01 \\ 1.01 \pm 0.01 \\ 1.01 \pm 0.01 \\ 1.00 \pm 0.00 \\ \end{array}$ | $\begin{array}{c} 0.19 \pm 0.00 \\ 1.00 \pm 0.00 \\ 0.45 \pm 0.04 \\ 0.18 \pm 0.05 \\ 0.20 \pm 0.05 \\ 0.01 \pm 0.00 \\ 1.00 \pm 0.$ | $\begin{array}{c} 0.42 \pm 0.00 \\ 1.00 \pm 0.00 \\ 0.10 \pm 0.01 \\ 0.03 \pm 0.01 \\ 0.01 \pm 0.02 \\ 0.00 \pm 0.00 \\ 1.00 \pm 0.01 \\ 0.99 \pm 0.00 \\ \end{array}$ | $ \begin{array}{ c c c c c c c c c c c c c c c c c c c$ | $\begin{array}{c} 0.47 \pm 0.00 \\ 1.00 \pm 0.00 \\ \hline 0.47 \pm 0.00 \\ 0.31 \pm 0.00 \\ 0.33 \pm 0.00 \\ 0.29 \pm 0.00 \\ \hline 0.98 \pm 0.01 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.11 \pm 0.00 \\ 1.00 \pm 0.00 \\ \hline 0.69 \pm 0.01 \\ 0.65 \pm 0.00 \\ 0.64 \pm 0.00 \\ 0.57 \pm 0.00 \\ \hline 1.00 \pm 0.01 \\ 0.99 \pm 0.02 \\ 1.01 \pm 0.00 \\ \hline 1.00 \pm 0.01 \\ \hline 0.90 \pm 0.02 \\ \hline 1.00 \pm 0.00 \\ \hline 0.00 $ | $\begin{array}{c} 0.23 \pm 0.00 \\ 1.00 \pm 0.00 \\ \hline \\ 0.57 \pm 0.00 \\ 0.62 \pm 0.00 \\ 0.57 \pm 0.02 \\ 0.37 \pm 0.00 \\ \hline \\ 1.00 \pm 0.10 \\ 0.99 \pm 0.08 \\ 1.00 \pm 0.01 \\ \hline \end{array}$ | $\begin{array}{c} 0.19 \pm 0.00 \\ 1.00 \pm 0.00 \\ \hline 0.82 \pm 0.01 \\ 0.32 \pm 0.04 \\ 0.01 \pm 0.00 \\ \hline 0.01 \pm 0.00 \\ \hline 1.00 \pm 0.01 \\ 1.01 \pm 0.01 \\ \hline \end{array}$ | $\begin{array}{c} 0.77 \pm 0.00 \\ 1.00 \pm 0.00 \\ \hline 0.85 \pm 0.00 \\ 0.84 \pm 0.00 \\ 0.82 \pm 0.00 \\ 0.43 \pm 0.00 \\ \hline 1.00 \pm 0.01 \\ 1.00 \pm 0.01 \\ 1.01 \pm 0.00 \\ \hline \end{array}$ | $\begin{array}{c} 0.28 \pm 0.21 \\ 1.00 \pm 0.00 \\ \hline \\ 0.62 \pm 0.23 \\ 0.48 \pm 0.28 \\ 0.44 \pm 0.30 \\ 0.31 \pm 0.26 \\ \hline \\ 1.00 \pm 0.04 \\ 1.00 \pm 0.03 \\ 1.00 \pm 0.01 \\ \hline \end{array}$ | 1.00/1.00 0.00/0.00 1.00 / 1.00 1.00 / 1.00 1.00 / 1.00 1.00 / 1.00 0.88 / 0.88 0.75 / 0.75 0.50 / 0.50 |
| TIES SVD 10 diagonal (D) 10 full (F) | 10 50 100 500 SVD 2 SVD 4 SVD 8 SVD 8 SVD 16 64 D 128 D 256 D 128 D 256 C 16 F 32 F 64 F 64 F 128 F | | $\begin{array}{c} 0.57 \pm 0.02 \\ 0.41 \pm 0.00 \\ 0.40 \pm 0.00 \\ 0.26 \pm 0.00 \\ 1.07 \pm 0.02 \\ 1.04 \pm 0.01 \\ 1.02 \pm 0.01 \\ 1.00 \pm 0.00 \\ \hline \end{array}$ | $\begin{array}{c} 0.45 \pm 0.04 \\ 0.18 \pm 0.05 \\ 0.20 \pm 0.05 \\ 0.01 \pm 0.00 \\ 1.00 \pm 0.$ | $\begin{array}{c} 0.10 \pm 0.01 \\ 0.03 \pm 0.01 \\ 0.01 \pm 0.02 \\ 0.00 \pm 0.00 \\ 1.00 \pm 0.01 \\ 0.99 \pm 0.00 \\ 0.01 \\ 0.01 \pm 0.01 \\ 0.00 \\ 0.01 \\ 0.01 \\ 0.01 \\ 0.00 \\ 0.01 \\ 0.00 \\ 0.01 \\ 0.00 \\ 0.01 \\ 0.00 \\ 0.$ | $\left \begin{array}{c} 0.83 \pm 0.01 \\ 0.91 \pm 0.01 \\ 0.88 \pm 0.00 \\ 0.83 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.01 \\ 1.00 \pm 0.00 \\ 1.00 \pm $ | $\begin{array}{c} 0.47 \pm 0.00 \\ 0.31 \pm 0.00 \\ 0.33 \pm 0.00 \\ 0.29 \pm 0.00 \\ \hline 0.98 \pm 0.01 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ \hline \end{array}$ | $\begin{array}{c} 0.69 \pm 0.01 \\ 0.65 \pm 0.00 \\ 0.64 \pm 0.00 \\ 0.57 \pm 0.00 \\ \end{array}$ | $\begin{array}{c} 0.57 \pm 0.00 \\ 0.62 \pm 0.00 \\ 0.57 \pm 0.02 \\ 0.37 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.82 \pm 0.01 \\ 0.32 \pm 0.04 \\ 0.01 \pm 0.00 \\ 0.01 \pm 0.00 \\ \hline 1.00 \pm 0.01 \\ 0.99 \pm 0.01 \\ 1.01 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.85 \pm 0.00 \\ 0.84 \pm 0.00 \\ 0.82 \pm 0.00 \\ 0.43 \pm 0.00 \\ \hline 1.00 \pm 0.01 \\ 1.00 \pm 0.01 \\ 1.01 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.62 \pm 0.23 \\ 0.48 \pm 0.28 \\ 0.44 \pm 0.30 \\ 0.31 \pm 0.26 \end{array}$ $\begin{array}{c} 1.00 \pm 0.04 \\ 1.00 \pm 0.03 \\ 1.00 \pm 0.01 \end{array}$ | 1.00 / 1.00 1.00 / 1.00 1.00 / 1.00 1.00 / 1.00 0.88 / 0.88 0.75 / 0.75 0.50 / 0.50 |
| SVD 10 diagonal (D) 10 full (F) | SVD 2 SVD 4 SVD 8 SVD 16 16 D 32 D 64 D 256 D 16 F 32 F 64 F 128 F | $ \begin{array}{ c c c c c c c c c c c c c c c c c c c$ | $\begin{array}{c} 1.07 \pm 0.02 \\ 1.04 \pm 0.01 \\ 1.02 \pm 0.01 \\ 1.00 \pm 0.00 \\ \hline \\ 1.01 \pm 0.01 \\ 1.03 \pm 0.01 \\ 1.03 \pm 0.01 \\ 1.01 \pm 0.01 \\ 1.00 \pm 0.00 \\ \hline \end{array}$ | $\begin{array}{c} 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ \hline 1.00 \pm 0.01 \\ 0.99 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.00 \pm 0.00 \\ 1.00 \pm 0.01 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.98 \pm 0.01 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.00 \pm 0.01 \\ 0.99 \pm 0.02 \\ 1.01 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.00 \pm 0.10 \\ 0.99 \pm 0.08 \\ 1.00 \pm 0.01 \end{array}$ | $\begin{array}{c} 1.00 \pm 0.01 \\ 0.99 \pm 0.01 \\ 1.01 \pm 0.01 \end{array}$ | $\begin{array}{c} 1.00 \pm 0.01 \\ 1.00 \pm 0.01 \\ 1.01 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.00 \pm 0.04 \\ 1.00 \pm 0.03 \\ 1.00 \pm 0.01 \end{array}$ | 0.88 / 0.88 0.75 / 0.75 0.50 / 0.50 |
| 10 diagonal (D) | 16 D 32 D 64 D 128 D 256 D 16 F 32 F 64 F 128 F | $\begin{array}{c} 1.02 \pm 0.01 \\ 1.01 \pm 0.01 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.02 \pm 0.00 \\ \end{array}$ | $\begin{array}{c} 1.01 \pm 0.01 \\ 1.05 \pm 0.01 \\ 1.03 \pm 0.01 \\ 1.01 \pm 0.01 \\ 1.00 \pm 0.00 \end{array}$ | 1.00 ± 0.00 1.00 ± 0.00 1.00 ± 0.00 1.00 ± 0.00 | $\begin{array}{c} 1.00 \pm 0.01 \\ 0.99 \pm 0.00 \end{array}$ | 0.99 ± 0.00 | | 1.00 ± 0.00 | 1.00 ± 0.00 | 1.00 ± 0.00 | 1.00 ± 0.00 | 1.00 ± 0.00 | 0.00 / 0.00 |
| 10 full (F) | 16 F 32 F 64 F 128 F | 1.02 ± 0.00 | | 1.00 ± 0.00 1.00 ± 0.00 | 1.00 ± 0.00 1.00 ± 0.00 1.00 ± 0.00 | $\begin{array}{c} 1.01 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.96 \pm 0.00 \\ 0.99 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.02 \pm 0.02 \\ 0.97 \pm 0.01 \\ 1.01 \pm 0.01 \\ 1.01 \pm 0.01 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.13 \pm 0.03 \\ 1.05 \pm 0.03 \\ 0.99 \pm 0.01 \\ 0.99 \pm 0.01 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.99 \pm 0.02 \\ 1.00 \pm 0.01 \\ 1.01 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.98 \pm 0.01 \\ 1.00 \pm 0.01 \\ 1.01 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.01 \pm 0.05 \\ 1.00 \pm 0.03 \\ 1.00 \pm 0.01 \\ 1.00 \pm 0.01 \\ 1.00 \pm 0.00 \end{array}$ | 1.00 / 0.90 1.00 / 0.80 1.00 / 0.60 1.00 / 0.20 1.00 / -0.60 |
| | 256 F | $\begin{array}{c} 1.02 \pm 0.01 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.06 \pm 0.01 \\ 1.04 \pm 0.01 \\ 1.03 \pm 0.01 \\ 1.01 \pm 0.01 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.99 \pm 0.01 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.98 \pm 0.00 \\ 0.99 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.01 \pm 0.02 \\ 0.96 \pm 0.01 \\ 1.01 \pm 0.01 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.07 \pm 0.00 \\ 1.00 \pm 0.02 \\ 0.98 \pm 0.01 \\ 0.99 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.01 \pm 0.01 \\ 1.00 \pm 0.01 \\ 1.01 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.00 \pm 0.00 \\ 1.01 \pm 0.00 \\ 1.01 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.01 \pm 0.03 \\ 1.00 \pm 0.02 \\ 1.00 \pm 0.01 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | 1.00/0.90 0.99/0.79 0.97/0.57 0.88/0.07 0.50/-1.10 |
| 50 diagonal (D) | 16 D 32 D 64 D 128 D 256 D | $\begin{array}{c} 0.98 \pm 0.04 \\ 1.00 \pm 0.02 \\ 1.02 \pm 0.00 \\ 1.01 \pm 0.01 \\ 1.01 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.98 \pm 0.01 \\ 1.02 \pm 0.02 \\ 1.05 \pm 0.02 \\ 1.08 \pm 0.01 \\ 1.03 \pm 0.01 \end{array}$ | $\begin{array}{c} 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.92 \pm 0.06 \\ 0.99 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.84 \pm 0.07 \\ 0.96 \pm 0.01 \\ 0.99 \pm 0.01 \\ 0.99 \pm 0.01 \\ 1.00 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.92 \pm 0.02 \\ 0.95 \pm 0.02 \\ 0.97 \pm 0.00 \\ 0.98 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.68 \pm 0.05 \\ 0.84 \pm 0.02 \\ 0.99 \pm 0.01 \\ 0.98 \pm 0.01 \\ 0.97 \pm 0.03 \end{array}$ | $\begin{array}{c} 0.87 \pm 0.10 \\ 1.00 \pm 0.13 \\ 1.09 \pm 0.03 \\ 1.11 \pm 0.03 \\ 1.01 \pm 0.03 \end{array}$ | $\begin{array}{c} 0.88 \pm 0.07 \\ 0.98 \pm 0.01 \\ 1.01 \pm 0.01 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.83 \pm 0.02 \\ 0.88 \pm 0.01 \\ 0.90 \pm 0.01 \\ 1.00 \pm 0.01 \\ 1.01 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.89 \pm 0.10 \\ 0.96 \pm 0.07 \\ 1.00 \pm 0.05 \\ 1.01 \pm 0.04 \\ 1.00 \pm 0.02 \end{array}$ | 1.00 / 0.98 1.00 / 0.96 1.00 / 0.92 1.00 / 0.84 1.00 / 0.68 |
| 50 full (F) | 16 F 32 F 64 F 128 F 256 F | $\begin{array}{c} 0.99 \pm 0.04 \\ 1.02 \pm 0.00 \\ 1.02 \pm 0.01 \\ 1.02 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.00 \pm 0.01 \\ 1.00 \pm 0.02 \\ 1.06 \pm 0.02 \\ 1.06 \pm 0.01 \\ 1.02 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.00 \pm 0.01 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.96 \pm 0.01 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.95 \pm 0.02 \\ 0.98 \pm 0.01 \\ 0.99 \pm 0.01 \\ 1.00 \pm 0.01 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.94 \pm 0.01 \\ 0.96 \pm 0.00 \\ 0.98 \pm 0.01 \\ 0.98 \pm 0.00 \\ 0.99 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.64 \pm 0.10 \\ 0.95 \pm 0.01 \\ 1.03 \pm 0.01 \\ 0.98 \pm 0.01 \\ 1.01 \pm 0.01 \end{array}$ | $\begin{array}{c} 1.01 \pm 0.15 \\ 1.09 \pm 0.02 \\ 1.11 \pm 0.00 \\ 1.03 \pm 0.04 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.97 \pm 0.02 \\ 1.01 \pm 0.02 \\ 1.00 \pm 0.01 \\ 1.00 \pm 0.01 \\ 1.01 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.87 \pm 0.00 \\ 0.89 \pm 0.01 \\ 0.98 \pm 0.02 \\ 1.00 \pm 0.00 \\ 1.01 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.93 \pm 0.12 \\ 0.99 \pm 0.05 \\ 1.02 \pm 0.04 \\ 1.01 \pm 0.03 \\ 1.00 \pm 0.01 \end{array}$ | 1.00/0.98 0.99/0.95 0.97/0.89 0.88/0.72 0.50/0.18 |
| 100 diagonal (D) | 16 D 32 D 64 D 128 D 256 D | $\begin{array}{c} 0.80 \pm 0.07 \\ 0.95 \pm 0.06 \\ 1.01 \pm 0.03 \\ 1.01 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.89 \pm 0.06 \\ 0.98 \pm 0.01 \\ 1.01 \pm 0.01 \\ 1.02 \pm 0.01 \\ 1.06 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.93 \pm 0.03 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.96 \pm 0.01 \\ 0.91 \pm 0.06 \\ 0.98 \pm 0.02 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.50 \pm 0.09 \\ 0.80 \pm 0.14 \\ 0.96 \pm 0.01 \\ 0.99 \pm 0.01 \\ 0.99 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.78 \pm 0.01 \\ 0.89 \pm 0.06 \\ 0.94 \pm 0.01 \\ 0.97 \pm 0.00 \\ 0.98 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.28 \pm 0.07 \\ 0.60 \pm 0.10 \\ 0.88 \pm 0.05 \\ 1.00 \pm 0.03 \\ 1.00 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.52 \pm 0.10 \\ 0.77 \pm 0.26 \\ 1.11 \pm 0.08 \\ 1.11 \pm 0.02 \\ 1.11 \pm 0.03 \end{array}$ | $\begin{array}{c} 0.78 \pm 0.03 \\ 0.91 \pm 0.02 \\ 0.96 \pm 0.02 \\ 0.99 \pm 0.01 \\ 1.00 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.81 \pm 0.02 \\ 0.83 \pm 0.02 \\ 0.87 \pm 0.03 \\ 0.89 \pm 0.02 \\ 0.98 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.72 \pm 0.22 \\ 0.86 \pm 0.14 \\ 0.97 \pm 0.07 \\ 1.00 \pm 0.05 \\ 1.01 \pm 0.04 \end{array}$ | 1.00 / 0.99 1.00 / 0.98 1.00 / 0.96 1.00 / 0.92 1.00 / 0.84 |
| 100 full (F) | 16 F 32 F 64 F 128 F 256 F | $\begin{array}{c} 0.95 \pm 0.01 \\ 1.00 \pm 0.02 \\ 1.02 \pm 0.00 \\ 1.01 \pm 0.01 \\ 1.01 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.97 \pm 0.03 \\ 0.99 \pm 0.01 \\ 1.00 \pm 0.02 \\ 1.05 \pm 0.01 \\ 1.03 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.97 \pm 0.03 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.97 \pm 0.03 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \\ 0.99 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.93 \pm 0.01 \\ 0.97 \pm 0.01 \\ 0.98 \pm 0.00 \\ 1.00 \pm 0.00 \\ 1.01 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.92 \pm 0.01 \\ 0.95 \pm 0.00 \\ 0.96 \pm 0.00 \\ 0.98 \pm 0.00 \\ 0.99 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.64 \pm 0.03 \\ 0.86 \pm 0.03 \\ 0.99 \pm 0.01 \\ 1.03 \pm 0.01 \\ 0.98 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.89 \pm 0.16 \\ 1.12 \pm 0.03 \\ 1.09 \pm 0.01 \\ 1.10 \pm 0.01 \\ 1.00 \pm 0.03 \end{array}$ | $\begin{array}{c} 0.87 \pm 0.02 \\ 0.96 \pm 0.01 \\ 0.99 \pm 0.02 \\ 1.01 \pm 0.00 \\ 1.01 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.83 \pm 0.01 \\ 0.87 \pm 0.00 \\ 0.89 \pm 0.00 \\ 0.99 \pm 0.01 \\ 1.01 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.89 \pm 0.11 \\ 0.97 \pm 0.07 \\ 0.99 \pm 0.05 \\ 1.02 \pm 0.04 \\ 1.00 \pm 0.01 \end{array}$ | 1.00/0.99 0.99/0.97 0.97/0.93 0.88/0.80 0.50/0.34 |
| 100 w/clusters (C) | 16 C 5 16 C 7 | 1.12 1.12 | 1.02 1.02 | 1.00 1.00 | 1.00 1.00 | 0.98 1.00 | 0.96 0.97 | 1.00 1.01 | 1.20 1.30 | 1.04 1.03 | 0.90 0.92 | 1.02 1.04 | 1.00/0.95 1.00/0.93 |
| 500 diagonal (D) | 16 D 32 D 64 D 128 D 256 D | $\begin{array}{c} 0.57 \pm 0.07 \\ 0.61 \pm 0.12 \\ 0.73 \pm 0.02 \\ 0.84 \pm 0.00 \\ 0.99 \pm 0.03 \end{array}$ | $\begin{array}{c} 0.55 \pm 0.03 \\ 0.55 \pm 0.08 \\ 0.63 \pm 0.11 \\ 0.92 \pm 0.02 \\ 0.99 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.83 \pm 0.04 \\ 0.83 \pm 0.02 \\ 0.89 \pm 0.04 \\ 0.97 \pm 0.03 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.78 \pm 0.16 \\ 0.84 \pm 0.12 \\ 0.97 \pm 0.00 \\ 0.98 \pm 0.01 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.85 \pm 0.04 \\ 0.91 \pm 0.02 \\ 0.94 \pm 0.00 \\ 0.94 \pm 0.00 \\ 0.96 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.68 \pm 0.07 \\ 0.71 \pm 0.05 \\ 0.83 \pm 0.05 \\ 0.88 \pm 0.02 \\ 0.92 \pm 0.03 \end{array}$ | $\begin{array}{c} 0.24 \pm 0.01 \\ 0.29 \pm 0.05 \\ 0.45 \pm 0.09 \\ 0.60 \pm 0.15 \\ 0.66 \pm 0.06 \end{array}$ | $\begin{array}{c} 0.43 \pm 0.01 \\ 0.47 \pm 0.08 \\ 0.50 \pm 0.07 \\ 0.53 \pm 0.01 \\ 0.84 \pm 0.14 \end{array}$ | $\begin{array}{c} 0.76 \pm 0.06 \\ 0.79 \pm 0.04 \\ 0.82 \pm 0.02 \\ 0.85 \pm 0.05 \\ 0.92 \pm 0.02 \end{array}$ | $\begin{array}{c} 0.79 \pm 0.01 \\ 0.79 \pm 0.01 \\ 0.80 \pm 0.02 \\ 0.80 \pm 0.02 \\ 0.84 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.65 \pm 0.20 \\ 0.68 \pm 0.20 \\ 0.76 \pm 0.18 \\ 0.83 \pm 0.15 \\ 0.91 \pm 0.11 \end{array}$ | 1.00 / 1.00 1.00 / 1.00 1.00 / 0.99 1.00 / 0.98 1.00 / 0.97 |
| 500 full (F) | 16 F 32 F 64 F 128 F 256 F | $\begin{array}{c} 0.57 \pm 0.01 \\ 0.79 \pm 0.05 \\ 1.02 \pm 0.00 \\ 1.03 \pm 0.01 \\ 1.03 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.43 \pm 0.07 \\ 0.54 \pm 0.04 \\ 0.96 \pm 0.01 \\ 0.97 \pm 0.02 \\ 1.03 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.78 \pm 0.01 \\ 0.93 \pm 0.02 \\ 0.94 \pm 0.01 \\ 0.99 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.97 \pm 0.00 \\ 0.98 \pm 0.00 \\ 1.00 \pm 0.01 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.96 \pm 0.00 \\ 0.97 \pm 0.00 \\ 0.96 \pm 0.00 \\ 0.98 \pm 0.00 \\ 0.99 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.83 \pm 0.01 \\ 0.90 \pm 0.01 \\ 0.97 \pm 0.01 \\ 0.96 \pm 0.00 \\ 0.97 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.64 \pm 0.00 \\ 0.69 \pm 0.01 \\ 0.73 \pm 0.01 \\ 0.87 \pm 0.01 \\ 0.99 \pm 0.02 \end{array}$ | $\begin{array}{c} 0.53 \pm 0.03 \\ 0.50 \pm 0.00 \\ 0.54 \pm 0.01 \\ 1.07 \pm 0.02 \\ 1.03 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.83 \pm 0.01 \\ 0.86 \pm 0.02 \\ 0.91 \pm 0.01 \\ 0.98 \pm 0.00 \\ 1.00 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.83 \pm 0.00 \\ 0.83 \pm 0.01 \\ 0.86 \pm 0.00 \\ 0.87 \pm 0.00 \\ 0.87 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.75 \pm 0.17 \\ 0.81 \pm 0.16 \\ 0.89 \pm 0.14 \\ 0.97 \pm 0.06 \\ 0.99 \pm 0.05 \end{array}$ | 1.00/1.00 0.99/0.99 0.97/0.96 0.88/0.86 0.50/0.47 |
| 500 w/clusters (C) | 16 C 7 16 C 10 16 C 25 64 C 5 64 C 7 | 1.09 1.10 1.10 1.09 1.12 | 1.00 1.01 1.00 0.98 1.02 | 0.99 1.00 1.00 1.00 1.00 | 1.00 0.99 0.99 1.00 1.00 | 0.98 0.97 0.99 0.99 1.00 | 0.95 0.93 0.96 0.96 0.96 | 0.72 0.70 0.98 0.99 0.99 | 0.87 1.30 1.31 1.18 1.22 | 0.98 1.02 1.03 1.04 1.04 | 0.90 0.88 0.91 0.87 0.93 | 0.95 0.99 1.03 1.01 1.03 | 1.00/0.98 1.00/0.98 1.00/0.95 0.97/0.93 0.97/0.91 |

Table 4: Relative In-Distribution ROUGE-L scores for various tasks and methods

G.7 AGREEMENT AND COMPRESSION RATE

Table 10 provides full results for *agreement*, which shows the same trends as the results for relative performance of Rouge-L (Table 4). Note that *agreement* measures the exact match in task generations between the uncompressed LoRA model and the compressed LoRA model, rather than comparing to the task's ground truth data. The comparison is very strict and requires an exact match between the generations of the two models (LoRA and the compressed LoRA), comparing each sample one at a time.

1231

1220

1221 1222 1223

1224

1232 G.8 RECONSTRUCTION ERROR AND COMPRESSION RATE

Table 11 provides the full results of the experiments behind Figure 3 for every evaluation task.

1236

1237 G.9 RECONSTRUCTION ERROR: TRAINED VS. RANDOM

Table 12 provides the reconstruction error on random (untrained) LoRA matrices. Comparing with
Table 11, we find that reconstruction error is consistently higher on random (untrained LoRA) matrices than on trained LoRA matrices. This demonstrates that after training, LoRAs have a shared structure that JD exploits.

| 2 | Madal Trans | Mathead Trans | 1 | | | | Tas | ks | | | | | A | Deres Cound |
|---|--------------------|--|--|--|--|---|--|---|--|--|--|--|--|---|
| | Model Type | Method Type | task039 | task190 | task280 | task290 | task391 | task442 | task620 | task1342 | task1391 | task1598 | Average | Para. Saved |
| | | base lora | $\begin{array}{c} 24.44 \pm 0.00 \\ 95.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.60 \pm 0.00 \\ 86.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 19.13 \pm 0.00 \\ 99.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 39.22 \pm 0.00 \\ 93.67 \pm 0.00 \end{array}$ | $\begin{array}{c} 10.27 \pm 0.00 \\ 94.33 \pm 0.00 \end{array}$ | $\begin{array}{c} 35.46 \pm 0.00 \\ 74.88 \pm 0.00 \end{array}$ | $\begin{array}{c} 7.85 \pm 0.00 \\ 74.40 \pm 0.00 \end{array}$ | $\begin{array}{c} 6.22 \pm 0.00 \\ 26.68 \pm 0.00 \end{array}$ | $\begin{array}{c} 17.82 \pm 0.00 \\ 95.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 38.87 \pm 0.00 \\ 50.32 \pm 0.00 \end{array}$ | $\begin{array}{c} 20.24 \pm {\scriptstyle 13.27} \\ 78.87 \pm {\scriptstyle 22.56} \end{array}$ | 1.00/1.00 0.00/0.00 |
| | TIES | 10 50 100 500 | $\begin{array}{c} 76.50 \pm 0.00 \\ 55.80 \pm 0.00 \\ 52.43 \pm 0.00 \\ 35.18 \pm 0.00 \end{array}$ | $\begin{array}{c} 49.00 \pm 1.73 \\ 35.00 \pm 0.00 \\ 34.00 \pm 0.00 \\ 22.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 44.33 \pm 4.04 \\ 18.00 \pm 5.20 \\ 19.67 \pm 4.62 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 9.80 \pm 0.58 \\ 2.42 \pm 0.50 \\ 1.09 \pm 1.66 \\ 0.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 78.56 \pm 0.96 \\ 85.78 \pm 0.96 \\ 83.33 \pm 0.00 \\ 78.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 35.24 \pm 0.00 \\ 23.03 \pm 0.00 \\ 24.89 \pm 0.00 \\ 21.46 \pm 0.00 \end{array}$ | $\begin{array}{c} 51.37 \pm 0.67 \\ 48.03 \pm 0.00 \\ 47.52 \pm 0.00 \\ 42.22 \pm 0.04 \end{array}$ | $\begin{array}{c} 15.26 \pm 0.12 \\ 16.50 \pm 0.00 \\ 15.18 \pm 0.42 \\ 9.93 \pm 0.13 \end{array}$ | $\begin{array}{c} 77.67 \pm 1.15 \\ 30.00 \pm 3.46 \\ 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 42.72 \pm 0.01 \\ 42.47 \pm 0.02 \\ 41.19 \pm 0.03 \\ 21.50 \pm 0.03 \end{array}$ | $\begin{array}{c} 48.05 \pm 23.61 \\ 35.70 \pm 23.01 \\ 32.03 \pm 24.50 \\ 23.27 \pm 23.64 \end{array}$ | 1.00 / 1.00 1.00 / 1.00 1.00 / 1.00 1.00 / 1.00 |
| | SVD | SVD 2 SVD 4 SVD 8 SVD 16 | $\begin{array}{c} 93.15 \pm 2.77 \\ 94.01 \pm 3.60 \\ 95.00 \pm 0.00 \\ 95.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 92.24 \pm 1.85 \\ 89.21 \pm 0.71 \\ 87.40 \pm 0.59 \\ 86.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 99.09 \pm 0.18 \\ 99.05 \pm 0.09 \\ 99.05 \pm 0.09 \\ 99.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 93.44 \pm 0.14 \\ 93.65 \pm 0.03 \\ 93.65 \pm 0.03 \\ 93.67 \pm 0.00 \end{array}$ | $\begin{array}{c} 93.89 \pm 0.35 \\ 94.66 \pm 0.63 \\ 94.36 \pm 0.38 \\ 94.33 \pm 0.00 \end{array}$ | $\begin{array}{c} 73.74 \pm 0.51 \\ 74.89 \pm 0.33 \\ 74.58 \pm 0.12 \\ 74.90 \pm 0.03 \end{array}$ | $\begin{array}{c} 74.55 \pm 0.98 \\ 73.61 \pm 1.15 \\ 75.07 \pm 0.00 \\ 74.23 \pm 0.18 \end{array}$ | $\begin{array}{ c c c c c c c c c c c c c c c c c c c$ | $\begin{array}{c} 95.06 \pm 1.35 \\ 93.98 \pm 0.77 \\ 95.51 \pm 1.09 \\ 95.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 50.21 \pm 0.44 \\ 50.47 \pm 0.54 \\ 50.89 \pm 0.07 \\ 50.30 \pm 0.02 \end{array}$ | $\begin{array}{c} 79.11 \pm 22.72 \\ 78.90 \pm 22.68 \\ 81.01 \pm 21.74 \\ 78.36 \pm 22.97 \end{array}$ | 0.88 / 0.88 0.75 / 0.75 0.50 / 0.50 0.00 / 0.00 |
| | 10 diagonal (D) | 16 D 32 D 64 D 128 D 256 D | $\begin{array}{c} 96.67 \pm 0.58 \\ 95.67 \pm 0.58 \\ 95.00 \pm 0.00 \\ 95.00 \pm 0.00 \\ 95.00 \pm 0.00 \\ \end{array}$ | $ \begin{array}{ c c c c c c c c c c c c c c c c c c c$ | $\begin{array}{c} 99.00 \pm 0.00 \\ \end{array}$ | $\begin{array}{c} 94.00 \pm 0.67 \\ 93.00 \pm 0.33 \\ 93.67 \pm 0.00 \\ 93.67 \pm 0.00 \\ 03.67 \pm 0.00 \end{array}$ | $\begin{array}{c} 93.11 \pm 0.38 \\ 94.89 \pm 0.51 \\ 94.78 \pm 0.38 \\ 94.33 \pm 0.00 \\ 94.33 \pm 0.00 \end{array}$ | $\begin{array}{c} 72.08 \pm 0.06 \\ 73.86 \pm 0.31 \\ 74.61 \pm 0.13 \\ 74.92 \pm 0.13 \\ 74.88 \pm 0.00 \end{array}$ | 76.26 ± 1.19 71.92 ± 0.84 74.97 ± 0.58 74.96 ± 0.51 74.40 ± 0.99 | $\begin{array}{c} 30.11 \pm 0.79 \\ 27.89 \pm 0.70 \\ 26.35 \pm 0.25 \\ 26.45 \pm 0.23 \\ 26.68 \pm 0.00 \\ \end{array}$ | $\begin{array}{c} 94.00 \pm 1.73 \\ 94.67 \pm 0.58 \\ 96.00 \pm 0.00 \\ 95.00 \pm 0.00 \\ 95.00 \pm 0.00 \\ \end{array}$ | $\begin{array}{c} 49.30 \pm 0.46 \\ 50.36 \pm 0.26 \\ 50.99 \pm 0.06 \\ 50.21 \pm 0.12 \\ 50.27 \pm 0.02 \end{array}$ | $\begin{array}{c} 79.15 \pm 22.18 \\ 79.13 \pm 22.75 \\ 79.37 \pm 22.94 \\ 79.02 \pm 22.84 \\ 78.02 \pm 22.84 \end{array}$ | 1.00 / 0.90 1.00 / 0.80 1.00 / 0.60 1.00 / 0.20 |
| | 10 full (F) | 16 F 32 F 64 F 128 F 256 F | 97.00 ± 0.00 96.67 ± 0.58 95.00 ± 0.00 95.00 ± 0.00 95.00 ± 0.00 | $\begin{array}{c} 91.00 \pm 1.00 \\ 89.33 \pm 0.58 \\ 88.67 \pm 0.58 \\ 86.67 \pm 0.58 \\ 86.00 \pm 0.00 \\ \end{array}$ | $\begin{array}{c} 99.00 \pm 0.00 \\ \end{array}$ | $\begin{array}{c} 93.56 \pm 0.19 \\ 93.22 \pm 0.19 \\ 93.67 \pm 0.00 \\ 93.67 \pm 0.00 \\ 93.67 \pm 0.00 \\ \end{array}$ | $\begin{array}{c} 93.56 \pm 0.69 \\ 94.44 \pm 0.19 \\ 94.56 \pm 0.38 \\ 94.33 \pm 0.00 \\ 94.33 \pm 0.00 \end{array}$ | 73.60 ± 0.36 74.11 ± 0.19 74.56 ± 0.13 75.04 ± 0.03 74.90 ± 0.03 | 74.94 ± 1.25 71.74 ± 0.59 75.47 ± 0.58 74.40 ± 0.00 74.29 ± 0.19 | $\begin{array}{c} 28.66 \pm 0.03 \\ 26.74 \pm 0.50 \\ 26.26 \pm 0.34 \\ 26.53 \pm 0.13 \\ 26.68 \pm 0.00 \end{array}$ | 96.00 ± 1.00 94.67 ± 0.58 96.00 ± 0.00 95.00 ± 0.00 95.00 ± 0.00 | $\begin{array}{c} 50.15 \pm 0.20 \\ 50.63 \pm 0.24 \\ 50.89 \pm 0.17 \\ 50.36 \pm 0.03 \\ 50.30 \pm 0.03 \\ \end{array}$ | $\begin{array}{c} 79.75 \pm 22.72 \\ 79.06 \pm 23.01 \\ 79.41 \pm 22.97 \\ 79.00 \pm 22.81 \\ 78.92 \pm 22.77 \end{array}$ | 1.00/0.90 0.99/0.79 0.97/0.57 0.88/0.07 0.50/-1.10 |
| | 50 diagonal (D) | 16 D 32 D 64 D 128 D 256 D | $\begin{array}{c} 92.76 \pm 3.53 \\ 95.33 \pm 2.08 \\ 97.00 \pm 0.00 \\ 96.33 \pm 0.58 \\ 95.67 \pm 0.58 \end{array}$ | $\begin{array}{ c c c c c c c c c c c c c c c c c c c$ | $\begin{array}{c} 99.00 \pm 0.00 \\ 99.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 86.17 \pm 5.81 \\ 92.60 \pm 0.29 \\ 93.78 \pm 0.19 \\ 93.56 \pm 0.19 \\ 93.56 \pm 0.19 \\ \end{array}$ | $\begin{array}{c} 79.68 \pm 6.21 \\ 90.32 \pm 1.04 \\ 93.00 \pm 0.58 \\ 93.00 \pm 0.58 \\ 94.67 \pm 0.67 \end{array}$ | $\begin{array}{c} 69.07 \pm 1.54 \\ 71.16 \pm 1.47 \\ 72.37 \pm 0.35 \\ 73.32 \pm 0.24 \\ 74.82 \pm 0.24 \end{array}$ | $\begin{array}{c} 50.65 \pm 3.97 \\ 62.51 \pm 1.64 \\ 73.39 \pm 0.93 \\ 73.03 \pm 1.09 \\ 72.36 \pm 2.07 \end{array}$ | $\begin{array}{ c c c c c c c c c c c c c c c c c c c$ | $\begin{array}{c} 83.90 \pm 6.43 \\ 93.33 \pm 1.15 \\ 95.67 \pm 0.58 \\ 95.00 \pm 0.00 \\ 95.33 \pm 0.58 \end{array}$ | $\begin{array}{c} 41.86 \pm 0.96 \\ 44.35 \pm 0.41 \\ 45.43 \pm 0.34 \\ 50.16 \pm 0.74 \\ 50.73 \pm 0.46 \end{array}$ | $\begin{array}{c} 71.10 \pm 23.99 \\ 76.25 \pm 23.81 \\ 78.90 \pm 23.29 \\ 79.56 \pm 22.51 \\ 79.14 \pm 22.90 \end{array}$ | 1.00 / 0.98 1.00 / 0.96 1.00 / 0.92 1.00 / 0.84 1.00 / 0.68 |
| | 50 full (F) | 16 F 32 F 64 F 128 F 256 F | $\begin{array}{c} 94.06 \pm 3.54 \\ 97.00 \pm 0.00 \\ 96.67 \pm 0.58 \\ 97.00 \pm 0.00 \\ 95.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 85.67 \pm 1.15 \\ 85.67 \pm 1.53 \\ 91.00 \pm 2.00 \\ 91.00 \pm 1.00 \\ 88.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 98.67 \pm 0.58 \\ 99.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 90.35 \pm 1.37 \\ 93.67 \pm 0.00 \\ 93.56 \pm 0.19 \\ 93.33 \pm 0.00 \\ 93.67 \pm 0.00 \end{array}$ | $\begin{array}{c} 89.90 \pm 1.91 \\ 92.22 \pm 0.69 \\ 93.22 \pm 0.51 \\ 94.11 \pm 0.51 \\ 94.44 \pm 0.19 \end{array}$ | $\begin{array}{c} 70.32 \pm 0.66 \\ 71.88 \pm 0.30 \\ 73.16 \pm 0.41 \\ 73.51 \pm 0.23 \\ 74.25 \pm 0.21 \end{array}$ | $\begin{array}{c} 47.62 \pm 7.28 \\ 71.01 \pm 1.02 \\ 76.28 \pm 0.51 \\ 73.17 \pm 0.58 \\ 74.97 \pm 0.58 \end{array}$ | $\begin{array}{c} 26.88 \pm 3.96 \\ 29.07 \pm 0.65 \\ 29.67 \pm 0.12 \\ 27.53 \pm 1.12 \\ 26.79 \pm 0.09 \end{array}$ | $\begin{array}{c} 92.33 \pm 1.53 \\ 95.67 \pm 1.53 \\ 95.33 \pm 0.58 \\ 95.00 \pm 1.00 \\ 96.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 43.68 \pm 0.24 \\ 44.97 \pm 0.41 \\ 49.31 \pm 1.00 \\ 50.56 \pm 0.06 \\ 50.86 \pm 0.19 \end{array}$ | $\begin{array}{c} 73.95 \pm 24.73 \\ 78.02 \pm 23.18 \\ 79.72 \pm 22.50 \\ 79.42 \pm 22.93 \\ 79.30 \pm 22.82 \end{array}$ | 1.00/0.98 0.99/0.95 0.97/0.89 0.88/0.72 0.50/0.18 |
| | 100 diagonal (D) | 16 D 32 D 64 D 128 D 256 D | $\begin{array}{c} 76.43 \pm 7.07 \\ 90.10 \pm 5.85 \\ 95.56 \pm 2.49 \\ 96.00 \pm 0.00 \\ 95.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 76.67 \pm 4.93 \\ 84.00 \pm 1.00 \\ 86.67 \pm 0.58 \\ 87.33 \pm 1.15 \\ 91.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 91.61 \pm 2.75 \\ 99.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 89.99 \pm 1.07 \\ 85.52 \pm 5.34 \\ 92.24 \pm 1.68 \\ 93.89 \pm 0.19 \\ 93.56 \pm 0.19 \end{array}$ | $\begin{array}{c} 47.55 \pm 8.56 \\ 75.69 \pm 12.75 \\ 90.89 \pm 1.17 \\ 93.00 \pm 0.58 \\ 93.11 \pm 0.19 \end{array}$ | $\begin{array}{c} 58.08 \pm 0.72 \\ 66.62 \pm 4.18 \\ 70.35 \pm 0.45 \\ 72.70 \pm 0.30 \\ 73.05 \pm 0.20 \end{array}$ | $\begin{array}{c} 20.77 \pm 5.50 \\ 44.66 \pm 7.26 \\ 65.62 \pm 4.03 \\ 74.34 \pm 2.07 \\ 74.52 \pm 0.95 \end{array}$ | $\begin{array}{c} 13.90 \pm 2.79 \\ 20.49 \pm 7.07 \\ 29.58 \pm 2.02 \\ 29.66 \pm 0.54 \\ 29.67 \pm 0.67 \end{array}$ | $\begin{array}{c} 73.93 \pm 3.13 \\ 86.67 \pm 1.86 \\ 91.67 \pm 2.31 \\ 93.67 \pm 0.58 \\ 95.33 \pm 0.58 \end{array}$ | $\begin{array}{c} 40.74 \pm 0.85 \\ 42.01 \pm 0.94 \\ 43.64 \pm 1.36 \\ 44.82 \pm 0.89 \\ 49.42 \pm 0.65 \end{array}$ | $\begin{array}{c} 58.97 \pm 26.83 \\ 69.48 \pm 25.14 \\ 76.52 \pm 23.02 \\ 78.44 \pm 22.87 \\ 79.37 \pm 22.38 \end{array}$ | 1.00 / 0.99 1.00 / 0.98 1.00 / 0.96 1.00 / 0.92 1.00 / 0.84 |
| | 100 full (F) | 16 F 32 F 64 F 128 F 256 F | $\begin{array}{c} 90.70 \pm 1.07 \\ 95.33 \pm 1.53 \\ 97.00 \pm 0.00 \\ 96.33 \pm 0.58 \\ 96.33 \pm 0.58 \end{array}$ | $\begin{array}{c} 83.00 \pm 2.65 \\ 85.00 \pm 1.00 \\ 85.67 \pm 1.53 \\ 90.33 \pm 0.58 \\ 88.67 \pm 0.58 \end{array}$ | $\begin{array}{c} 96.00 \pm 3.00 \\ 99.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 91.22 \pm 2.94 \\ 93.50 \pm 0.22 \\ 93.78 \pm 0.19 \\ 93.00 \pm 0.00 \\ 93.67 \pm 0.00 \end{array}$ | $\begin{array}{c} 87.94 \pm 0.54 \\ 91.44 \pm 0.84 \\ 92.56 \pm 0.19 \\ 93.89 \pm 0.19 \\ 94.89 \pm 0.19 \end{array}$ | $\begin{array}{c} 68.72 \pm 1.05 \\ 70.94 \pm 0.02 \\ 72.11 \pm 0.08 \\ 73.11 \pm 0.36 \\ 74.40 \pm 0.16 \end{array}$ | $\begin{array}{c} 47.57 \pm 2.54 \\ 63.64 \pm 1.98 \\ 73.29 \pm 0.64 \\ 76.50 \pm 1.01 \\ 72.90 \pm 0.12 \end{array}$ | $\begin{array}{c} 23.75 \pm 4.33 \\ 29.82 \pm 0.81 \\ 29.15 \pm 0.24 \\ 29.45 \pm 0.35 \\ 26.77 \pm 0.68 \end{array}$ | $\begin{array}{c} 82.33 \pm 2.08 \\ 91.67 \pm 0.58 \\ 94.33 \pm 1.53 \\ 96.00 \pm 0.00 \\ 96.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 41.51 \pm 0.67 \\ 43.94 \pm 0.18 \\ 44.97 \pm 0.05 \\ 49.81 \pm 0.34 \\ 50.83 \pm 0.09 \end{array}$ | $\begin{array}{c} 71.27 \pm 24.23 \\ 76.43 \pm 23.01 \\ 78.18 \pm 23.03 \\ 79.74 \pm 22.47 \\ 79.35 \pm 23.04 \end{array}$ | 1.00/0.99 0.99/0.97 0.97/0.93 0.88/0.80 0.50/0.34 |
| | 100 w/clusters (C) | 16 C 5 16 C 7 | 98.00 98.00 | 88.00 88.00 | 99.00 99.00 | 93.38 93.67 | 91.67 94.00 | 72.02 72.97 | 76.80 76.83 | 27.74 29.91 | 96.00 95.00 | 46.06 47.33 | 78.87 79.47 | 1.00/0.95 1.00/0.93 |
| | 500 diagonal (D) | 16 D 32 D 64 D 128 D 256 D | $\begin{array}{c} 54.44 \pm 6.87 \\ 58.08 \pm 11.52 \\ 69.21 \pm 2.03 \\ 79.77 \pm 0.37 \\ 93.83 \pm 2.52 \end{array}$ | $\begin{array}{c} 47.00 \pm 2.83 \\ 47.00 \pm 7.07 \\ 54.50 \pm 9.19 \\ 79.50 \pm 2.12 \\ 85.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 82.21 \pm 3.59 \\ 82.06 \pm 1.69 \\ 88.33 \pm 4.04 \\ 95.89 \pm 2.83 \\ 99.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 73.38 \pm {\rm 14.97} \\ 78.62 \pm {\rm 11.23} \\ 91.11 \pm {\rm 0.38} \\ 91.89 \pm {\rm 1.39} \\ 93.78 \pm {\rm 0.19} \end{array}$ | $\begin{array}{c} 80.08 \pm 3.71 \\ 85.57 \pm 1.48 \\ 88.78 \pm 0.38 \\ 88.67 \pm 0.00 \\ 90.56 \pm 0.38 \end{array}$ | $\begin{array}{c} 51.02 \pm 5.31 \\ 52.98 \pm 3.81 \\ 62.36 \pm 3.52 \\ 65.92 \pm 1.79 \\ 68.95 \pm 1.92 \end{array}$ | $\begin{array}{c} 17.49 \pm 1.10 \\ 21.73 \pm 3.95 \\ 33.36 \pm 6.69 \\ 44.98 \pm 10.98 \\ 49.39 \pm 4.36 \end{array}$ | $\begin{array}{ c c c c c c c c c c c c c c c c c c c$ | $\begin{array}{c} 72.67 \pm 6.03 \\ 75.33 \pm 4.04 \\ 77.67 \pm 2.31 \\ 81.00 \pm 5.00 \\ 87.33 \pm 2.31 \end{array}$ | $\begin{array}{c} 39.65 \pm 0.28 \\ 39.78 \pm 0.42 \\ 40.42 \pm 0.98 \\ 40.34 \pm 0.80 \\ 42.15 \pm 0.73 \end{array}$ | $\begin{array}{c} 53.16 \pm 24.97 \\ 55.66 \pm 25.48 \\ 62.16 \pm 26.05 \\ 67.82 \pm 26.35 \\ 72.83 \pm 25.93 \end{array}$ | 1.00 / 1.00 1.00 / 1.00 1.00 / 0.99 1.00 / 0.98 1.00 / 0.97 |
| | 500 full (F) | 16 F 32 F 64 F 128 F 256 F | $\begin{array}{c} 54.30 \pm 1.13 \\ 75.10 \pm 4.92 \\ 96.94 \pm 0.42 \\ 97.67 \pm 0.58 \\ 98.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 37.00 \pm 5.66 \\ 46.50 \pm 3.54 \\ 82.50 \pm 0.71 \\ 83.50 \pm 2.12 \\ 88.50 \pm 0.71 \end{array}$ | $\begin{array}{c} 77.67 \pm 0.58 \\ 91.67 \pm 1.53 \\ 93.33 \pm 0.58 \\ 98.00 \pm 0.00 \\ 99.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 91.00 \pm 0.00 \\ 91.56 \pm 0.19 \\ 93.89 \pm 0.69 \\ 93.56 \pm 0.19 \\ 93.78 \pm 0.19 \end{array}$ | $\begin{array}{c} 90.56 \pm 0.19 \\ 91.56 \pm 0.38 \\ 90.67 \pm 0.00 \\ 92.00 \pm 0.00 \\ 93.00 \pm 0.88 \end{array}$ | $\begin{array}{c} 62.47 \pm 0.79 \\ 67.37 \pm 0.83 \\ 72.30 \pm 0.71 \\ 71.92 \pm 0.19 \\ 72.45 \pm 0.38 \end{array}$ | $\begin{array}{c} 47.56 \pm 0.29 \\ 51.17 \pm 0.81 \\ 54.63 \pm 0.79 \\ 65.02 \pm 0.81 \\ 73.77 \pm 1.21 \end{array}$ | $\begin{array}{c} 14.18 \pm 0.67 \\ 13.44 \pm 0.02 \\ 14.49 \pm 0.27 \\ 28.49 \pm 0.55 \\ 27.59 \pm 0.39 \end{array}$ | $\begin{array}{c} 79.00 \pm 1.00 \\ 81.67 \pm 1.53 \\ 86.33 \pm 0.58 \\ 93.00 \pm 0.00 \\ 95.33 \pm 0.58 \end{array}$ | $\begin{array}{c} 41.58 \pm 0.23 \\ 41.92 \pm 0.42 \\ 43.16 \pm 0.08 \\ 43.85 \pm 0.12 \\ 43.81 \pm 0.17 \end{array}$ | $\begin{array}{c} 60.31 \pm 24.42 \\ 65.84 \pm 25.64 \\ 72.49 \pm 26.64 \\ 76.47 \pm 23.77 \\ 78.18 \pm 24.16 \end{array}$ | 1.00/1.00 0.99/0.99 0.97/0.96 0.88/0.86 0.50/0.47 |
| | 500 w/clusters (C) | 16 C 7 16 C 10 16 C 25 64 C 5 64 C 7 | 95.00 96.00 95.00 98.00 | 86.00 87.00 86.00 84.00 88.00 | 98.00 99.00 99.00 99.00 99.00 99.00 | 93.67 93.00 92.71 93.67 94.00 | 91.67 91.33 93.00 92.67 93.33 | 71.19 69.93 72.13 72.32 72.18 | 54.69 53.48 74.59 75.60 75.83 | 20.03 30.09 30.21 27.17 28.14 | 90.00 94.00 95.00 96.00 96.00 | 46.34 44.89 46.66 44.43 47.68 | 74.66 75.87 78.53 77.99 79.22 | 1.00/0.98 1.00/0.98 1.00/0.95 0.97/0.93 0.97/0.91 |

Table 5: Absolute In-Distribution ROUGE-L scores for various tasks and methods

1274 G.10 CONVERGENCE

Table 13 presents outcomes where the JD-Full algorithm is executed until convergence. Our convergence criterion is defined as follows:

$$\max\left(\|U_{t+1} - U_t U_t^\top U_{t+1}\|_{\mathrm{Fro}} / \|U_{t+1}\|_{\mathrm{Fro}}, \|V_{t+1} - V_t V_t^\top V_{t+1}\|_{\mathrm{Fro}} / \|V_{t+1}\|_{\mathrm{Fro}} \right) < \tau$$
(19)

where the tolerance threshold τ is set to 0.001. Due to the slow per-iteration computation times of the primary JD-Full algorithm, which quickly reaches an approximate optimum but then has a long tail of convergence for final digits of precision, we devised an alternative eigenvalue iteration algorithm (Appendix A.2) optimized for GPU acceleration. Our analysis indicates that adherence to this convergence criterion does not significantly alter the results.

1284 1285

1286

1271

1272 1273

1278

G.11 OUT-OF-DISTRIBUTION PERFORMANCE (LORA-HUB)

For completeness, we incorporate results using the protocol of LoRA-hub (Huang et al., 2024). That 1287 is, 100 LoRA-adapters are sampled, independent of the evaluation task, representing a measure of 1288 out-of-distribution performance. This also means that each result on a task is averaged across all 1289 100 LoRA-adapters (as there is no a priori LoRA-to-task mapping). These results were obtained 1290 without normalizing the LoRA-adapters before applying the JD algorithms, a step we later identified 1291 as beneficial. We present performance comparison in Table 15. Table 14 presents the average 1292 agreement between uncompressed and compressed LoRA across 10 evaluation tasks. Results per 1293 task for JD-diagonal and JD-full are shown in Table 16 and Table 17, respectively. 1294

1295 From these tables, we find that the JD algorithms successfully maintain performance in this out-ofdistribution context.

| 1299 | | | | | | | | | | | | | | |
|------|--------------------|-------------------|---|---|---|---|---|---|---|---|---|---|---|----------------------------|
| 1300 | | | | | | | | | | | | | | |
| 1301 | | | | | | | | | | | | | | |
| 1302 | | | | | | | | | | | | | | |
| 1303 | | | | | | | | | | | | | | |
| 1304 | | | | | | | | | | | | | | |
| 1305 | | | | | | | | | | | | | | |
| 1306 | | | | | | | | | | | | | | |
| 1307 | Model Type | Method Type | | | | | Ta | sks | | | | | Average | Para. Saved |
| 1308 | | base | task039 | task190 | task280 0.19 ± 0.00 | task290 0.42 ± 0.00 | task391 0.11 ± 0.00 | task442 | task620 | task1342 0.26 ± 0.00 | task1391 0.19 ± 0.00 | task1598 | 0.29 ± 0.22 | 1.00/1.00 |
| 1309 | | lora | 1.00 ± 0.00 | 1.00 ± 0.00 | 1.00 ± 0.00 | 1.00 ± 0.00 | 1.00 ± 0.00 | 1.00 ± 0.00 | 1.00 ± 0.00 | 1.00 ± 0.00 | 1.00 ± 0.00 | 1.00 ± 0.00 | 1.00 ± 0.00 | 0.00/0.00 |
| 1310 | TIEC | 10 50 | 0.81 ± 0.00 0.59 ± 0.00 | 0.57 ± 0.02 0.41 ± 0.00 | 0.45 ± 0.04 0.18 ± 0.05 | 0.10 ± 0.01 0.03 ± 0.01 | 0.83 ± 0.01 0.91 ± 0.01 | 0.52 ± 0.00 0.34 ± 0.00 | 0.71 ± 0.01 0.67 ± 0.00 | 0.58 ± 0.00 0.62 ± 0.00 | 0.82 ± 0.01 0.32 ± 0.04 | 0.80 ± 0.00 0.78 ± 0.00 | 0.62 ± 0.22 0.48 ± 0.27 | 1.00 / 1.00 1.00 / 1.00 |
| 1311 | 11ES | 500 | $\begin{array}{c} 0.55 \pm 0.00 \\ 0.37 \pm 0.00 \end{array}$ | 0.40 ± 0.00 0.26 ± 0.00 | 0.20 ± 0.05 0.01 ± 0.00 | 0.01 ± 0.02 0.00 ± 0.00 | $\begin{array}{c} 0.88 \pm 0.00 \\ 0.83 \pm 0.00 \end{array}$ | 0.36 ± 0.00 0.31 ± 0.00 | 0.65 ± 0.00 0.58 ± 0.00 | $\begin{array}{c} 0.37 \pm 0.02 \\ 0.37 \pm 0.00 \end{array}$ | 0.01 ± 0.00 0.01 ± 0.00 | 0.78 ± 0.00 0.41 ± 0.00 | $0.44 \pm 0.29 \\ 0.32 \pm 0.26$ | 1.00 / 1.00 |
| 1312 | SVD | SVD 2 SVD 4 | $\begin{array}{c} 0.98 \pm 0.03 \\ 0.99 \pm 0.04 \end{array}$ | $\begin{array}{c} 1.07 \pm 0.02 \\ 1.04 \pm 0.01 \end{array}$ | $\begin{array}{c} 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.00 \pm 0.00 \\ 1.00 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.99 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.01 \pm 0.01 \\ 0.99 \pm 0.01 \end{array}$ | $\begin{array}{c} 1.00 \pm 0.10 \\ 0.99 \pm 0.08 \end{array}$ | $\begin{array}{c} 1.00 \pm 0.01 \\ 0.99 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.99 \pm 0.01 \\ 1.01 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.00 \pm 0.04 \\ 1.00 \pm 0.03 \end{array}$ | 0.88 / 0.88 0.75 / 0.75 |
| 1313 | 370 | SVD 8 SVD 16 | $\begin{array}{c} 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.02 \pm 0.01 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.01 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.00 \pm 0.01 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.01 \pm 0.01 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.01 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.00 \pm 0.01 \\ 1.00 \pm 0.00 \end{array}$ | 0.50 / 0.50 0.00 / 0.00 |
| 1314 | | 16 D | 1.02 ± 0.01 | 1.01 ± 0.01 | 1.00 ± 0.00 | 1.00 ± 0.01 | 0.99 ± 0.00 | 0.97 ± 0.00 | 1.03 ± 0.02 | 1.12 ± 0.03 | 0.99 ± 0.02 | 0.99 ± 0.00 | 1.01 ± 0.04 | 1.00/0.90 |
| 1315 | 10 diagonal (D) | 64 D 128 D | 1.00 ± 0.00 1.00 ± 0.00 | 1.03 ± 0.01 1.01 ± 0.01 | 1.00 ± 0.00 1.00 ± 0.00 1.00 ± 0.00 | 1.00 ± 0.00 1.00 ± 0.00 | 1.00 ± 0.00 1.00 ± 0.00 | 1.00 ± 0.00 1.00 ± 0.00 | 1.01 ± 0.01 1 01 + 0.01 | 0.99 ± 0.01 0.99 + 0.01 | 1.00 ± 0.01 1.01 ± 0.00 1.00 ± 0.00 | 1.01 ± 0.00 1.01 ± 0.00 1.00 ± 0.00 | 1.00 ± 0.01 1.00 ± 0.01 1.00 ± 0.01 | 1.00 / 0.60 |
| 1316 | | 256 D | 1.00 ± 0.00 | 1.00 ± 0.00 | 1.00 ± 0.00 | 1.00 ± 0.00 | 1.00 ± 0.00 | 1.00 ± 0.00 | 1.00 ± 0.00 | 1.00 ± 0.00 | 1.00 ± 0.00 | 1.00 ± 0.00 | 1.00 ± 0.00 | 1.00 / -0.60 |
| 1317 | 10 full (E) | 16 F 32 F | 1.02 ± 0.00 1.02 ± 0.01 1.00 ± 0.00 | 1.06 ± 0.01 1.04 ± 0.01 1.02 ± 0.01 | 1.00 ± 0.00 1.00 ± 0.00 1.00 ± 0.00 | 1.00 ± 0.00 1.00 ± 0.00 1.00 ± 0.00 | 0.99 ± 0.01 1.00 ± 0.00 1.00 ± 0.00 | 0.99 ± 0.00 0.99 ± 0.00 1.00 ± 0.00 | 1.01 ± 0.02 0.96 ± 0.01 | 1.07 ± 0.00 1.00 ± 0.02 | 1.01 ± 0.01 1.00 ± 0.01 1.01 ± 0.00 | 1.00 ± 0.00 1.01 ± 0.00 1.01 ± 0.00 | 1.02 ± 0.03 1.00 ± 0.02 1.00 ± 0.01 | 0.99/0.79 |
| 1318 | 10 101 (1-) | 128 F 256 F | 1.00 ± 0.00 1.00 ± 0.00 1.00 ± 0.00 | 1.03 ± 0.01 1.01 ± 0.01 1.00 ± 0.00 | 1.00 ± 0.00 1.00 ± 0.00 1.00 ± 0.00 | 1.00 ± 0.00 1.00 ± 0.00 1.00 ± 0.00 | 1.00 ± 0.00 1.00 ± 0.00 1.00 ± 0.00 | 1.00 ± 0.00 1.00 ± 0.00 1.00 ± 0.00 | 1.00 ± 0.00 1.00 ± 0.00 1.00 ± 0.00 | 0.98 ± 0.01 0.99 ± 0.00 1.00 ± 0.00 | 1.00 ± 0.00 1.00 ± 0.00 1.00 ± 0.00 | 1.00 ± 0.00 1.00 ± 0.00 1.00 ± 0.00 | 1.00 ± 0.01 1.00 ± 0.00 1.00 ± 0.00 | 0.88/0.07 |
| 1310 | | 16 D | 0.98 ± 0.04 | 0.98 ± 0.01 | 1.00 ± 0.00 | 0.92 ± 0.06 | 0.85 ± 0.06 | 0.94 ± 0.02 | 0.69 ± 0.05 | 0.88 ± 0.10 | 0.88 ± 0.07 | 0.86 ± 0.01 | 0.90 ± 0.10 | 1.00 / 0.98 |
| 1220 | 50 diagonal (D) | 32 D 64 D | 1.00 ± 0.02 1.02 ± 0.00 | 1.02 ± 0.02 1.05 ± 0.02 | 1.00 ± 0.00 1.00 ± 0.00 | 0.99 ± 0.00 1.00 ± 0.00 | 0.96 ± 0.01 0.99 ± 0.01 | $\begin{array}{c} 0.96 \pm 0.02 \\ 0.97 \pm 0.01 \end{array}$ | 0.85 ± 0.02 0.99 ± 0.01 | 1.00 ± 0.12 1.09 ± 0.03 | 0.98 ± 0.01 1.01 ± 0.01 | $\begin{array}{c} 0.90 \pm 0.00 \\ 0.94 \pm 0.00 \end{array}$ | 0.97 ± 0.06 1.01 ± 0.04 | 1.00 / 0.96 1.00 / 0.92 |
| 1221 | | 128 D 256 D | 1.01 ± 0.01 1.01 ± 0.01 | 1.08 ± 0.01 1.03 ± 0.01 | $1.00 \pm 0.00 \\ 1.00 \pm 0.00$ | $1.00 \pm 0.00 \\ 1.00 \pm 0.00$ | $\begin{array}{c} 0.99 \pm 0.01 \\ 1.00 \pm 0.01 \end{array}$ | 0.98 ± 0.00 1.00 ± 0.00 | $\begin{array}{c} 0.98 \pm 0.02 \\ 0.97 \pm 0.03 \end{array}$ | 1.10 ± 0.03 1.00 ± 0.03 | 1.00 ± 0.00 1.00 ± 0.01 | 1.01 ± 0.01 1.01 ± 0.00 | $1.02 \pm 0.04 \\ 1.00 \pm 0.02$ | 1.00 / 0.84 |
| 1321 | | 16 F 32 F | $\begin{array}{c} 0.99 \pm 0.04 \\ 1.02 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.00 \pm 0.01 \\ 1.00 \pm 0.02 \end{array}$ | $\begin{array}{c} 1.00 \pm 0.01 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.96 \pm 0.01 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.95 \pm 0.02 \\ 0.98 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.95 \pm 0.01 \\ 0.97 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.65 \pm 0.09 \\ 0.96 \pm 0.01 \end{array}$ | $\begin{array}{c} 1.01 \pm 0.15 \\ 1.09 \pm 0.03 \end{array}$ | $\begin{array}{c} 0.97 \pm 0.02 \\ 1.01 \pm 0.02 \end{array}$ | $\begin{array}{c} 0.88 \pm 0.01 \\ 0.93 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.94 \pm 0.11 \\ 0.99 \pm 0.04 \end{array}$ | 1.00/0.98 0.99/0.95 |
| 1000 | 50 full (F) | 64 F 128 F | $1.02 \pm 0.01 \\ 1.02 \pm 0.00$ | $\begin{array}{c} 1.06 \pm 0.02 \\ 1.06 \pm 0.01 \end{array}$ | $1.00 \pm 0.00 \\ 1.00 \pm 0.00$ | $\begin{array}{c} 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.99 \pm 0.01 \\ 1.00 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.98 \pm 0.00 \\ 0.98 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.03 \pm 0.01 \\ 0.98 \pm 0.01 \end{array}$ | $\begin{array}{c} 1.11 \pm 0.00 \\ 1.03 \pm 0.04 \end{array}$ | $\begin{array}{c} 1.00 \pm 0.01 \\ 1.00 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.99 \pm 0.01 \\ 1.01 \pm 0.00 \end{array}$ | $1.02 \pm 0.04 \\ 1.01 \pm 0.02$ | 0.97/0.89 0.88/0.72 |
| 1020 | | 256 F | 1.00 ± 0.00 0.80 ± 0.07 | 1.02 ± 0.00 | 1.00 ± 0.00 0.93 ± 0.03 | 1.00 ± 0.00 0.96 ± 0.01 | 1.00 ± 0.00 | 0.99 ± 0.00 | 1.01 ± 0.01 0.30 ± 0.07 | 1.00 ± 0.00 | 1.01 ± 0.00 0.78 ± 0.03 | 1.01 ± 0.00 | 1.00 ± 0.01 0.73 ± 0.21 | 0.50/0.18 |
| 1024 | 100 diagonal (D) | 32 D 64 D | $\begin{array}{c} 0.95 \pm 0.06 \\ 1.01 \pm 0.03 \end{array}$ | $\begin{array}{c} 0.98 \pm 0.01 \\ 1.01 \pm 0.01 \end{array}$ | 1.00 ± 0.00 1.00 ± 0.00 | $\begin{array}{c} 0.91 \pm 0.06 \\ 0.98 \pm 0.02 \end{array}$ | 0.80 ± 0.13 0.96 ± 0.01 | $0.91 \pm 0.05 \\ 0.95 \pm 0.01$ | 0.62 ± 0.10 0.90 ± 0.05 | 0.78 ± 0.25 1.11 ± 0.07 | 0.91 ± 0.02 0.96 ± 0.02 | 0.85 ± 0.01 0.88 ± 0.02 | 0.87 ± 0.14 0.98 ± 0.07 | 1.00 / 0.98 1.00 / 0.96 |
| 1020 | 0 | 128 D 256 D | $\begin{array}{c} 1.01 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.02 \pm 0.01 \\ 1.06 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.99 \pm 0.01 \\ 0.99 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.98 \pm 0.00 \\ 0.98 \pm 0.00 \end{array}$ | 1.00 ± 0.03 1.00 ± 0.01 | 1.11 ± 0.02 1.11 ± 0.03 | $\begin{array}{c} 0.99 \pm 0.01 \\ 1.00 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.92 \pm 0.00 \\ 0.99 \pm 0.02 \end{array}$ | $\begin{array}{c} 1.00 \pm 0.05 \\ 1.01 \pm 0.04 \end{array}$ | 1.00 / 0.92 1.00 / 0.84 |
| 1320 | | 16 F 32 F | 0.95 ± 0.01 1 00 ± 0.02 | 0.97 ± 0.03 0.99 ± 0.01 | 0.97 ± 0.03 1.00 ± 0.00 | 0.97 ± 0.03 1.00 ± 0.00 | 0.93 ± 0.01 0.97 ± 0.01 | 0.93 ± 0.01 0.96 ± 0.00 | 0.66 ± 0.03 0.87 ± 0.03 | 0.90 ± 0.16 1 12 ± 0.03 | 0.87 ± 0.02 0.96 ± 0.01 | 0.85 ± 0.01 0.89 ± 0.00 | 0.90 ± 0.10 0.98 ± 0.07 | 1.00/0.99 |
| 1327 | 100 full (F) | 64 F 128 F | 1.00 ± 0.02 1.02 ± 0.00 1.01 ± 0.01 | 1.00 ± 0.02 1.05 ± 0.01 | 1.00 ± 0.00 1.00 ± 0.00 1.00 ± 0.00 | 1.00 ± 0.00 0.99 ± 0.00 | 0.98 ± 0.00 1.00 ± 0.00 | 0.97 ± 0.00 0.98 ± 0.00 | 0.99 ± 0.01 1.03 ± 0.01 | 1.10 ± 0.01 1.10 ± 0.01 | 0.99 ± 0.02 1.01 ± 0.00 | 0.93 ± 0.01 1.00 ± 0.00 | 1.00 ± 0.04 1.02 ± 0.03 | 0.97/0.93 |
| 1328 | | 256 F | 1.01 ± 0.01 | 1.03 ± 0.01 | 1.00 ± 0.00 | 1.00 ± 0.00 | 1.01 ± 0.00 | 1.00 ± 0.00 | 0.98 ± 0.00 | 1.00 ± 0.03 | 1.01 ± 0.00 | 1.01 ± 0.00 | 1.00 ± 0.01 | 0.50/0.34 |
| 1329 | 100 w/clusters (C) | 16 C J 16 C 7 | 1.12 | 1.02 | 1.00 | 1.00 | 1.00 | 0.97 | 1.01 | 1.19 | 1.04 | 0.94 | 1.03 | 1.00/0.93 |
| 1004 | | 16 D 32 D | $\begin{array}{c} 0.57 \pm 0.07 \\ 0.61 \pm 0.12 \end{array}$ | $\begin{array}{c} 0.55 \pm 0.03 \\ 0.55 \pm 0.08 \end{array}$ | $\begin{array}{c} 0.83 \pm 0.04 \\ 0.83 \pm 0.02 \end{array}$ | $\begin{array}{c} 0.78 \pm 0.16 \\ 0.84 \pm 0.12 \end{array}$ | $\begin{array}{c} 0.85 \pm 0.04 \\ 0.91 \pm 0.02 \end{array}$ | $\begin{array}{c} 0.73 \pm 0.07 \\ 0.75 \pm 0.05 \end{array}$ | $\begin{array}{c} 0.24 \pm 0.02 \\ 0.30 \pm 0.05 \end{array}$ | $\begin{array}{c} 0.45 \pm 0.01 \\ 0.49 \pm 0.07 \end{array}$ | $\begin{array}{c} 0.76 \pm 0.06 \\ 0.79 \pm 0.04 \end{array}$ | $\begin{array}{c} 0.81 \pm 0.00 \\ 0.82 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.66 \pm 0.20 \\ 0.69 \pm 0.20 \end{array}$ | 1.00 / 1.00 1.00 / 1.00 |
| 1331 | 500 diagonal (D) | 64 D 128 D | $\begin{array}{c} 0.73 \pm 0.02 \\ 0.84 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.63 \pm 0.11 \\ 0.92 \pm 0.02 \end{array}$ | $\begin{array}{c} 0.89 \pm 0.04 \\ 0.97 \pm 0.03 \end{array}$ | $\begin{array}{c} 0.97 \pm 0.00 \\ 0.98 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.94 \pm 0.00 \\ 0.94 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.86 \pm 0.03 \\ 0.90 \pm 0.02 \end{array}$ | $\begin{array}{c} 0.46 \pm 0.09 \\ 0.62 \pm 0.14 \end{array}$ | $\begin{array}{c} 0.51 \pm 0.07 \\ 0.54 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.82 \pm 0.02 \\ 0.85 \pm 0.05 \end{array}$ | $\begin{array}{c} 0.83 \pm 0.01 \\ 0.83 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.77 \pm 0.18 \\ 0.84 \pm 0.15 \end{array}$ | 1.00 / 0.99 1.00 / 0.98 |
| 1332 | | 256 D 16 F | 0.99 ± 0.03 0.57 ± 0.01 | 0.99 ± 0.00 0.43 ± 0.07 | 1.00 ± 0.00 0.78 ± 0.01 | 1.00 ± 0.00 0.97 ± 0.00 | 0.96 ± 0.00 0.96 ± 0.00 | 0.93 ± 0.02 0.86 ± 0.01 | 0.68 ± 0.05 0.65 ± 0.00 | 0.85 ± 0.14 0.55 ± 0.02 | 0.92 ± 0.02 0.83 ± 0.01 | 0.85 ± 0.00 0.84 ± 0.00 | 0.92 ± 0.11 0.76 ± 0.17 | 1.00/0.97 |
| 1333 | 500 full (F) | 32 F 64 F | $\begin{array}{c} 0.79 \pm 0.05 \\ 1.02 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.54 \pm 0.04 \\ 0.96 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.93 \pm 0.02 \\ 0.94 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.98 \pm 0.00 \\ 1.00 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.97 \pm 0.00 \\ 0.96 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.92 \pm 0.00 \\ 0.97 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.70 \pm 0.01 \\ 0.74 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.52 \pm 0.00 \\ 0.55 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.86 \pm 0.02 \\ 0.91 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.85 \pm 0.00 \\ 0.87 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.81 \pm 0.16 \\ 0.89 \pm 0.14 \end{array}$ | 0.99/0.99 0.97/0.96 |
| 1334 | | 128 F 256 F | $\begin{array}{c} 1.03 \pm 0.01 \\ 1.03 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.97 \pm 0.02 \\ 1.03 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.99 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.98 \pm 0.00 \\ 0.99 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.97 \pm 0.00 \\ 0.97 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.88 \pm 0.01 \\ 1.00 \pm 0.02 \end{array}$ | $\begin{array}{c} 1.07 \pm 0.02 \\ 1.04 \pm 0.02 \end{array}$ | $\begin{array}{c} 0.98 \pm 0.00 \\ 1.00 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.90 \pm 0.00 \\ 0.93 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.98 \pm 0.05 \\ 1.00 \pm 0.03 \end{array}$ | 0.88/0.86 0.50/0.47 |
| 1335 | | 16 C 7 16 C 10 | 1.09 | 1.00 | 0.99 1.00 | 1.00 0.99 | 0.98 0.97 | 0.96 | 0.72 0.72 | 0.88 | 0.98 1.02 | 0.93 0.92 | 0.95 | 1.00/0.98 1.00/0.98 |
| 1336 | 500 w/clusters (C) | 16 C 25 64 C 5 | 1.10 1.09 | 1.00 0.98 | 1.00 1.00 | 0.99 1.00 | 0.99 0.99 | 0.97 0.97 | 0.98 0.99 | 1.30 1.17 | 1.03 1.04 | 0.96 0.93 | 1.03 1.02 | 1.00/0.95 0.97/0.93 |
| 1337 | | 64 C 7 | 1.12 | 1.02 | 1.00 | 1.00 | 1.00 | 0.97 | 1.00 | 1.22 | 1.04 | 0.99 | 1.04 | 0.97/0.91 |

Table 6: Relative In-Distribution ROUGE-1 scores for various tasks and methods

| Model Type | Method Type | | t1-100 | 41-290 | t1-200 | Tas | ks | t1-(20 | t1-1242 | t1-1201 | t1-1500 | Average | Para. Saved |
|-------------------|-----------------------|---|---|--|---|---|--|--|---|--|--|---|----------------------------|
| | base | 24.44 ± 0.00 | 1.60 ± 0.00 | 19.13 ± 0.00 | 39.22 ± 0.00 | 10.42 ± 0.00 | task442 39.88 ± 0.00 | 8.05 ± 0.00 | 6.96 ± 0.00 | 17.82 ± 0.00 | task1598 55.03 ± 0.00 | 22.43 ± 16.49 | 1.00 / 1.00 |
| | lora | 95.00 ± 0.00 | 86.00 ± 0.00 | 99.00 ± 0.00 | 93.67 ± 0.00 | 94.33 ± 0.00 | 78.43 ± 0.00 | 74.90 ± 0.00 | 26.87 ± 0.00 | 95.00 ± 0.00 | 68.66 ± 0.00 | 81.14 ± 20.67 | 0.00/0.00 |
| TIFS | 50 | 55.80 ± 0.00 52.43 ± 0.00 | 35.00 ± 0.00 34.00 ± 0.00 | $\frac{14.05 \pm 4.04}{18.00 \pm 5.20}$ 19.67 ± 4.62 | 2.42 ± 0.50 1.09 ± 1.66 | 85.78 ± 0.96 83.33 ± 0.96 | 26.75 ± 0.00 28.57 ± 0.00 | 49.96 ± 0.00 48.89 ± 0.00 | 15.46 ± 0.12 16.73 ± 0.00 15.18 ± 0.42 | 30.00 ± 3.46 1.00 ± 0.00 | 53.87 ± 0.02 53.44 ± 0.02 | 37.43 ± 23.49 33.76 ± 25.22 | 1.00 / 1.00 |
| TILS | 500 | 35.18 ± 0.00 | 22.00 ± 0.00 | 1.00 ± 0.00 | 0.00 ± 0.00 | 78.00 ± 0.00 | 24.32 ± 0.00 | 43.80 ± 0.04 | 9.96 ± 0.13 | 1.00 ± 0.00 1.00 ± 0.00 | 27.90 ± 0.03 | 24.40 ± 23.79 | 1.00 / 1.00 |
| SVD | SVD 2 SVD 4 | $\begin{array}{c} 93.15 \pm 2.77 \\ 94.01 \pm 3.60 \end{array}$ | $\begin{array}{c} 92.24 \pm 1.85 \\ 89.21 \pm 0.71 \end{array}$ | $\begin{array}{c} 99.09 \pm 0.18 \\ 99.05 \pm 0.09 \end{array}$ | $\begin{array}{c} 93.44 \pm 0.14 \\ 93.65 \pm 0.03 \end{array}$ | $\begin{array}{c} 93.89 \pm 0.35 \\ 94.66 \pm 0.63 \end{array}$ | $\begin{array}{c} 77.33 \pm 0.29 \\ 78.42 \pm 0.23 \end{array}$ | $\begin{array}{c} 75.40 \pm 1.01 \\ 74.09 \pm 1.12 \end{array}$ | $\begin{array}{c} 26.90 \pm 2.68 \\ 26.47 \pm 2.06 \end{array}$ | $\begin{array}{c} 95.06 \pm 1.35 \\ 93.98 \pm 0.77 \end{array}$ | $\begin{array}{c} 67.71 \pm 0.49 \\ 69.37 \pm 0.21 \end{array}$ | $\begin{array}{c} 81.33 \pm 20.85 \\ 81.22 \pm 20.80 \end{array}$ | 0.88 / 0.88 0.75 / 0.75 |
| | SVD 8 SVD 16 | 95.00 ± 0.00 95.00 ± 0.00 | 87.40 ± 0.59 86.00 ± 0.00 | $\begin{array}{c} 99.05 \pm 0.09 \\ 99.00 \pm 0.00 \end{array}$ | 93.65 ± 0.03 93.67 ± 0.00 | $\begin{array}{c} 94.36 \pm 0.38 \\ 94.33 \pm 0.00 \end{array}$ | $78.21 \pm 0.03 \\ 78.44 \pm 0.03$ | 75.57 ± 0.00 74.73 ± 0.18 | 26.88 ± 0.27 26.87 ± 0.00 | 95.51 ± 1.09 95.00 ± 0.00 | $\begin{array}{c} 69.33 \pm 0.08 \\ 68.62 \pm 0.04 \end{array}$ | 83.02 ± 19.87 80.76 ± 21.05 | 0.50/0.50 0.00/0.00 |
| | 16 D 32 D | 96.67 ± 0.58 95.67 ± 0.58 | 87.00 ± 1.00 90.00 ± 1.00 | $\begin{array}{c} 99.00 \pm 0.00 \\ 99.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 94.00 \pm 0.67 \\ 93.00 \pm 0.33 \end{array}$ | 93.11 ± 0.38 94.89 ± 0.51 | 76.08 ± 0.17 77.46 ± 0.24 | 77.26 ± 1.47 72.53 ± 1.00 | 30.15 ± 0.72 27.98 ± 0.71 | $\begin{array}{c} 94.00 \pm 1.73 \\ 94.67 \pm 0.58 \end{array}$ | $\begin{array}{c} 68.25 \pm 0.18 \\ 69.16 \pm 0.41 \end{array}$ | 81.55 ± 20.03 81.44 ± 20.80 | 1.00 / 0.90 1.00 / 0.80 |
| 10 diagonal (D) |) 64 D 128 D | $\begin{array}{c} 95.00 \pm 0.00 \\ 95.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 88.33 \pm 0.58 \\ 86.67 \pm 0.58 \end{array}$ | $\begin{array}{c} 99.00 \pm 0.00 \\ 99.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 93.67 \pm 0.00 \\ 93.67 \pm 0.00 \end{array}$ | $\begin{array}{c} 94.78 \pm 0.38 \\ 94.33 \pm 0.00 \end{array}$ | $\begin{array}{c} 78.28 \pm 0.07 \\ 78.45 \pm 0.16 \end{array}$ | $\begin{array}{c} 75.47 \pm 0.58 \\ 75.46 \pm 0.51 \end{array}$ | $\begin{array}{c} 26.53 \pm 0.25 \\ 26.64 \pm 0.23 \end{array}$ | $\begin{array}{c} 96.00 \pm 0.00 \\ 95.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 69.36 \pm 0.05 \\ 68.70 \pm 0.14 \end{array}$ | $\begin{array}{c} 81.64 \pm {\scriptstyle 21.06} \\ 81.29 \pm {\scriptstyle 20.92} \end{array}$ | 1.00 / 0.60 1.00 / 0.20 |
| | 256 D | 95.00 ± 0.00 | 86.00 ± 0.00 | 99.00 ± 0.00 99.00 ± 0.00 | 93.67 ± 0.00 93.56 ± 0.19 | 94.33 ± 0.00 | 78.43 ± 0.00 | 74.90 ± 0.00 | 26.87 ± 0.00 28.71 ± 0.09 | 95.00 ± 0.00 96.00 ± 1.00 | 68.59 ± 0.03 68.69 ± 0.08 | 81.18 ± 20.86 82.09 ± 20.68 | 1.00/-0.60 |
| 10 full (F) | 32 F 64 F | 96.67 ± 0.58 95.00 ± 0.00 | $\begin{array}{c} 89.33 \pm 0.58 \\ 88.67 \pm 0.58 \end{array}$ | 99.00 ± 0.00 99.00 ± 0.00 | 93.22 ± 0.19 93.67 ± 0.00 | 94.44 ± 0.19 94.56 ± 0.38 | 77.84 ± 0.21 78.19 ± 0.08 | 72.24 ± 0.59 75.97 ± 0.58 | 26.84 ± 0.50 26.43 ± 0.34 | $\begin{array}{c} 94.67 \pm 0.58 \\ 96.00 \pm 0.00 \end{array}$ | 69.55 ± 0.08 69.38 ± 0.11 | 81.38 ± 21.11 81.69 ± 21.07 | 0.99/0.79 0.97/0.57 |
| | 128 F 256 F | $\begin{array}{c} 95.00 \pm 0.00 \\ 95.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 86.67 \pm 0.58 \\ 86.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 99.00 \pm 0.00 \\ 99.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 93.67 \pm 0.00 \\ 93.67 \pm 0.00 \end{array}$ | $\begin{array}{c} 94.33 \pm 0.00 \\ 94.33 \pm 0.00 \end{array}$ | $\begin{array}{c} 78.46 \pm 0.03 \\ 78.44 \pm 0.03 \end{array}$ | $\begin{array}{c} 74.90 \pm 0.00 \\ 74.79 \pm 0.19 \end{array}$ | $\begin{array}{c} 26.72 \pm 0.13 \\ 26.87 \pm 0.00 \end{array}$ | $\begin{array}{c} 95.00 \pm 0.00 \\ 95.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 68.65 \pm 0.03 \\ 68.64 \pm 0.03 \end{array}$ | $\begin{array}{c} 81.24 \pm 20.91 \\ 81.17 \pm 20.86 \end{array}$ | 0.88/0.07 0.50/-1.10 |
| | 16 D 32 D | 92.76 ± 3.53 95.33 ± 2.08 | 84.67 ± 1.15 87 33 ± 2.08 | 99.00 ± 0.00 99.00 ± 0.00 | 86.17 ± 5.81 92.60 ± 0.29 | 79.83 ± 6.08 90.35 ± 1.00 | 73.55 ± 1.39 75.43 ± 1.33 | 51.72 ± 3.78 | 23.75 ± 2.66 26.97 ± 3.21 | 83.90 ± 6.43 93.33 ± 1.15 | 59.05 ± 0.94 61.94 ± 0.32 | 73.44 ± 22.08 | 1.00/0.98 |
| 50 diagonal (D) |) 64 D 128 D | 97.00 ± 0.00 96.33 ± 0.58 | 90.33 ± 1.53 92.67 ± 0.58 | 99.00 ± 0.00 99.00 ± 0.00 | 93.78 ± 0.19 93.56 ± 0.19 | 93.00 ± 0.58 93.00 ± 0.58 | 76.27 ± 0.49 77.24 ± 0.19 | 74.39 ± 0.90 73.76 ± 1.25 | 29.28 ± 0.81 29.58 ± 0.93 | 95.67 ± 0.58 95.00 ± 0.00 | 64.84 ± 0.27 69.04 ± 0.54 | 81.36 ± 20.83 81.92 ± 20.44 | 1.00 / 0.92 1.00 / 0.84 |
| | 256 D | 95.67 ± 0.58 | 88.33 ± 0.58 | 99.00 ± 0.00 98.67 ± 0.58 | 93.56 ± 0.19 | 94.67 ± 0.67 | 78.45 ± 0.14 | 72.86 ± 2.07 | 27.00 ± 0.77 | 95.33 ± 0.58 | 69.61 ± 0.18 60.26 ± 1.03 | 81.45 ± 21.00 | 1.00 / 0.68 |
| 50 full (F) | 32 F 64 F | 97.00 ± 0.00 96.67 ± 0.58 | 85.67 ± 1.53 91.00 ± 2.00 | 99.00 ± 0.00 99.00 ± 0.00 | 93.67 ± 0.00 93.56 ± 0.19 | 92.22 ± 0.69 93.22 ± 0.51 | 75.86 ± 0.22 77.17 ± 0.38 | 71.68 ± 0.65 77.11 ± 0.51 | 29.26 ± 0.70 29.75 ± 0.03 | 95.67 ± 1.53 95.33 ± 0.58 | 63.88 ± 0.10 68.13 ± 0.75 | 80.39 ± 20.81 82.09 ± 20.33 | 0.99/0.95 |
| | 128 F 256 F | $\begin{array}{c} 97.00 \pm 0.00 \\ 95.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 91.00 \pm 1.00 \\ 88.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 99.00 \pm 0.00 \\ 99.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 93.33 \pm 0.00 \\ 93.67 \pm 0.00 \end{array}$ | $\begin{array}{c} 94.11 \pm 0.51 \\ 94.44 \pm 0.19 \end{array}$ | $\begin{array}{c} 77.23 \pm 0.17 \\ 77.97 \pm 0.24 \end{array}$ | $\begin{array}{c} 73.67 \pm 0.58 \\ 75.47 \pm 0.58 \end{array}$ | $\begin{array}{c} 27.62 \pm 1.12 \\ 26.96 \pm 0.09 \end{array}$ | $\begin{array}{c} 95.00 \pm 1.00 \\ 96.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 69.40 \pm 0.16 \\ 69.28 \pm 0.05 \end{array}$ | $\begin{array}{c} 81.74 \pm 20.97 \\ 81.58 \pm 20.92 \end{array}$ | 0.88/0.72 0.50/0.18 |
| | 16 D | 76.43 ± 7.07 | 76.67 ± 4.93 | 91.61 ± 2.75 99.00 + 0.00 | 89.99 ± 1.07 | 47.89 ± 8.62 | 63.17 ± 1.31 | 22.23 ± 5.27 | 14.46 ± 2.89 | 73.93 ± 3.13 | 57.17 ± 1.05 | 61.35 ± 25.78 | 1.00 / 0.99 |
| 100 diagonal (D | 0) 64 D 128 D | 95.56 ± 2.49 96.00 + 0.00 | 86.67 ± 0.58 87.33 + 1.15 | 99.00 ± 0.00 99.00 ± 0.00 | 92.24 ± 1.68 93.89 ± 0.19 | 90.89 ± 1.17 93.00 + 0.58 | 74.57 ± 0.50 76.68 ± 0.18 | 67.07 ± 3.81 74.84 + 2.23 | 29.78 ± 1.92 29.79 ± 0.50 | 91.67 ± 2.31 93.67 ± 0.58 | 60.28 ± 1.51 63.49 ± 0.34 | 78.77 ± 20.77 80.77 + 20.47 | 1.00 / 0.96 |
| | 256 D | 95.00 ± 0.00 | 91.00 ± 0.00 | 99.00 ± 0.00 | 93.56 ± 0.19 | 93.11 ± 0.19 | 76.93 ± 0.23 | 75.13 ± 0.84 | 29.75 ± 0.73 | 95.33 ± 0.58 | 67.89 ± 1.34 | 81.67 ± 20.28 | 1.00 / 0.84 |
| 100 full (F) | 32 F 64 F | 95.33 ± 1.53 97.00 ± 0.00 | 85.00 ± 2.65 85.00 ± 1.00 85.67 ± 1.52 | 99.00 ± 3.00 99.00 ± 0.00 99.00 ± 0.00 | 93.50 ± 0.22 93.78 ± 0.10 | 91.44 ± 0.84 92.56 ± 0.10 | 75.00 ± 0.93 75.00 ± 0.19 76.01 ± 0.12 | $+5.41 \pm 2.04$ 65.09 ± 2.23 73.96 ± 0.89 | 24.17 ± 4.22 30.20 ± 0.81 29.46 ± 0.21 | 91.67 ± 0.58 94.33 ± 1.52 | 50.10 ± 0.44 60.92 ± 0.26 64.07 ± 0.37 | 78.72 ± 20.72 80.58 ± 20.59 | 0.99/0.97 |
| 100 Iun (I') | 128 F 256 F | 96.33 ± 0.58 96.33 ± 0.58 | 90.33 ± 0.58 88.67 ± 0.58 | 99.00 ± 0.00 99.00 ± 0.00 | 93.00 ± 0.00 93.67 ± 0.00 | 93.89 ± 0.19 94.89 ± 0.19 | 77.04 ± 0.30 78.16 ± 0.18 | 77.33 ± 1.01 73.40 ± 0.12 | 29.49 ± 0.35 26.86 ± 0.68 | 96.00 ± 0.00 96.00 ± 0.00 | 68.76 ± 0.25 69.47 ± 0.23 | 82.12 ± 20.35 81.64 ± 21.15 | 0.88/0.80 0.50/0.34 |
| 100 w/clusters (0 | C) 16 C 5 | 98.00 98.00 | 88.00 88.00 | 99.00 99.00 | 93.38 93.67 | 91.67 94.00 | 75.97 | 77.63 77.67 | 27.91 | 96.00 95.00 | 64.18 66.78 | 81.17 | 1.00/0.95 |
| | 16 D | 54.44 ± 6.87 | 47.00 ± 2.83 | 82.21 ± 3.59 | 73.38 ± 14.97 | 80.13 ± 3.68 | 57.42 ± 5.29 | 18.33 ± 1.33 | 12.19 ± 0.30 | 72.67 ± 6.03 | 55.79 ± 0.20 | 55.64 ± 24.25 | 1.00 / 1.00 |
| 500 diagonal (D |) 32 D 64 D | 58.08 ± 11.52 69.21 ± 2.03 | 47.00 ± 7.07 54.50 ± 9.19 | 82.06 ± 1.69 88.33 ± 4.04 | 78.62 ± 11.23 91.11 ± 0.38 | 85.57 ± 1.48 88.78 ± 0.38 | 59.19 ± 3.70 67.71 ± 2.59 | $\begin{array}{c} 22.76 \pm 3.95 \\ 34.79 \pm 6.86 \end{array}$ | 13.15 ± 1.94 13.80 ± 1.95 | 75.33 ± 4.04 77.67 ± 2.31 | 56.07 ± 0.52 56.78 ± 0.73 | 58.16 ± 24.56 64.61 ± 24.79 | 1.00 / 1.00 1.00 / 0.99 |
| | 128 D 256 D | /9.7/ ± 0.37 93.83 ± 2.52 | /9.50 ± 2.12 85.00 ± 0.00 | $\begin{array}{c} 95.89 \pm 2.83 \\ 99.00 \pm 0.00 \end{array}$ | $\begin{array}{r} 91.89 \pm 1.39 \\ 93.78 \pm 0.19 \end{array}$ | $\begin{array}{c} 88.6/\pm0.00\\ 90.56\pm0.38\end{array}$ | /0.27 ± 1.73 73.25 ± 1.86 | $\begin{array}{c} 46.64 \pm 10.58 \\ 51.14 \pm 3.86 \end{array}$ | $\begin{array}{c} 14.63 \pm 0.25 \\ 22.93 \pm 3.86 \end{array}$ | $\begin{array}{c} 81.00 \pm 5.00 \\ 87.33 \pm 2.31 \end{array}$ | $\begin{array}{c} 56.88 \pm 0.55 \\ 58.48 \pm 0.20 \end{array}$ | 70.20 ± 24.63 75.20 ± 23.90 | 1.00 / 0.98 1.00 / 0.97 |
| | 16 F 32 F | $\begin{array}{c} 54.30 \pm {\scriptstyle 1.13} \\ 75.10 \pm {\scriptstyle 4.92} \end{array}$ | $\begin{array}{c} 37.00 \pm 5.66 \\ 46.50 \pm 3.54 \end{array}$ | $\begin{array}{c} 77.67 \pm 0.58 \\ 91.67 \pm 1.53 \end{array}$ | $\begin{array}{c} 91.00 \pm 0.00 \\ 91.56 \pm 0.19 \end{array}$ | $\begin{array}{c} 90.56 \pm 0.19 \\ 91.56 \pm 0.38 \end{array}$ | $\begin{array}{c} 67.63 \pm 0.45 \\ 72.03 \pm 0.15 \end{array}$ | $\begin{array}{c} 48.81 \pm 0.35 \\ 52.63 \pm 0.86 \end{array}$ | $\begin{array}{c} 14.70 \pm 0.65 \\ 13.93 \pm 0.02 \end{array}$ | $\begin{array}{c} 79.00 \pm 1.00 \\ 81.67 \pm 1.53 \end{array}$ | $\begin{array}{c} 57.66 \pm 0.19 \\ 58.50 \pm 0.20 \end{array}$ | $\begin{array}{c} 62.69 \pm {23.46} \\ 68.24 \pm {24.29} \end{array}$ | 1.00/1.00 0.99/0.99 |
| 500 full (F) | 64 F 128 F | 96.94 ± 0.42 97.67 ± 0.58 | 82.50 ± 0.71 83.50 ± 2.12 | $\begin{array}{c} 93.33 \pm 0.58 \\ 98.00 \pm 0.00 \\ \end{array}$ | 93.89 ± 0.69 93.56 ± 0.19 | 90.67 ± 0.00 92.00 ± 0.00 | $\begin{array}{c} 75.99 \pm 0.64 \\ 75.80 \pm 0.16 \\ 75.22 \end{array}$ | 55.63 ± 1.07 66.19 ± 0.81 | 14.74 ± 0.27 28.67 ± 0.49 | $\begin{array}{c} 86.33 \pm 0.58 \\ 93.00 \pm 0.00 \\ 05.22 \end{array}$ | $\begin{array}{c} 59.43 \pm 0.05 \\ 61.53 \pm 0.13 \\ 62.50 \end{array}$ | 74.69 ± 25.01 78.84 ± 21.50 | 0.97/0.96 |
| | 256 F 16 C 7 | 98.00 ± 0.00 95.00 | 88.50 ± 0.71 86.00 | 99.00 ± 0.00 98.00 | 93.78 ± 0.19 93.67 | 93.00 ± 0.88 91.67 | 76.33 ± 0.29 75.10 | 74.60 ± 1.21 55.52 | 27.82 ± 0.42 20.50 | 95.33 ± 0.58 90.00 | 63.70 ± 0.14 63.57 | 80.75 ± 21.60 76.90 | 0.50/0.47 |
| 500 w/clusters (0 | C) 16 C 10 16 C 25 | 96.00 96.00 | 87.00 86.00 | 99.00 99.00 | 93.00 92.71 | 91.33 93.00 | 74.17 76.42 | 55.14 75.42 | 30.29 30.40 | 94.00 95.00 | 63.09 66.07 | 78.30 81.00 | 1.00/0.98 1.00/0.95 |
| | 64 C 5 64 C 7 | 95.00 98.00 | 84.00 88.00 | 99.00 99.00 | 93.67 94.00 | 92.67 93.33 | 76.45 76.42 | 76.43 76.67 | 27.49 28.48 | 96.00 96.00 | 64.10 68.00 | 80.48 81.79 | 0.97/0.93 0.97/0.91 |
| | | | | | | | | | | | | | |

Table 7: Absolute In-Distribution ROUGE-1 scores for various tasks and methods

| 1407 | | | | | | | | | | | | | | |
|-------|--------------------|-------------------|---|---|---|---|---|---|---|---|---|---|---|----------------------------|
| 1408 | | | | | | | | | | | | | | |
| 1409 | | | | | | | | | | | | | | |
| 1410 | | | | | | | | | | | | | | |
| 1411 | | | | | | | | | | | | | | |
| 1412 | | | | | | | | | | | | | | |
| 1413 | | | | | | | | | | | | | | |
| 1/1/ | | | | | | | | | | | | | | |
| 1/15 | Model Type | Method Type | | | | | Ta | sks | | | | | Average | Para Saved |
| 1415 | | linearou Type | task039 | task190 | task280 | task290 | task391 | task442 | task620 | task1342 | task1391 | task1598 | | |
| 1410 | | lora | $\begin{array}{c c} 0.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $0.02 \pm 0.00 \\ 1.00 \pm 0.00$ | $\begin{array}{c} 0.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | 0.00 ± 0.00 1.00 ± 0.00 | 0.00 ± 0.01 1.00 ± 0.00 | 1.00 / 1.00 0.00 / 0.00 |
| 1417 | | 10 | 0.69 ± 0.00 | 0.57 ± 0.02 | 0.45 ± 0.04 | 0.10 ± 0.01 | 0.57 ± 0.03 | 0.00 ± 0.00 | 0.39 ± 0.01 | 0.21 ± 0.00 | 0.82 ± 0.01 | 0.00 ± 0.00 | 0.38 ± 0.28 | 1.00/1.00 |
| 1418 | TIES | 100 | 0.43 ± 0.00 0.41 ± 0.00 | 0.41 ± 0.00 0.40 ± 0.00 | 0.13 ± 0.05 0.20 ± 0.05 | 0.03 ± 0.01 0.01 ± 0.02 | 0.70 ± 0.02 0.65 ± 0.00 | 0.00 ± 0.00 0.00 ± 0.00 | 0.36 ± 0.00 0.36 ± 0.00 | 0.21 ± 0.00 0.21 ± 0.00 | 0.01 ± 0.00 | 0.00 ± 0.00 0.00 ± 0.00 | 0.27 ± 0.22 0.23 ± 0.22 | 1.00 / 1.00 |
| 1419 | | SVD 2 | 0.22 ± 0.00 | 1.07 ± 0.02 | 1.00 ± 0.00 | 0.00 ± 0.00 0.99 ± 0.01 | 0.00 ± 0.00 0.98 ± 0.01 | 0.00 ± 0.00 0.98 ± 0.03 | 0.32 ± 0.00 0.94 ± 0.01 | 1.03 ± 0.17 | 1.00 ± 0.01 | 0.00 ± 0.00 0.15 ± 0.29 | 0.13 ± 0.20 0.91 ± 0.28 | 0.88/0.88 |
| 1420 | SVD | SVD 4 SVD 8 | $0.99 \pm 0.04 \\ 1.00 \pm 0.00$ | 1.04 ± 0.01 1.02 ± 0.01 | 1.00 ± 0.00 1.00 ± 0.00 | $1.00 \pm 0.00 \\ 1.00 \pm 0.00$ | 1.01 ± 0.02 1.00 ± 0.01 | 1.11 ± 0.00 1.02 ± 0.05 | $0.97 \pm 0.02 \\ 1.00 \pm 0.00$ | 0.99 ± 0.13 1.00 ± 0.00 | $\begin{array}{c} 0.99 \pm 0.01 \\ 1.01 \pm 0.01 \end{array}$ | 0.90 ± 0.17 1.00 ± 0.00 | 1.00 ± 0.08 1.00 ± 0.02 | 0.75 / 0.75 0.50 / 0.50 |
| 1421 | | SVD 16 | 1.00 ± 0.00 | 1.00 ± 0.00 | 1.00 ± 0.00 | 1.00 ± 0.00 | 1.00 ± 0.00 | 1.00 ± 0.00 | 0.99 ± 0.01 | 1.00 ± 0.00 | 1.00 ± 0.00 | 1.00 ± 0.00 | 1.00 ± 0.00 | 0.00 / 0.00 |
| 1422 | | 16 D 32 D | 1.02 ± 0.01 1.01 ± 0.01 | $\begin{array}{c} 1.01 \pm 0.01 \\ 1.05 \pm 0.01 \end{array}$ | $\begin{array}{c} 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.01 \pm 0.02 \\ 0.98 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.96 \pm 0.01 \\ 1.02 \pm 0.02 \end{array}$ | 1.11 ± 0.11 1.11 ± 0.00 | $\begin{array}{c} 0.89 \pm 0.03 \\ 0.93 \pm 0.01 \end{array}$ | $\begin{array}{c} 1.19 \pm 0.04 \\ 1.10 \pm 0.04 \end{array}$ | $\begin{array}{c} 0.99 \pm 0.02 \\ 1.00 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.33 \pm 0.58 \\ 0.67 \pm 0.58 \end{array}$ | $\begin{array}{c} 0.95 \pm 0.27 \\ 0.98 \pm 0.19 \end{array}$ | 1.00 / 0.90 1.00 / 0.80 |
| 1423 | 10 diagonal (D) | 64 D 128 D | 1.00 ± 0.00 1.00 ± 0.00 | $1.03 \pm 0.01 \\ 1.01 \pm 0.01$ | 1.00 ± 0.00 1.00 ± 0.00 | $\begin{array}{c} 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | 1.02 ± 0.01 1.00 ± 0.00 | $1.11 \pm 0.00 \\ 1.00 \pm 0.00$ | $\begin{array}{c} 0.99 \pm 0.01 \\ 1.00 \pm 0.01 \end{array}$ | 1.00 ± 0.00 1.00 ± 0.00 | $\begin{array}{c} 1.01 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.67 \pm 0.58 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.98 \pm 0.19 \\ 1.00 \pm 0.00 \end{array}$ | 1.00 / 0.60 1.00 / 0.20 |
| 1424 | | 256 D | 1.00 ± 0.00 1.02 ± 0.00 | 1.00 ± 0.00 1.06 ± 0.01 | 1.00 ± 0.00 1.00 ± 0.00 | 1.00 ± 0.00 1.00 ± 0.01 | 1.00 ± 0.00 0.97 ± 0.03 | 1.00 ± 0.00 1.15 ± 0.06 | 1.00 ± 0.00 0.92 ± 0.02 | 1.00 ± 0.00 1.17 ± 0.04 | 1.00 ± 0.00 1.01 ± 0.01 | 1.00 ± 0.00 0.67 ± 0.58 | 1.00 ± 0.00 1.00 ± 0.20 | 1.00/0.90 |
| 1425 | 10 full (F) | 32 F 64 F | 1.02 ± 0.01 1.00 ± 0.00 | 1.04 ± 0.01 1.03 ± 0.01 | 1.00 ± 0.00 1.00 ± 0.00 | 0.98 ± 0.01 1.00 ± 0.00 | 1.00 ± 0.01 1.01 ± 0.01 | 1.11 ± 0.00 1.07 ± 0.06 | 0.92 ± 0.01 1.01 ± 0.01 | 1.02 ± 0.04 1.00 ± 0.00 | 1.00 ± 0.01 1.01 ± 0.00 | 1.00 ± 0.00 1.00 ± 0.00 | 1.01 ± 0.05 1.01 ± 0.03 | 0.99/0.79 |
| 1426 | | 128 F 256 F | $1.00 \pm 0.00 \\ 1.00 \pm 0.00$ | $\begin{array}{c} 1.01 \pm 0.01 \\ 1.00 \pm 0.00 \end{array}$ | $1.00 \pm 0.00 \\ 1.00 \pm 0.00$ | $1.00 \pm 0.00 \\ 1.00 \pm 0.00$ | $1.00 \pm 0.00 \\ 1.00 \pm 0.00$ | 1.00 ± 0.00 1.00 ± 0.00 | $1.00 \pm 0.00 \\ 1.00 \pm 0.01$ | $1.00 \pm 0.00 \\ 1.00 \pm 0.00$ | $1.00 \pm 0.00 \\ 1.00 \pm 0.00$ | $1.00 \pm 0.00 \\ 1.00 \pm 0.00$ | $1.00 \pm 0.00 \\ 1.00 \pm 0.00$ | 0.88/0.07 |
| 1427 | | 16 D | 0.91 ± 0.06 | 0.98 ± 0.01 | 1.00 ± 0.00 | 0.91 ± 0.09 | 0.78 ± 0.05 | 0.89 ± 0.29 | 0.34 ± 0.06 | 0.50 ± 0.45 | 0.86 ± 0.07 | 0.00 ± 0.00 | 0.72 ± 0.35 | 1.00 / 0.98 |
| 1/100 | 50 diagonal (D) | 32 D 64 D | 1.00 ± 0.02 1.02 ± 0.00 | 1.02 ± 0.02 1.05 ± 0.02 | 1.00 ± 0.00 1.00 ± 0.00 | 1.00 ± 0.01 1.00 ± 0.01 | $\begin{array}{c} 0.90 \pm 0.03 \\ 0.95 \pm 0.02 \end{array}$ | $\begin{array}{c} 0.85 \pm 0.42 \\ 1.15 \pm 0.17 \end{array}$ | $\begin{array}{c} 0.56 \pm 0.04 \\ 0.81 \pm 0.03 \end{array}$ | 0.98 ± 0.23 1.14 ± 0.00 | $\begin{array}{c} 0.98 \pm 0.01 \\ 1.01 \pm 0.01 \end{array}$ | 0.00 ± 0.00 0.00 ± 0.00 | $\begin{array}{c} 0.83 \pm 0.34 \\ 0.91 \pm 0.33 \end{array}$ | 1.00 / 0.96 1.00 / 0.92 |
| 1420 | | 128 D 256 D | 1.01 ± 0.01 1.01 ± 0.01 | 1.08 ± 0.01 1.03 ± 0.01 | 1.00 ± 0.00 1.00 ± 0.00 | 1.00 ± 0.01 1.00 ± 0.01 | $\begin{array}{c} 0.95 \pm 0.02 \\ 1.01 \pm 0.02 \end{array}$ | 1.04 ± 0.06 1.11 ± 0.00 | $\begin{array}{c} 0.92 \pm 0.03 \\ 0.95 \pm 0.04 \end{array}$ | 1.21 ± 0.07 1.02 ± 0.04 | 1.00 ± 0.00 1.00 ± 0.01 | 0.67 ± 0.58 1.00 ± 0.00 | 0.99 ± 0.20 1.01 ± 0.04 | 1.00 / 0.84 1.00 / 0.68 |
| 1429 | | 16 F 32 F | $\begin{array}{ c c c c c c c c c c c c c c c c c c c$ | $\begin{array}{c} 1.00 \pm 0.01 \\ 1.00 \pm 0.02 \end{array}$ | $\begin{array}{c} 1.00 \pm 0.01 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.95 \pm 0.04 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.87 \pm 0.01 \\ 0.92 \pm 0.03 \end{array}$ | 1.04 ± 0.06 1.15 ± 0.06 | $\begin{array}{c} 0.31 \pm 0.08 \\ 0.73 \pm 0.04 \end{array}$ | $\begin{array}{c} 0.98 \pm 0.23 \\ 1.17 \pm 0.04 \end{array}$ | $\begin{array}{c} 0.97 \pm 0.02 \\ 1.01 \pm 0.02 \end{array}$ | $\begin{array}{c} 0.00 \pm 0.00 \\ 0.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.81 \pm 0.35 \\ 0.90 \pm 0.33 \end{array}$ | 1.00/0.98 0.99/0.95 |
| 1430 | 50 full (F) | 64 F 128 F | 1.02 ± 0.01 1.02 ± 0.00 | $\begin{array}{c} 1.06 \pm 0.02 \\ 1.06 \pm 0.01 \end{array}$ | $\begin{array}{c} 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.00 \pm 0.01 \\ 0.99 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.96 \pm 0.02 \\ 0.99 \pm 0.02 \end{array}$ | 1.22 ± 0.00 1.15 ± 0.06 | $\begin{array}{c} 0.94 \pm 0.01 \\ 0.92 \pm 0.01 \end{array}$ | $\begin{array}{c} 1.17 \pm 0.04 \\ 1.10 \pm 0.08 \end{array}$ | 1.00 ± 0.01 1.00 ± 0.01 | $\begin{array}{c} 0.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.94 \pm 0.33 \\ 1.02 \pm 0.07 \end{array}$ | 0.97/0.89 0.88/0.72 |
| 1431 | | 256 F | 1.00 ± 0.00 | 1.02 ± 0.00 | 1.00 ± 0.00 | 1.00 ± 0.00 | 1.00 ± 0.01 | 1.04 ± 0.06 | 0.99 ± 0.00 | 1.00 ± 0.00 | 1.01 ± 0.00 | 1.00 ± 0.00 | 1.01 ± 0.02 | 0.50/0.18 |
| 1432 | | 16 D 32 D | 0.54 ± 0.16 0.85 ± 0.15 | $\begin{array}{c} 0.89 \pm 0.06 \\ 0.98 \pm 0.01 \end{array}$ | 0.90 ± 0.04 1.00 ± 0.00 | 0.89 ± 0.05 0.86 ± 0.13 | 0.42 ± 0.08 0.70 ± 0.14 | 0.44 ± 0.00 0.74 ± 0.28 | 0.08 ± 0.02 0.28 ± 0.07 | 0.00 ± 0.00 0.48 ± 0.55 | $\begin{array}{c} 0.76 \pm 0.05 \\ 0.91 \pm 0.02 \end{array}$ | 0.00 ± 0.00 0.00 ± 0.00 | 0.49 ± 0.36 0.68 ± 0.36 | 1.00 / 0.99 1.00 / 0.98 |
| 1433 | 100 diagonal (D) | 64 D 128 D | 1.00 ± 0.04 1.01 ± 0.00 | 1.01 ± 0.01 1.02 ± 0.01 | 1.00 ± 0.00 1.00 ± 0.00 | 0.98 ± 0.02 1.01 ± 0.01 | $\begin{array}{c} 0.88 \pm 0.04 \\ 0.95 \pm 0.02 \end{array}$ | 1.07 ± 0.06 1.11 ± 0.00 | 0.58 ± 0.09 0.81 ± 0.06 | 1.10 ± 0.04 1.21 ± 0.00 | 0.96 ± 0.02 0.99 ± 0.01 | 0.00 ± 0.00 0.00 ± 0.00 | 0.86 ± 0.32 0.91 ± 0.33 | 1.00/0.96 |
| 1434 | | 256 D 16 F | 1.00 ± 0.00 0.85 ± 0.03 | 1.06 ± 0.00 0.97 ± 0.03 | 1.00 ± 0.00 0.97 ± 0.03 | 1.00 ± 0.01 0.95 ± 0.06 | 0.96 ± 0.01 0.80 ± 0.02 | 1.11 ± 0.11 0.81 ± 0.17 | 0.92 ± 0.02 0.29 ± 0.04 | 1.21 ± 0.07 0.60 ± 0.34 | 1.00 ± 0.01 0.87 ± 0.02 | 0.00 ± 0.00 0.00 ± 0.00 | 0.93 ± 0.33 0.71 ± 0.33 | 1.00/0.99 |
| 1435 | 100 full (F) | 32 F 64 F | $\begin{array}{c} 0.99 \pm 0.02 \\ 1.02 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.99 \pm 0.01 \\ 1.00 \pm 0.02 \end{array}$ | $\begin{array}{c} 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.00 \pm 0.01 \\ 1.00 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.90 \pm 0.03 \\ 0.94 \pm 0.01 \end{array}$ | $\begin{array}{c} 1.04 \pm 0.06 \\ 1.04 \pm 0.06 \end{array}$ | $\begin{array}{c} 0.55 \pm 0.04 \\ 0.78 \pm 0.01 \end{array}$ | $\begin{array}{c} 1.07 \pm 0.07 \\ 1.14 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.96 \pm 0.01 \\ 0.99 \pm 0.02 \end{array}$ | $\begin{array}{c} 0.00 \pm 0.00 \\ 0.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.85 \pm 0.32 \\ 0.89 \pm 0.31 \end{array}$ | 0.99/0.97 0.97/0.93 |
| 1436 | | 128 F 256 F | 1.01 ± 0.01 1.01 ± 0.01 | $\begin{array}{c} 1.05 \pm 0.01 \\ 1.03 \pm 0.01 \end{array}$ | $\begin{array}{c} 1.00 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.98 \pm 0.00 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.98 \pm 0.01 \\ 1.02 \pm 0.01 \end{array}$ | $\begin{array}{c} 1.15 \pm 0.06 \\ 1.19 \pm 0.06 \end{array}$ | $\begin{array}{c} 0.94 \pm 0.01 \\ 0.93 \pm 0.01 \end{array}$ | $1.21 \pm 0.00 \\ 1.02 \pm 0.04$ | $\begin{array}{c} 1.01 \pm 0.00 \\ 1.01 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.33 \pm 0.58 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.97 \pm 0.28 \\ 1.02 \pm 0.06 \end{array}$ | 0.88/0.80 0.50/0.34 |
| 1437 | 100 w/clusters (C) | 16 C 5 | 1.13 | 1.02 | 1.00 | 1.01 | 0.92 | 1.24 | 0.86 | 1.33 | 1.04 | 0.00 | 0.96 | 1.00/0.95 |
| 1438 | | 16 D | 0.22 ± 0.10 | 0.55 ± 0.03 | 0.81 ± 0.05 | 0.33 ± 0.49 | 0.70 ± 0.03 | 0.15 ± 0.17 | 0.03 ± 0.01 | 0.00 ± 0.00 | 0.76 ± 0.06 | 0.00 ± 0.00 | 0.35 ± 0.35 | 1.00/1.00 |
| 1439 | 500 diagonal (D) | 32 D 64 D | $\begin{array}{c} 0.27 \pm 0.18 \\ 0.40 \pm 0.04 \end{array}$ | $\begin{array}{c} 0.55 \pm 0.08 \\ 0.63 \pm 0.11 \end{array}$ | $\begin{array}{c} 0.82 \pm 0.02 \\ 0.89 \pm 0.04 \end{array}$ | $\begin{array}{c} 0.49 \pm 0.37 \\ 0.91 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.75 \pm 0.01 \\ 0.80 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.22 \pm 0.11 \\ 0.48 \pm 0.06 \end{array}$ | $\begin{array}{c} 0.05 \pm 0.05 \\ 0.13 \pm 0.04 \end{array}$ | $\begin{array}{c} 0.02 \pm 0.04 \\ 0.05 \pm 0.08 \end{array}$ | $\begin{array}{c} 0.79 \pm 0.04 \\ 0.82 \pm 0.02 \end{array}$ | $\begin{array}{c} 0.00 \pm 0.00 \\ 0.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.39 \pm 0.34 \\ 0.51 \pm 0.35 \end{array}$ | 1.00 / 1.00 1.00 / 0.99 |
| 1440 | | 128 D 256 D | $\begin{array}{c} 0.61 \pm 0.04 \\ 0.95 \pm 0.02 \end{array}$ | $\begin{array}{c} 0.92 \pm 0.02 \\ 0.99 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.97 \pm 0.03 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.93 \pm 0.05 \\ 1.00 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.80 \pm 0.00 \\ 0.86 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.74 \pm 0.17 \\ 0.85 \pm 0.28 \end{array}$ | $\begin{array}{c} 0.22 \pm 0.11 \\ 0.28 \pm 0.06 \end{array}$ | $\begin{array}{c} 0.12 \pm 0.08 \\ 0.55 \pm 0.39 \end{array}$ | $\begin{array}{c} 0.85 \pm 0.05 \\ 0.92 \pm 0.02 \end{array}$ | $\begin{array}{c} 0.00 \pm 0.00 \\ 0.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.61 \pm 0.36 \\ 0.73 \pm 0.36 \end{array}$ | 1.00 / 0.98 1.00 / 0.97 |
| 1441 | | 16 F 32 F | 0.21 ± 0.02 0.54 ± 0.02 | 0.43 ± 0.07 0.54 ± 0.04 | 0.78 ± 0.01 0.93 ± 0.02 | 0.90 ± 0.00 0.92 ± 0.01 | 0.86 ± 0.01 0.90 + 0.01 | 0.59 ± 0.06 0.63 ± 0.12 | 0.21 ± 0.01 0.26 ± 0.02 | 0.12 ± 0.04 0.14 + 0.00 | 0.83 ± 0.01 0.86 ± 0.02 | 0.00 ± 0.00 0.00 ± 0.00 | 0.50 ± 0.34 0.57 ± 0.34 | 1.00/1.00 |
| 1442 | 500 full (F) | 64 F 128 F | 0.99 ± 0.03 1.02 ± 0.03 | 0.96 ± 0.01 0.97 ± 0.02 | 0.93 ± 0.02 0.94 ± 0.01 0.99 ± 0.00 | 1.01 ± 0.03 1.00 + 0.03 | 0.87 ± 0.00 0.92 ± 0.00 | 1.04 ± 0.17 1.15 ± 0.04 | 0.36 ± 0.02 0.61 ± 0.01 | 0.14 ± 0.00 1.07 ± 0.00 | 0.91 ± 0.01 0.98 + 0.00 | 0.00 ± 0.00 0.00 ± 0.00 0.00 ± 0.00 | 0.71 ± 0.39 0.87 ± 0.39 | 0.97/0.96 |
| 1443 | | 256 F | 1.02 ± 0.01 | 1.03 ± 0.01 | 1.00 ± 0.00 | 1.00 ± 0.01 1.00 ± 0.01 | 0.92 ± 0.00 0.95 ± 0.03 | 1.00 ± 0.00 | 0.78 ± 0.01 | 1.07 ± 0.00 1.07 ± 0.00 | 1.00 ± 0.01 | 0.00 ± 0.00 | 0.88 ± 0.31 | 0.50/0.47 |
| 1/// | | 16 C 7 16 C 10 | 1.08 | 1.00 | 0.99 1.00 | 1.00 0.98 | 0.92 0.91 | 1.01 | 0.39 0.37 | 0.62 | 0.98 1.02 | 0.00 | 0.80 0.89 | 1.00/0.98 1.00/0.98 |
| 1444 | 500 w/clusters (C) | 16 C 25 64 C 5 | 1.10 | 1.00 0.98 | 1.00 | 0.99 | 0.97 0.96 | 1.12 | 0.81 0.83 | 1.42 1.33 | 1.03 | 0.00 | 0.95 0.94 | 1.00/0.95 0.97/0.93 |
| 1445 | | 64 C 7 | 1.13 | 1.02 | 1.00 | 1.01 | 0.98 | 1.12 | 0.90 | 1.42 | 1.04 | 0.00 | 0.96 | 0.97/0.91 |

Table 8: Relative In-Distribution exact match scores for various tasks and methods

| Madal Tura | Mathad Tona | | | | | Ta | sks | | | | | A | Dama Caused |
|--------------------|--|---|---|---|---|---|---|---|---|---|---|---|-----------------------------|
| Model Type | Metilou Type | task039 | task190 | task280 | task290 | task391 | task442 | task620 | task1342 | task1391 | task1598 | Average | Para. Saveu |
| | base lora | $\begin{array}{c} 8.59 \pm 0.08 \\ 0.36 \pm 0.01 \end{array}$ | $\begin{array}{c} 9.15 \pm 0.00 \\ 0.17 \pm 0.00 \end{array}$ | $\begin{array}{c} 2.55 \pm 0.00 \\ 0.01 \pm 0.00 \end{array}$ | $\begin{array}{c} 2.88 \pm 0.00 \\ 0.12 \pm 0.00 \end{array}$ | $\begin{array}{c} 2.34 \pm 0.00 \\ 0.11 \pm 0.00 \end{array}$ | $\begin{array}{c} 3.46 \pm 0.04 \\ 0.76 \pm 0.02 \end{array}$ | $\begin{array}{c} 6.40 \pm 0.18 \\ 1.17 \pm 0.07 \end{array}$ | $\begin{array}{c} 5.55 \pm 0.00 \\ 1.94 \pm 0.00 \end{array}$ | $\begin{array}{c} 8.60 \pm 0.00 \\ 0.16 \pm 0.00 \end{array}$ | $\begin{array}{c} 2.67 \pm 0.00 \\ 0.85 \pm 0.00 \end{array}$ | $\begin{array}{c} 5.19 \pm 2.65 \\ 0.57 \pm 0.59 \end{array}$ | 1.00 / 1.00 0.00 / 0.00 |
| | SVD 2 | 0.32 ± 0.01 | 0.15 ± 0.00 | 0.01 ± 0.00 | 0.12 ± 0.00 | 0.10 ± 0.00 | 0.76 ± 0.02 | 1.13 ± 0.08 | 1.94 ± 0.00 | 0.13 ± 0.00 | 0.97 ± 0.00 | 0.57 ± 0.60 | 0.88 / 0.88 |
| SVD | SVD 4 SVD 8 | 0.33 ± 0.01 0.35 ± 0.01 | 0.16 ± 0.00 0.17 ± 0.00 | 0.01 ± 0.00 0.01 ± 0.00 | 0.12 ± 0.00 0.12 ± 0.00 | 0.11 ± 0.00 0.11 ± 0.00 | 0.76 ± 0.02 0.77 ± 0.02 | 1.14 ± 0.08 1.16 ± 0.07 | 1.94 ± 0.00 1.94 ± 0.00 | 0.14 ± 0.00 0.15 ± 0.00 | 0.86 ± 0.00 0.84 ± 0.00 | 0.56 ± 0.59 0.51 ± 0.58 | 0.75 / 0.75 0.50 / 0.50 |
| | SVD 16 | 0.36 ± 0.01 | 0.17 ± 0.00 | 0.01 ± 0.00 | 0.12 ± 0.00 | 0.11 ± 0.00 | 0.76 ± 0.02 | 1.14 ± 0.06 | 1.94 ± 0.00 | 0.16 ± 0.00 | 0.85 ± 0.00 | 0.56 ± 0.59 | 0.00 / 0.00 |
| | 16 D 32 D | 0.33 ± 0.01 0.33 ± 0.01 | 0.15 ± 0.01 0.16 ± 0.00 | 0.01 ± 0.00 0.01 + 0.00 | 0.12 ± 0.00 0.12 + 0.00 | 0.10 ± 0.00 0.10 ± 0.00 | 0.76 ± 0.03 0.75 ± 0.02 | 1.13 ± 0.08 1.11 + 0.07 | 1.95 ± 0.01 1.93 ± 0.00 | 0.14 ± 0.00 0.14 ± 0.01 | 1.00 ± 0.02 0.88 ± 0.00 | 0.57 ± 0.61 0.55 ± 0.60 | 1.00/0.90 |
| 10 diagonal (D) | 64 D | 0.35 ± 0.01 0.35 ± 0.01 | 0.17 ± 0.00 0.17 ± 0.00 | 0.01 ± 0.00 0.01 ± 0.00 | 0.12 ± 0.00 0.12 ± 0.00 | 0.11 ± 0.00 0.11 ± 0.00 | 0.75 ± 0.02 0.75 ± 0.02 | 1.11 ± 0.07 1.11 ± 0.07 | 1.93 ± 0.00 1.94 ± 0.00 | 0.14 ± 0.01 0.15 ± 0.00 | 0.84 ± 0.00 | 0.55 ± 0.59 0.55 ± 0.59 | 1.00 / 0.60 |
| | 128 D 256 D | 0.35 ± 0.01 0.36 ± 0.01 | 0.17 ± 0.00 0.17 ± 0.00 | 0.01 ± 0.00 0.01 ± 0.00 | 0.12 ± 0.00 0.12 ± 0.00 | 0.11 ± 0.00 0.11 ± 0.00 | 0.75 ± 0.02 0.75 ± 0.02 | 1.11 ± 0.07 1.12 ± 0.07 | 1.94 ± 0.00 1.94 ± 0.00 | 0.16 ± 0.00 0.16 ± 0.00 | 0.84 ± 0.00 0.85 ± 0.00 | 0.56 ± 0.59 0.56 ± 0.59 | 1.00 / 0.20 1.00 / -0.60 |
| | 16 F | 0.33 ± 0.00 | 0.15 ± 0.00 | 0.01 ± 0.00 | 0.12 ± 0.00 | 0.10 ± 0.00 | 0.76 ± 0.02 | 1.20 ± 0.02 | 1.95 ± 0.00 | 0.13 ± 0.00 | 0.97 ± 0.00 | 0.57 ± 0.61 | 1.00/0.90 |
| 10 full (F) | 52 F 64 F | 0.33 ± 0.01 0.34 ± 0.01 | 0.16 ± 0.00 0.16 ± 0.00 | 0.01 ± 0.00 0.01 ± 0.00 | 0.12 ± 0.00 0.12 ± 0.00 | 0.10 ± 0.00 0.11 ± 0.00 | 0.75 ± 0.02 0.75 ± 0.02 | 1.11 ± 0.07 1.11 ± 0.07 | 1.94 ± 0.00 1.94 ± 0.00 | 0.14 ± 0.00 0.15 ± 0.00 | 0.86 ± 0.00 0.84 ± 0.00 | 0.55 ± 0.60 0.55 ± 0.59 | 0.99/0.79 |
| | 128 F 256 F | 0.35 ± 0.01 0.36 ± 0.01 | 0.17 ± 0.00 0.17 ± 0.00 | 0.01 ± 0.00 0.01 ± 0.00 | 0.12 ± 0.00 0.12 ± 0.00 | 0.11 ± 0.00 0.11 ± 0.00 | 0.75 ± 0.02 0.75 ± 0.02 | 1.12 ± 0.07 1.12 ± 0.07 | 1.94 ± 0.00 1.94 ± 0.00 | 0.16 ± 0.00 0.16 ± 0.00 | 0.84 ± 0.00 0.85 ± 0.00 | 0.56 ± 0.59 0.56 ± 0.59 | 0.88/0.07 0.50/-1.10 |
| | 16 D | 0.61 ± 0.06 | 0.19 ± 0.02 | 0.03 ± 0.01 | 0.29 ± 0.04 | 0.36 ± 0.04 | 0.95 ± 0.05 | 1.73 ± 0.21 | 2.66 ± 0.22 | 0.32 ± 0.11 | 1.98 ± 0.01 | 0.91 ± 0.88 | 1.00/0.98 |
| 50 diagonal (D) | 32 D | 0.37 ± 0.02 | 0.16 ± 0.00 | 0.01 ± 0.00 | 0.19 ± 0.03 0.12 + 0.03 | 0.18 ± 0.01 | 0.85 ± 0.05 | 1.37 ± 0.14 | 2.12 ± 0.05 | 0.16 ± 0.00 | 1.65 ± 0.03 | 0.71 ± 0.73 | 1.00/0.96 |
| 50 diagonai (D) | 128 D | 0.33 ± 0.02 0.33 ± 0.01 | 0.15 ± 0.00 0.15 ± 0.00 | 0.01 ± 0.00 0.01 ± 0.00 | 0.12 ± 0.00 0.12 ± 0.00 | 0.10 ± 0.00 0.10 ± 0.00 | 0.79 ± 0.02 0.76 ± 0.03 | 1.12 ± 0.08 1.10 ± 0.05 | 1.97 ± 0.01 1.93 ± 0.01 | 0.13 ± 0.01 0.14 ± 0.00 | 1.13 ± 0.03 0.93 ± 0.01 | 0.59 ± 0.63 0.56 ± 0.60 | 1.00 / 0.92 |
| | 256 D | 0.34 ± 0.01 | 0.16 ± 0.00 | 0.01 ± 0.00 | 0.12 ± 0.00 | 0.10 ± 0.00 | 0.76 ± 0.03 | 1.11 ± 0.05 | 1.93 ± 0.00 | 0.15 ± 0.00 | 0.85 ± 0.00 | 0.55 ± 0.59 | 1.00/0.68 |
| | 16 F 32 F | 0.47 ± 0.06 0.36 ± 0.02 | 0.17 ± 0.00 0.16 ± 0.00 | 0.02 ± 0.00 0.01 ± 0.00 | 0.20 ± 0.02 0.14 ± 0.00 | 0.19 ± 0.04 0.11 ± 0.00 | 0.80 ± 0.03 0.80 ± 0.03 | 1.71 ± 0.10 1.14 ± 0.08 | 2.20 ± 0.04 2.00 ± 0.01 | 0.17 ± 0.01 0.14 ± 0.00 | 1.84 ± 0.07 1.32 ± 0.02 | 0.78 ± 0.80 0.62 ± 0.65 | 0.99/0.95 |
| 50 full (F) | 64 F 128 F | 0.33 ± 0.01 0.33 ± 0.01 | 0.15 ± 0.00 0.16 ± 0.00 | 0.01 ± 0.00 0.00 + 0.00 | 0.12 ± 0.00 0.12 + 0.00 | 0.10 ± 0.00 0.10 ± 0.00 | 0.77 ± 0.03 0.76 ± 0.03 | 1.10 ± 0.06 1.11 ± 0.05 | 1.94 ± 0.00 1.93 ± 0.00 | 0.13 ± 0.00 0.14 ± 0.00 | 1.02 ± 0.00 0.87 ± 0.00 | 0.57 ± 0.61 0.55 ± 0.60 | 0.97/0.89 |
| | 256 F | 0.35 ± 0.01 0.35 ± 0.01 | 0.16 ± 0.00 0.16 ± 0.00 | 0.00 ± 0.00 0.01 ± 0.00 | 0.12 ± 0.00 0.12 ± 0.00 | 0.10 ± 0.00 0.11 ± 0.00 | 0.76 ± 0.03 0.76 ± 0.03 | 1.11 ± 0.05 1.11 ± 0.05 | 1.94 ± 0.00 | 0.14 ± 0.00 0.15 ± 0.00 | 0.84 ± 0.00 | 0.55 ± 0.59 | 0.50/0.12 |
| | 16 D 32 D | 1.69 ± 0.49 | 0.26 ± 0.04 | 0.18 ± 0.07 | 0.34 ± 0.02 | 1.01 ± 0.20 | 1.45 ± 0.10 | 3.59 ± 0.25 | 3.72 ± 0.72 | 0.44 ± 0.20 | 2.37 ± 0.09 | 1.51 ± 1.32 | 1.00/0.99 |
| 100 diagonal (D) | 64 D | 0.39 ± 0.06 | 0.18 ± 0.01 0.16 ± 0.00 | 0.00 ± 0.03 0.01 ± 0.00 | 0.131 ± 0.08 0.18 ± 0.02 | 0.13 ± 0.08 0.14 ± 0.01 | 0.86 ± 0.02 | 1.37 ± 0.13 1.39 ± 0.07 | 2.88 ± 0.04 2.18 ± 0.04 | 0.22 ± 0.01 0.17 ± 0.00 | 1.79 ± 0.02 | 0.93 ± 0.98 0.73 ± 0.76 | 1.00 / 0.98 |
| | 128 D 256 D | $0.32 \pm 0.00 \\ 0.32 \pm 0.00$ | 0.15 ± 0.00 0.15 ± 0.00 | 0.01 ± 0.00 0.01 ± 0.00 | 0.12 ± 0.00 0.12 ± 0.00 | 0.10 ± 0.00 0.10 ± 0.00 | 0.79 ± 0.02 0.77 ± 0.02 | 1.19 ± 0.02 1.16 ± 0.00 | 2.00 ± 0.01 1.94 ± 0.00 | $\begin{array}{c} 0.14 \pm 0.01 \\ 0.13 \pm 0.00 \end{array}$ | 1.24 ± 0.04 0.96 ± 0.01 | 0.61 ± 0.65 0.56 ± 0.61 | 1.00 / 0.92 1.00 / 0.84 |
| | 16 F | 0.66 ± 0.07 | 0.19 ± 0.01 | 0.03 ± 0.01 | 0.25 ± 0.02 | 0.29 ± 0.02 | 0.99 ± 0.07 | 2.50 ± 0.51 | 2.63 ± 0.03 | 0.24 ± 0.02 | 2.21 ± 0.08 | 1.00 ± 1.01 | 1.00/0.99 |
| 100 full (F) | 32 F 64 F | 0.40 ± 0.00 0.34 ± 0.01 | 0.17 ± 0.00 0.15 ± 0.00 | 0.01 ± 0.00 0.01 ± 0.00 | 0.15 ± 0.01 0.12 ± 0.00 | 0.13 ± 0.01 0.11 ± 0.00 | 0.85 ± 0.02 0.79 ± 0.01 | 1.53 ± 0.12 1.23 ± 0.07 | 2.17 ± 0.06 1.98 ± 0.01 | 0.15 ± 0.01 0.15 ± 0.00 | 1.93 ± 0.04 1.26 ± 0.01 | 0.75 ± 0.80 0.61 ± 0.65 | 0.99/0.97 0.97/0.93 |
| | 128 F 256 F | 0.32 ± 0.00 0.33 ± 0.00 | 0.15 ± 0.00 0.16 ± 0.00 | 0.01 ± 0.00 0.00 + 0.00 | 0.12 ± 0.00 0.12 + 0.00 | 0.10 ± 0.00 0.10 ± 0.00 | 0.77 ± 0.02 0.76 ± 0.02 | 1.16 ± 0.01 1.15 ± 0.01 | 1.94 ± 0.00 1.93 ± 0.00 | 0.13 ± 0.00 0.14 ± 0.00 | 0.99 ± 0.01 0.86 ± 0.00 | 0.57 ± 0.61 0.56 ± 0.60 | 0.88/0.80 |
| 100 w/clusters (C) | 16 C 5 | 0.33 | 0.15 | 0.01 | 0.14 | 0.11 | 0.76 | 1.16 | 1.97 | 0.13 | 1.12 | 0.59 | 1.00/0.95 |
| | 16 D | 2.95 ± 0.28 | 0.73 ± 0.29 | 0.27 ± 0.09 | 0.67 ± 0.28 | 0.52 ± 0.07 | 2.06 ± 0.30 | 4.85 ± 0.31 | 3.94 ± 0.42 | 0.50 ± 0.05 | 2.50 ± 0.03 | 1.94 ± 1.59 | 1.00 / 1.00 |
| 500 diagonal (D) | 32 D | 2.33 ± 0.30 | 0.62 ± 0.17 | 0.24 ± 0.05 | 0.50 ± 0.16 | 0.37 ± 0.07 | 1.86 ± 0.25 | 4.73 ± 0.35 | 3.81 ± 0.59 | 0.39 ± 0.04 | 2.46 ± 0.05 | 1.77 ± 1.57 | 1.00 / 1.00 |
| 500 uragonar (D) | 128 D | 1.07 ± 0.18 1.12 ± 0.02 | 0.43 ± 0.00 0.23 ± 0.00 | 0.13 ± 0.04 0.04 ± 0.03 | 0.29 ± 0.02 0.21 ± 0.04 | 0.23 ± 0.02 0.22 ± 0.03 | 1.02 ± 0.28 1.08 ± 0.06 | 3.05 ± 0.36 3.05 ± 0.87 | 3.09 ± 0.37 | 0.52 ± 0.05 0.26 ± 0.03 | 2.33 ± 0.11 2.31 ± 0.04 | 1.45 ± 1.39 1.19 ± 1.21 | 1.00 / 0.99 |
| | 256 D | 0.54 ± 0.03 | 0.18 ± 0.01 | 0.01 ± 0.00 | 0.16 ± 0.01 | 0.15 ± 0.01 | 0.92 ± 0.08 | 2.42 ± 0.14 | 2.51 ± 0.13 | 0.19 ± 0.01 | 2.09 ± 0.02 | 0.94 ± 0.99 | 1.00 / 0.97 |
| | 16 F 32 F | 2.14 ± 0.06 1.17 ± 0.07 | 0.70 ± 0.04 0.48 ± 0.03 | 0.28 ± 0.00 0.08 ± 0.04 | 0.27 ± 0.01 0.21 ± 0.01 | 0.21 ± 0.00 0.17 ± 0.00 | $1.14 \pm 0.04 \\ 0.99 \pm 0.04$ | 3.06 ± 0.27 2.69 ± 0.10 | 2.71 ± 0.01 2.47 ± 0.02 | $\begin{array}{c} 0.34 \pm 0.01 \\ 0.25 \pm 0.02 \end{array}$ | 2.21 ± 0.01 2.11 ± 0.04 | 1.33 ± 1.09 1.08 ± 0.99 | 1.00/1.00 0.99/0.99 |
| 500 full (F) | 64 F 128 F | 0.51 ± 0.03 0.39 ± 0.01 | 0.21 ± 0.00 0.16 ± 0.00 | 0.02 ± 0.00 0.01 ± 0.00 | 0.17 ± 0.01 0.13 ± 0.00 | 0.14 ± 0.00 0.11 ± 0.00 | 0.88 ± 0.04 0.81 ± 0.02 | 2.19 ± 0.14 1.42 ± 0.07 | 2.34 ± 0.03 2.03 ± 0.01 | 0.20 ± 0.00 0.16 ± 0.00 | 1.97 ± 0.02 1.71 ± 0.02 | 0.89 ± 0.91 0.71 ± 0.74 | 0.97/0.96 |
| | 256 F | 0.32 ± 0.01 | 0.15 ± 0.00 | 0.01 ± 0.00 0.01 ± 0.00 | 0.12 ± 0.00 | 0.10 ± 0.00 | 0.01 ± 0.03 0.77 ± 0.01 | 1.18 ± 0.04 | 1.96 ± 0.00 | 0.14 ± 0.01 | 1.25 ± 0.00 | 0.61 ± 0.65 | 0.50/0.47 |
| | 16 C 7 16 C 10 | 0.40 | 0.18 | 0.01 | 0.15 | 0.13 | 0.90 | 2.03 | 2.21 | 0.16 | 1.50 | 0.77 | 1.00/0.98 |
| 500 w/clusters (C) | 16 C 25 | 0.32 | 0.16 | 0.01 | 0.13 | 0.10 | 0.81 | 1.28 | 1.96 | 0.12 | 1.07 | 0.60 | 1.00/0.95 |
| | 64 C 5 64 C 7 | 0.36 0.34 | 0.16 0.15 | 0.01 0.01 | 0.12 0.12 | 0.10 0.10 | 0.80 | 1.17 1.14 | 1.98 1.96 | 0.14 0.13 | 1.17 1.08 | 0.60 | 0.97/0.93 |
| 500 w/clusters (C) | 16 C 7 16 C 10 16 C 25 64 C 5 64 C 7 | 0.40 0.36 0.32 0.36 0.34 | 0.18 0.16 0.16 0.16 0.15 | 0.01 0.01 0.01 0.01 0.01 0.01 | 0.12 ± 0.00 0.15 0.14 0.13 0.12 0.12 | 0.13 0.13 0.10 0.10 0.10 0.10 | 0.90 0.87 0.81 0.80 0.79 | 2.03 2.19 1.28 1.17 1.14 | 2.21 2.04 1.96 1.98 1.96 | 0.16 0.15 0.12 0.14 0.13 | 1.50 1.38 1.07 1.17 1.08 | 0.77 0.74 0.60 0.60 0.58 | 1. 1. 1. 0. 0. |

Table 9: Absolute In-Distribution test loss for various tasks and methods

| 1514 | | | | | | | | | | | | | | |
|------|--------------------|------------------------|--|--|--|---|---|---|---|---|---|---|--|---|
| 1515 | | | | | | | | | | | | | | |
| 1516 | | | | | | | | | | | | | | |
| 1517 | | | | | | | | | | | | | | |
| 1518 | | | | | | | | | | | | | | |
| 1519 | | | | | | | | | | | | | | |
| 1520 | | | | | | | | | | | | | | |
| 1521 | | | | | | | | | | | | | | |
| 1522 | | | | | | | | | | | | | | |
| 1523 | | | | | | | | | | | | | | |
| 1524 | | | | | | | | | | | | | | |
| 1525 | Model Type | Method Type | task039 | task190 | task280 | task290 | Ta task391 | sks task442 | task620 | task1342 | task1391 | task1598 | Average | Para. Saved |
| 1526 | | base lora | $\begin{array}{c} 0.00 \pm 0.00 \\ 100.00 \pm 0.00 \end{array}$ | 0.00 ± 0.00 100.00 ± 0.00 | 1.00 ± 0.00 100.00 ± 0.00 | $\begin{array}{c} 0.00 \pm 0.00 \\ 100.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.00 \pm 0.00 \\ 100.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.00 \pm 0.00 \\ 100.00 \pm 0.00 \end{array}$ | 0.00 ± 0.00 100.00 ± 0.00 | $\begin{array}{c} 0.00 \pm 0.00 \\ 100.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.00 \pm 0.00 \\ 100.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.00 \pm 0.00 \\ 100.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.10 \pm 0.30 \\ 100.00 \pm 0.00 \end{array}$ | 1.00 / 1.00 0.00 / 0.00 |
| 1527 | | 10 50 | 41.00 ± 0.00 24.00 ± 0.00 | 53.67 ± 0.58 38.67 ± 0.58 | 44.33 ± 4.04 17.67 ± 4.62 | 10.33 ± 0.58 2.00 ± 0.00 | 46.33 ± 4.04 56.33 ± 0.58 | 1.00 ± 0.00 1.00 ± 0.00 | 8.00 ± 0.00 8.00 ± 0.00 | 8.00 ± 0.00 8.00 ± 0.00 | 76.67 ± 1.15 29.67 ± 2.89 | 1.00 ± 0.00 0.00 ± 0.00 | 29.03 ± 25.69 18.53 ± 18.07 | 1.00 / 1.00 |
| 1528 | TIES | 100 500 | $\begin{array}{c} 22.00 \pm 0.00 \\ 8.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 38.00 \pm 0.00 \\ 25.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 18.67 \pm 4.62 \\ 1.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.00 \pm {\scriptstyle 1.73} \\ 0.00 \pm {\scriptstyle 0.00} \end{array}$ | $\begin{array}{c} 51.67 \pm 4.62 \\ 59.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 1.00 \pm 0.00 \\ 0.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 8.00 \pm 0.00 \\ 3.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 7.33 \pm 0.58 \\ 6.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 2.00 \pm 0.00 \\ 2.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.00 \pm 0.00 \\ 0.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 14.97 \pm {\scriptstyle 17.20} \\ 9.90 \pm {\scriptstyle 18.12} \end{array}$ | 1.00 / 1.00 1.00 / 1.00 |
| 1529 | SVD | SVD 2 SVD 4 | $\begin{array}{c} 88.33 \pm 0.65 \\ 93.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 91.91 \pm 0.94 \\ 96.64 \pm 0.50 \end{array}$ | $\begin{array}{c} 100.00 \pm 0.00 \\ 100.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 97.25 \pm 0.45 \\ 100.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 92.83 \pm 0.39 \\ 96.75 \pm 0.87 \end{array}$ | $\begin{array}{c} 76.50 \pm 1.51 \\ 88.83 \pm 1.53 \end{array}$ | 66.00 ± 1.41 90.67 ± 1.23 | $\begin{array}{c} 58.08 \pm 1.16 \\ 72.17 \pm 0.58 \end{array}$ | $\begin{array}{c} 98.67 \pm 0.49 \\ 98.67 \pm 0.49 \end{array}$ | $\begin{array}{c} 5.83 \pm 0.94 \\ 16.67 \pm 1.78 \end{array}$ | 77.42 ± 27.69 85.24 ± 24.39 | 0.88/0.88 0.75/0.75 |
| 1530 | 310 | SVD 8 SVD 16 | $\begin{array}{c} 98.89 \pm 0.60 \\ 100.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 98.55 \pm 0.52 \\ 100.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 100.00 \pm 0.00 \\ 100.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 100.00 \pm 0.00 \\ 100.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 99.42 \pm 0.51 \\ 100.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 93.44 \pm 0.73 \\ 99.67 \pm 0.50 \end{array}$ | $\begin{array}{c} 97.22 \pm 0.44 \\ 99.50 \pm 0.55 \end{array}$ | $\begin{array}{c} 82.78 \pm 1.64 \\ 99.67 \pm 0.50 \end{array}$ | $\begin{array}{c} 99.00 \pm 0.00 \\ 100.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 60.00 \pm 0.87 \\ 98.11 \pm 0.78 \end{array}$ | $\begin{array}{c} 93.70 \pm 11.59 \\ 99.69 \pm 0.68 \end{array}$ | 0.50 / 0.50 0.00 / 0.00 |
| 1531 | | 16 D 32 D | $\begin{array}{c} 83.33 \pm 1.53 \\ 93.00 \pm 1.00 \end{array}$ | $\begin{array}{c} 88.33 \pm 0.58 \\ 95.33 \pm 0.58 \end{array}$ | $\begin{array}{c} 100.00 \pm 0.00 \\ 100.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 97.00 \pm 2.00 \\ 98.00 \pm 1.00 \end{array}$ | $\begin{array}{c} 88.33 \pm 1.15 \\ 93.67 \pm 1.53 \end{array}$ | $\begin{array}{c} 57.00 \pm 1.00 \\ 80.67 \pm 2.31 \end{array}$ | $\begin{array}{c} 48.67 \pm 3.21 \\ 78.67 \pm 1.15 \end{array}$ | $\begin{array}{c} 50.67 \pm 4.93 \\ 68.00 \pm 1.73 \end{array}$ | $\begin{array}{c} 97.67 \pm 1.53 \\ 98.33 \pm 0.58 \end{array}$ | $\begin{array}{c} 5.33 \pm 1.15 \\ 14.67 \pm 2.52 \end{array}$ | $\begin{array}{c} 71.63 \pm _{29.53} \\ 82.03 \pm _{24.99} \end{array}$ | 1.00/0.90 1.00/0.80 |
| 1532 | 10 diagonal (D) | 64 D 128 D 256 D | 99.00 ± 0.00 100.00 ± 0.00 100.00 ± 0.00 | 97.00 ± 1.00 99.33 ± 0.58 100.00 ± 0.00 | 100.00 ± 0.00 100.00 ± 0.00 100.00 ± 0.00 | 100.00 ± 0.00 100.00 ± 0.00 100.00 ± 0.00 | 98.00 ± 0.00 100.00 ± 0.00 100.00 ± 0.00 | 90.67 ± 1.53 96.67 ± 1.53 100.00 ± 0.00 | 95.33 ± 1.15 98.33 ± 1.15 100.00 ± 0.00 | 79.67 ± 1.53 95.67 ± 2.31 99.33 ± 1.15 | 99.00 ± 0.00 100.00 ± 0.00 100.00 ± 0.00 | 55.00 ± 4.36 91.67 ± 3.51 95.00 ± 1.00 | 91.37 ± 13.78 98.17 ± 2.94 99.43 ± 1.57 | 1.00/0.60 1.00/0.20 1.00/-0.60 |
| 1533 | | 16 F 32 F | 83.00 ± 2.00 91.33 ± 0.58 | 93.00 ± 1.00 96.00 ± 1.00 | 100.00 ± 0.00 100.00 ± 0.00 | 98.33 ± 0.58 98.33 ± 0.58 | 91.67 ± 0.58 94.33 ± 0.58 | 64.33 ± 3.21 84.00 ± 2.00 | 59.33 ± 1.15 83.00 ± 1.73 | 52.67 ± 1.53 70.33 ± 1.53 | 98.33 ± 0.58 99.00 ± 1.00 | 6.33 ± 1.15 22.00 ± 2.65 | 74.70 ± 28.71 83.83 ± 22.82 | 1.00/0.90 |
| 1534 | 10 full (F) | 64 F 128 F 256 F | 99.00 ± 0.00 99.67 ± 0.58 | 97.33 ± 0.58 99.33 ± 0.58 100.00 ± 0.00 | 100.00 ± 0.00 100.00 ± 0.00 100.00 ± 0.00 | $\begin{array}{c} 100.00 \pm 0.00 \\ 100.00 \pm 0.00 \\ 100.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 99.33 \pm 1.15 \\ 100.00 \pm 0.00 \\ 100.00 \pm 0.00 \end{array}$ | 91.33 ± 1.53 97.67 ± 1.15 00.67 ± 0.59 | 96.33 ± 0.58 100.00 ± 0.00 00.67 ± 0.00 | 81.67 ± 2.31 95.67 ± 1.15 00.67 ± 0.50 | $\begin{array}{c} 99.00 \pm 0.00 \\ 100.00 \pm 0.00 \\ 100.00 \pm 0.00 \end{array}$ | 58.33 ± 1.53 91.00 ± 1.00 | 92.23 ± 12.76 98.33 ± 2.89 90.70 ± 0.70 | 0.97/0.57 0.88/0.07 |
| 1535 | | 16 D | 52.67 ± 4.51 | 86.67 ± 3.06 | 100.00 ± 0.00 | 85.00 ± 3.46 | 65.33 ± 3.79 | 25.33 ± 5.03 | 10.00 ± 1.00 | 10.67 ± 10.26 | 81.00 ± 6.56 | 0.00 ± 0.00 | 51.67 ± 36.18 | 1.00 / 0.98 |
| 1536 | 50 diagonal (D) | 64 D 128 D | 69.67 ± 3.21 79.67 ± 2.52 90.00 ± 1.00 | 88.07 ± 1.53 91.00 ± 1.00 91.33 ± 0.58 | 100.00 ± 0.00 100.00 ± 0.00 100.00 ± 0.00 | 95.00 ± 2.00 97.67 ± 0.58 98.33 ± 0.58 | 80.00 ± 3.00 88.00 ± 1.00 90.67 ± 2.08 | 52.00 ± 1.00 73.67 ± 2.08 | 17.00 ± 2.63 36.67 ± 5.69 63.67 ± 1.53 | 26.33 ± 5.03 41.33 ± 1.15 56.33 ± 0.58 | 95.00 ± 2.00 96.00 ± 1.00 98.00 ± 0.00 | 0.00 ± 0.00 0.33 ± 0.58 7.33 ± 1.15 | 60.83 ± 36.02 68.27 ± 32.66 76.93 ± 27.79 | 1.00 / 0.98 1.00 / 0.92 1.00 / 0.84 |
| 1537 | | 256 D 16 F | 94.67 ± 0.58 61.67 ± 3.06 | 96.33 ± 0.58 89.67 ± 1.15 | 100.00 ± 0.00 99.67 ± 0.58 | 99.67 ± 0.58 90.67 ± 2.52 | 96.33 ± 1.15 78.33 ± 3.51 | 87.33 ± 0.58 34.00 ± 1.00 | 87.00 ± 2.65 | 71.67 ± 1.53 25.00 ± 6.24 | 99.67 ± 0.58 90.00 ± 1.00 | 31.67 ± 1.15 0.00 ± 0.00 | 86.43 ± 20.41 57.60 ± 36.59 | 1.00/0.68 |
| 1538 | 50 full (F) | 32 F 64 F | $\begin{array}{c} 71.00 \pm 1.00 \\ 81.67 \pm 0.58 \end{array}$ | $\begin{array}{c} 89.00 \pm 1.73 \\ 93.67 \pm 1.15 \\ \end{array}$ | $\begin{array}{c} 100.00 \pm 0.00 \\ 100.00 \pm 0.00 \\ 1000.00 \end{array}$ | $\begin{array}{c} 98.00 \pm 0.00 \\ 98.33 \pm 0.58 \end{array}$ | $\begin{array}{c} 85.00 \pm 1.00 \\ 90.67 \pm 2.08 \end{array}$ | $\begin{array}{c} 47.00 \pm 1.73 \\ 61.67 \pm 1.53 \end{array}$ | $\begin{array}{c} 29.00 \pm 3.00 \\ 54.33 \pm 1.53 \end{array}$ | $\begin{array}{c} 35.00 \pm 2.00 \\ 51.33 \pm 1.15 \end{array}$ | $\begin{array}{c} 98.00 \pm 1.00 \\ 98.33 \pm 0.58 \\ \end{array}$ | $\begin{array}{c} 0.00 \pm 0.00 \\ 3.33 \pm 0.58 \end{array}$ | $\begin{array}{c} 65.20 \pm 34.00 \\ 73.33 \pm 29.89 \\ 0.122 \end{array}$ | 0.99/0.95 0.97/0.89 |
| 1539 | | 128 F 256 F | 91.00 ± 1.00 97.00 ± 0.00 | 94.33 ± 0.58 98.00 ± 0.00 | 100.00 ± 0.00 100.00 ± 0.00 | 99.00 ± 0.00 100.00 ± 0.00 | 93.33 ± 1.53 99.67 ± 0.58 | 81.67 ± 0.58 92.00 ± 0.00 | 75.00 ± 1.73 94.33 ± 1.15 | 67.67 ± 2.08 79.67 ± 1.53 | 98.67 ± 0.58 99.00 ± 0.00 | 16.67 ± 0.58 57.33 ± 2.52 | 81.73 ± 24.46 91.70 ± 13.11 | 0.88/0.72 |
| 1540 | 100 disconst (D) | 16 D 32 D | 33.00 ± 8.19 51.00 ± 7.81 | 79.33 ± 5.69 90.00 ± 1.00 | 89.33 ± 5.51 100.00 ± 0.00 100.00 ± 0.00 | 80.00 ± 3.61 88.00 ± 7.00 04.22 ± 1.01 | 35.33 ± 6.03 58.67 ± 11.68 | 4.00 ± 1.73 17.67 ± 10.26 28.00 ± 2.00 | 3.00 ± 1.00 7.67 ± 2.89 10.67 ± 1.00 | 0.00 ± 0.00 9.33 ± 12.86 | 71.33 ± 4.51 86.33 ± 2.52 02.67 ± 1.15 | 0.00 ± 0.00 0.00 ± 0.00 0.22 ± 0.00 | 39.53 ± 36.15 50.87 ± 38.43 60.00 ± 34.01 | 1.00/0.99 1.00/0.98 |
| 1541 | 100 diagonal (D) | 128 D 256 D | 82.00 ± 2.00 90.00 ± 1.00 | 87.33 ± 1.33 90.00 ± 2.00 93.00 ± 2.00 | 100.00 ± 0.00 100.00 ± 0.00 100.00 ± 0.00 | 94.33 ± 4.04 97.33 ± 0.58 97.67 ± 0.58 | 80.33 ± 2.08 85.33 ± 0.58 91.67 ± 0.58 | 55.33 ± 2.08 71.67 ± 5.13 | 19.07 ± 4.31 34.33 ± 3.79 59.67 ± 1.53 | 28.33 ± 1.33 36.67 ± 2.52 58.00 ± 0.00 | 92.07 ± 1.15 94.67 ± 0.58 97.67 ± 0.58 | 0.33 ± 0.38 0.00 ± 0.00 4.00 ± 1.00 | 67.57 ± 32.95 76.33 ± 28.88 | 1.00 / 0.90 1.00 / 0.92 1.00 / 0.84 |
| 1542 | | 16 F 32 F | $\begin{array}{c} 49.00 \pm 2.00 \\ 65.00 \pm 3.46 \end{array}$ | 89.67 ± 3.21 90.33 ± 1.53 | $\begin{array}{c} 97.00 \pm 3.00 \\ 100.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 84.33 \pm {\scriptstyle 3.06} \\ 96.33 \pm {\scriptstyle 1.53} \end{array}$ | $\begin{array}{c} 65.33 \pm 2.52 \\ 80.00 \pm 2.65 \end{array}$ | $\begin{array}{c} 20.67 \pm 8.33 \\ 41.33 \pm 3.21 \end{array}$ | $\begin{array}{c} 6.33 \pm 2.08 \\ 16.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 8.33 \pm 4.73 \\ 29.33 \pm 2.08 \end{array}$ | $\begin{array}{c} 81.33 \pm 2.08 \\ 92.00 \pm 2.65 \end{array}$ | $\begin{array}{c} 0.00 \pm 0.00 \\ 0.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 50.20 \pm {}_{37.06} \\ 61.03 \pm {}_{35.43} \end{array}$ | 1.00/0.99 0.99/0.97 |
| 1543 | 100 full (F) | 64 F 128 F | 72.33 ± 0.58 84.33 ± 1.53 01.67 | 89.67 ± 1.53 92.33 ± 1.53 | $\begin{array}{c} 100.00 \pm 0.00 \\ 100.00 \pm 0.00 \\ 100.00 \pm 0.00 \end{array}$ | 97.67 ± 0.58 98.00 ± 0.00 | 86.00 ± 1.00 91.33 ± 0.58 | 53.00 ± 1.00 68.67 ± 0.58 | 35.33 ± 1.53 56.00 ± 1.00 78.00 ± 1.00 | 38.00 ± 1.73 57.67 ± 1.15 | 94.67 ± 0.58 99.00 ± 0.00 | 0.00 ± 0.00 5.33 ± 0.58 | 66.67 ± 32.54 75.27 ± 28.67 | 0.97/0.93 0.88/0.80 |
| 1544 | 100 w/clusters (C) | 16 C 5 | 74.00 78.00 | 92.00 92.00 | 100.00 ± 0.00 100.00 | 95.00 98.00 | 94.33 ± 0.38 86.00 93.00 | 54.00 58.00 | 41.00 51.00 | 47.00 48.00 | 99.00 ± 0.00 | 0.00 2.00 | 68.80 71.80 | 1.00/0.95 |
| 1545 | | 16 D | 8.00 ± 3.61 | 51.50 ± 3.54 | 79.67 ± 4.93 | 28.00 ± 42.44 | 56.67 ± 2.89 | 0.67 ± 1.15 | 0.33 ± 0.58 | 48.00 0.00 ± 0.00 | 98.00 71.67 ± 6.03 | 0.00 ± 0.00 | 28.90 ± 33.54 | 1.00/1.00 |
| 1546 | 500 diagonal (D) | 32 D 64 D 128 D | 14.55 ± 11.02 25.67 ± 1.15 38.33 ± 3.21 | 52.50 ± 9.19 62.50 ± 12.02 85.50 ± 2.12 | 80.67 ± 2.08 87.33 ± 4.04 96.00 ± 3.00 | 43.00 ± 31.48 78.33 ± 1.15 81.33 ± 2.31 | 60.67 ± 0.58 65.33 ± 3.06 65.67 ± 1.15 | 0.67 ± 1.15 5.33 ± 3.51 11.67 ± 4.73 | $1.6/ \pm 1.15$ 3.67 ± 1.15 5.33 ± 2.08 | 0.00 ± 0.00 1.00 ± 1.73 2.00 ± 1.00 | 74.33 ± 4.04 76.67 ± 2.31 80.00 ± 5.00 | 0.00 ± 0.00 0.00 ± 0.00 0.00 ± 0.00 | 32.10 ± 33.43 39.83 ± 35.80 45.24 ± 37.78 | 1.00/0.99 1.00/0.98 |
| 1547 | | 256 D 16 F | 53.33 ± 0.58 8.33 ± 2.08 | 91.00 ± 2.83 41.00 ± 5.66 | 100.00 ± 0.00 76.67 ± 0.58 | 89.00 ± 2.65 78.00 ± 0.00 | 76.00 ± 2.00 72.67 ± 0.58 | 20.67 ± 6.81 6.00 ± 0.00 | 6.00 ± 1.73 5.67 ± 0.58 | 12.00 ± 8.00 0.00 ± 0.00 | 86.33 ± 2.31 78.00 ± 1.00 | 0.00 ± 0.00 0.00 ± 0.00 | 52.14 ± 38.64 36.48 ± 35.46 | 1.00/0.97 |
| 1548 | 500 full (F) | 32 F 64 F | $\begin{array}{c} 33.67 \pm 4.16 \\ 56.00 \pm 2.65 \end{array}$ | 51.00 ± 1.41 85.50 ± 0.71 | $\begin{array}{c} 92.67 \pm 1.53 \\ 94.33 \pm 0.58 \end{array}$ | 77.00 ± 1.73 89.33 ± 2.89 | 75.00 ± 2.00 74.33 ± 1.15 | $\begin{array}{c} 14.33 \pm 1.53 \\ 36.33 \pm 1.15 \end{array}$ | 8.00 ± 0.00 9.00 ± 1.00 | 0.00 ± 0.00 2.67 ± 1.15 | $\begin{array}{c} 80.67 \pm 1.53 \\ 84.00 \pm 1.00 \end{array}$ | 0.00 ± 0.00 0.00 ± 0.00 | $\begin{array}{c} 42.97 \pm 35.80 \\ 52.03 \pm 36.92 \end{array}$ | 0.99/0.99 0.97/0.96 |
| 1549 | | 128 F 256 F | $\begin{array}{c} 69.33 \pm 0.58 \\ 79.67 \pm 0.58 \end{array}$ | 88.50 ± 0.71 89.50 ± 0.71 | $\begin{array}{c} 99.00 \pm 0.00 \\ 100.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 96.33 \pm 1.53 \\ 97.67 \pm 0.58 \end{array}$ | 80.33 ± 1.15 87.33 ± 0.58 | 45.00 ± 2.00 57.00 ± 1.00 | 16.33 ± 0.58 35.00 ± 1.00 | 31.00 ± 1.73 42.00 ± 1.00 | 92.00 ± 0.00 95.00 ± 1.00 | $\begin{array}{c} 0.00 \pm 0.00 \\ 0.00 \pm 0.00 \end{array}$ | $\begin{array}{c} 60.86 \pm 35.07 \\ 67.59 \pm 32.67 \end{array}$ | 0.88/0.86 0.50/0.47 |
| 1550 | 500 w/clusters (C) | 16 C 7 16 C 10 | 63.00 69.00 79.00 | 90.00 93.00 | 99.00 100.00 | 96.00 98.00 97.00 | 78.00 81.00 88.00 | 31.00 34.00 53.00 | 9.00 8.00 38.00 | 15.00 33.00 48.00 | 89.00 95.00 98.00 | 1.00 | 57.10 61.20 69.10 | 1.00/0.98 1.00/0.98 1.00/0.95 |
| 1551 | 500 wremsters (C) | 64 C 5 64 C 7 | 77.00 | 88.00 90.00 | 100.00 | 98.00 97.00 | 89.00 89.00 | 56.00 60.00 | 39.00 48.00 | 42.00 49.00 | 99.00 99.00 | 0.00 3.00 | 68.80 71.10 | 0.97/0.93 0.97/0.91 |
| 1552 | | | | | | | | | | | | | | |
| 1553 | | Table | 10: A | bsolute | e In-Di | stributi | on agre | ement | for var | ious ta | sks and | metho | ds | |
| 1554 | | | | | | | | | | | | | | |

| 1 | 5 | 6 | 6 | |
|---|---|---|---|--|
| 1 | 5 | 6 | 7 | |

| 568 <u>-</u> | | | | | | | | | | | | | |
|-----------------------|--------------------|---|---|---|---|---|------------------------------------|---|------------------------------------|---|---|---|------------------------------------|
| 69 | Model Type | Method Type | | | | | Ta | sks | | | | | Average |
| 70 - | | (1) (1) (1) (1) (1) (1) (1) (1) (1) (1) | task039 | task190 | task280 | task290 | task391 | task442 | task620 | task1342 | task1391 | task1598 | 0.25 |
| 70 | SVD | SVD 2 SVD 4 | 0.29 ± 0.00 0.16 ± 0.00 | 0.43 ± 0.00 0.24 ± 0.00 | 0.31 ± 0.00 0.16 ± 0.00 | 0.40 ± 0.00 0.25 ± 0.00 | 0.38 ± 0.00 0.23 ± 0.00 | 0.31 ± 0.00 0.17 ± 0.00 | 0.37 ± 0.00 0.22 ± 0.00 | 0.31 ± 0.00 0.16 ± 0.00 | 0.42 ± 0.00 0.25 ± 0.00 | 0.30 ± 0.00 0.16 ± 0.00 | 0.35 ± 0.05 0.20 ± 0.04 |
| | 37D | SVD 8 | 0.06 ± 0.00 | 0.09 ± 0.00 | 0.06 ± 0.00 | 0.11 ± 0.00 | 0.10 ± 0.00 | 0.07 ± 0.00 | 0.09 ± 0.00 | 0.06 ± 0.00 | 0.11 ± 0.00 | 0.06 ± 0.00 | 0.08 ± 0.02 |
| 72 | | 16 D 32 D | 0.37 ± 0.02 0.21 ± 0.01 | 0.51 ± 0.02 0.28 ± 0.00 | 0.36 ± 0.01 0.20 ± 0.01 | 0.57 ± 0.02 0.35 ± 0.00 | 0.55 ± 0.00 0.33 ± 0.01 | 0.39 ± 0.02 0.22 ± 0.01 | 0.49 ± 0.01 0.31 ± 0.01 | 0.36 ± 0.02 0.20 ± 0.01 | 0.53 ± 0.03 0.32 ± 0.01 | 0.39 ± 0.01 0.22 ± 0.00 | 0.45 ± 0.08 0.26 ± 0.06 |
| 573 | 10 diagonal (D) | 64 D | 0.10 ± 0.00 | 0.11 ± 0.01 | 0.09 ± 0.00 | 0.18 ± 0.00 | 0.18 ± 0.00 | 0.10 ± 0.00 | 0.15 ± 0.01 | 0.09 ± 0.00 | 0.14 ± 0.00 | 0.09 ± 0.00 | 0.12 ± 0.04 0.12 ± 0.04 |
| 74 | | 128 D 256 D | 0.02 ± 0.00 0.00 ± 0.00 | $0.01 \pm 0.00 \\ 0.00 \pm 0.00$ | 0.02 ± 0.00 0.00 ± 0.00 | $0.03 \pm 0.00 \\ 0.00 \pm 0.00$ | 0.04 ± 0.00 0.00 ± 0.00 | 0.02 ± 0.00 0.00 ± 0.00 | $0.03 \pm 0.00 \\ 0.00 \pm 0.00$ | $0.02 \pm 0.00 \\ 0.00 \pm 0.00$ | $0.02 \pm 0.00 \\ 0.00 \pm 0.00$ | $0.02 \pm 0.00 \\ 0.00 \pm 0.00$ | $0.03 \pm 0.01 \\ 0.00 \pm 0.00$ |
| 75 | | 16 F | 0.35 ± 0.00 | 0.46 ± 0.00 | 0.34 ± 0.00 | 0.51 ± 0.00 | 0.47 ± 0.01 | 0.36 ± 0.01 | 0.45 ± 0.01 | 0.35 ± 0.01 | 0.49 ± 0.00 | 0.35 ± 0.01 | 0.41 ± 0.06 |
| 76 | 10 full (F) | 64 F | 0.20 ± 0.00 0.10 ± 0.00 | 0.24 ± 0.00 0.10 ± 0.00 | 0.20 ± 0.00 0.09 ± 0.00 | 0.30 ± 0.00 0.13 ± 0.00 | 0.29 ± 0.00 0.13 ± 0.00 | 0.22 ± 0.00 0.10 ± 0.00 | 0.27 ± 0.00 0.12 ± 0.00 | 0.20 ± 0.00 0.09 ± 0.00 | 0.27 ± 0.00 0.12 ± 0.00 | 0.21 ± 0.00 0.10 ± 0.00 | 0.24 ± 0.04 0.11 ± 0.02 |
| 70 | | 128 F 256 F | 0.02 ± 0.00 0.00 ± 0.00 | 0.02 ± 0.00 0.00 ± 0.00 | 0.02 ± 0.00 0.00 ± 0.00 | 0.01 ± 0.00 0.00 ± 0.00 | 0.02 ± 0.00 0.00 ± 0.00 | 0.02 ± 0.00 0.00 ± 0.00 | 0.02 ± 0.00 0.00 ± 0.00 | 0.02 ± 0.00 0.00 ± 0.00 | 0.01 ± 0.00 0.00 ± 0.00 | 0.02 ± 0.00 0.00 ± 0.00 | 0.02 ± 0.00 0.00 ± 0.00 |
| | | 16 D | 0.66 ± 0.01 | 0.69 ± 0.01 | 0.88 ± 0.01 | 0.76 ± 0.03 | 0.95 ± 0.02 | 0.91 ± 0.01 | 0.83 ± 0.02 | 0.88 ± 0.03 | 0.72 ± 0.02 | 0.88 ± 0.02 | 0.82 ± 0.10 |
| 78 | 50 diagonal (D) | 32 D 64 D | 0.50 ± 0.01 0.34 ± 0.01 | 0.52 ± 0.02 0.37 ± 0.01 | 0.73 ± 0.01 0.52 ± 0.00 | 0.58 ± 0.03 0.38 ± 0.01 | 0.88 ± 0.03 0.71 ± 0.02 | 0.79 ± 0.03 0.58 ± 0.01 | 0.72 ± 0.01 0.54 ± 0.00 | 0.75 ± 0.01 0.56 ± 0.00 | 0.57 ± 0.02 0.44 ± 0.01 | 0.75 ± 0.01 0.58 ± 0.01 | 0.68 ± 0.12 0.50 ± 0.11 |
| 79 | | 128 D | 0.21 ± 0.01 | 0.22 ± 0.01 | 0.31 ± 0.00 | 0.22 ± 0.00 | 0.51 ± 0.01 | 0.42 ± 0.01 | 0.38 ± 0.00 | 0.39 ± 0.00 | 0.27 ± 0.00 | 0.40 ± 0.00 | 0.33 ± 0.10 |
| 80 - | | 230 D | 0.10 ± 0.00 0.57 ± 0.01 | 0.12 ± 0.00 0.60 ± 0.01 | 0.10 ± 0.00 0.86 ± 0.01 | 0.10 ± 0.00 0.71 ± 0.02 | 0.29 ± 0.01 0.95 ± 0.01 | 0.21 ± 0.00 0.88 + 0.01 | 0.19 ± 0.00 0.81 ± 0.00 | 0.23 ± 0.01 0.83 ± 0.01 | 0.13 ± 0.00 0.67 ± 0.01 | 0.20 ± 0.00 0.86 ± 0.01 | 0.18 ± 0.06 0.78 ± 0.12 |
| 581 | 50 6-11 (E) | 32 F | 0.47 ± 0.01 | 0.48 ± 0.01 | 0.71 ± 0.00 | 0.55 ± 0.01 | 0.78 ± 0.01 | 0.69 ± 0.01 | 0.69 ± 0.00 | 0.65 ± 0.01 | 0.53 ± 0.01 | 0.71 ± 0.00 | 0.63 ± 0.11 |
| 200 | 50 Iuli (F) | 128 F | $0.33 \pm 0.00 \\ 0.19 \pm 0.00$ | 0.35 ± 0.00 0.21 ± 0.00 | 0.45 ± 0.00 0.25 ± 0.00 | 0.36 ± 0.00 0.19 ± 0.00 | 0.36 ± 0.00 0.35 ± 0.00 | 0.30 ± 0.01 0.30 ± 0.00 | 0.47 ± 0.00 0.28 ± 0.00 | 0.49 ± 0.00 0.31 ± 0.00 | 0.39 ± 0.01 0.24 ± 0.00 | 0.49 ± 0.00 0.30 ± 0.00 | 0.44 ± 0.08 0.26 ± 0.05 |
| 02 00 ⁼ | | 256 F | 0.09 ± 0.00 | 0.10 ± 0.00 | 0.10 ± 0.00 | 0.08 ± 0.00 | 0.16 ± 0.00 | 0.13 ± 0.00 | 0.12 ± 0.00 | 0.15 ± 0.00 | 0.11 ± 0.00 | 0.13 ± 0.00 | 0.12 ± 0.02 |
| 83 | | 16 D 32 D | 0.90 ± 0.01 0.83 ± 0.02 | 0.85 ± 0.01 0.77 ± 0.00 | 0.87 ± 0.03 0.77 ± 0.01 | 0.88 ± 0.02 0.78 ± 0.00 | 0.68 ± 0.02 0.55 ± 0.02 | 0.91 ± 0.01 0.79 ± 0.01 | 0.97 ± 0.01 0.94 + 0.02 | 0.98 ± 0.01 0.94 ± 0.03 | 0.96 ± 0.01 0.87 ± 0.00 | 1.00 ± 0.00 0.98 + 0.01 | 0.90 ± 0.09 0.82 ± 0.12 |
| 84 | 100 diagonal (D) | 64 D | 0.67 ± 0.01 | 0.63 ± 0.00 | 0.59 ± 0.02 | 0.63 ± 0.01 | 0.40 ± 0.02 | 0.62 ± 0.01 | 0.86 ± 0.02 | 0.82 ± 0.02 | 0.71 ± 0.03 | 0.93 ± 0.00 | 0.68 ± 0.15 |
| 85 | | 128 D 256 D | 0.49 ± 0.01 0.32 ± 0.00 | 0.47 ± 0.00 0.31 ± 0.00 | 0.42 ± 0.01 0.26 ± 0.01 | $0.45 \pm 0.00 \\ 0.30 \pm 0.00$ | 0.27 ± 0.02 0.15 ± 0.01 | 0.44 ± 0.01 0.28 ± 0.00 | 0.73 ± 0.01 0.51 ± 0.02 | 0.69 ± 0.02 0.51 ± 0.02 | $0.59 \pm 0.02 \\ 0.40 \pm 0.01$ | 0.80 ± 0.02 0.61 ± 0.01 | $0.53 \pm 0.16 \\ 0.36 \pm 0.14$ |
| 86 | | 16 F | 0.88 ± 0.00 | 0.82 ± 0.00 | 0.84 ± 0.01 | 0.86 ± 0.00 | 0.67 ± 0.01 | 0.88 ± 0.01 | 0.99 ± 0.00 | 0.96 ± 0.01 | 0.91 ± 0.01 | 1.00 ± 0.00 | 0.88 ± 0.09 |
| 87 | 100 full (F) | 52 F 64 F | 0.78 ± 0.00 0.60 ± 0.00 | 0.72 ± 0.00 0.57 ± 0.00 | 0.73 ± 0.00 0.57 ± 0.00 | 0.74 ± 0.00 0.57 ± 0.00 | 0.32 ± 0.00 0.39 ± 0.00 | 0.74 ± 0.01 0.56 ± 0.00 | 0.94 ± 0.01 0.76 ± 0.00 | 0.89 ± 0.00 0.73 ± 0.00 | 0.77 ± 0.02 0.60 ± 0.00 | 0.99 ± 0.00 0.83 ± 0.01 | 0.78 ± 0.13 0.62 ± 0.12 |
| | | 128 F 256 F | 0.40 ± 0.00 0.21 ± 0.00 | 0.38 ± 0.00 0.20 ± 0.00 | 0.35 ± 0.00 0.18 ± 0.00 | 0.37 ± 0.00 0.19 ± 0.00 | 0.25 ± 0.00 0.13 ± 0.00 | 0.37 ± 0.00 0.19 ± 0.00 | 0.52 ± 0.00 0.30 ± 0.00 | 0.54 ± 0.00 0.34 ± 0.00 | 0.45 ± 0.00 0.26 ± 0.00 | 0.60 ± 0.00 0.38 ± 0.00 | 0.42 ± 0.10 0.24 ± 0.08 |
| ö _ | 100 w/clusters (C) | 16 C 5 | 0.46 | 0.46 | 0.44 | 0.47 | 0.62 | 0.64 | 0.61 | 0.63 | 0.45 | 0.60 | 0.54 |
| 9 - | Too menusiens (C) | 16 C 7 | 0.42 | 0.43 | 0.40 | 0.42 | 0.50 | 0.55 | 0.52 | 0.55 | 0.42 | 0.55 | 0.48 |
| 0 | | 16 D 32 D | 0.97 ± 0.00 0.96 ± 0.00 | 0.73 ± 0.00 0.70 ± 0.00 | 0.96 ± 0.00 0.92 ± 0.01 | 1.00 ± 0.00 0.98 ± 0.01 | 0.99 ± 0.01 0.96 ± 0.01 | 0.96 ± 0.01 0.93 ± 0.01 | 0.90 ± 0.00 0.86 ± 0.00 | 0.92 ± 0.00 0.89 ± 0.00 | 1.00 ± 0.00 1.00 ± 0.00 | 1.00 ± 0.00 1.00 ± 0.00 | 0.94 ± 0.08 0.92 ± 0.09 |
| 91 | 500 diagonal (D) | 64 D | 0.90 ± 0.01 | 0.65 ± 0.00 | 0.86 ± 0.01 | 0.96 ± 0.02 | 0.90 ± 0.01 | 0.87 ± 0.01 | 0.81 ± 0.00 | 0.83 ± 0.01 | 0.99 ± 0.01 | 1.00 ± 0.00 | 0.88 ± 0.10 |
| 92 | | 256 D | 0.82 ± 0.01 0.59 ± 0.02 | 0.60 ± 0.00 0.51 ± 0.00 | 0.76 ± 0.00 0.56 ± 0.01 | 0.90 ± 0.02 0.81 ± 0.02 | 0.83 ± 0.01 0.70 ± 0.02 | 0.78 ± 0.02 0.55 ± 0.01 | 0.74 ± 0.00 0.57 ± 0.01 | 0.74 ± 0.01 0.54 ± 0.01 | 0.97 ± 0.01 0.91 ± 0.01 | 1.00 ± 0.00 1.00 ± 0.01 | 0.81 ± 0.12 0.67 ± 0.17 |
| 93 | | 16 F 32 F | 0.94 ± 0.00 | 0.67 ± 0.00 | 0.88 ± 0.00 | 1.00 ± 0.00 | 0.98 ± 0.00 | 0.90 ± 0.00 | 0.82 ± 0.00 | 0.83 ± 0.00 | 1.00 ± 0.00 | 1.00 ± 0.00 | 0.90 ± 0.10 |
| 24 | 500 full (F) | 64 F | 0.80 ± 0.00 0.80 ± 0.00 | 0.51 ± 0.00 0.55 ± 0.00 | 0.01 ± 0.00 0.72 ± 0.00 | 0.97 ± 0.01 0.86 ± 0.00 | 0.94 ± 0.00 0.82 ± 0.01 | 0.76 ± 0.00 | 0.75 ± 0.00 0.67 ± 0.00 | 0.70 ± 0.00 | 0.99 ± 0.00 0.94 ± 0.00 | 0.99 ± 0.00 0.99 ± 0.00 | 0.30 ± 0.12 0.78 ± 0.13 |
| 94 | | 128 F 256 F | $\begin{array}{c} 0.64 \pm 0.00 \\ 0.43 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.46 \pm 0.00 \\ 0.35 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.60 \pm 0.00 \\ 0.44 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.74 \pm 0.00 \\ 0.55 \pm 0.00 \end{array}$ | $0.65 \pm 0.00 \\ 0.49 \pm 0.00$ | $\begin{array}{c} 0.63 \pm 0.00 \\ 0.45 \pm 0.00 \end{array}$ | $0.56 \pm 0.00 \\ 0.40 \pm 0.00$ | $\begin{array}{c} 0.58 \pm 0.00 \\ 0.42 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.85 \pm 0.00 \\ 0.67 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.96 \pm 0.00 \\ 0.84 \pm 0.00 \end{array}$ | 0.67 ± 0.14 0.50 ± 0.14 |
| 95 - | | 16 C 7 | 0.68 | 0.70 | 0.64 | 0.72 | 0.85 | 0.90 | 0.93 | 0.92 | 0.71 | 0.83 | 0.79 |
| 96 | 500 w/clusters (C) | 16 C 10 16 C 25 | 0.61 0.42 | 0.65 0.41 | 0.61 0.42 | 0.66 0.44 | 0.84 0.57 | 0.86 0.64 | 0.88 0.63 | 0.84 0.62 | 0.62 0.40 | 0.76 0.58 | 0.73 0.51 |
| 97 | | 64 C 5 64 C 7 | 0.49 | 0.49 | 0.45 | 0.51 | 0.64 | 0.66 | 0.62 | 0.67 | 0.50 | 0.65 | 0.57 |
| | | 0.07 | 0.10 | 0.15 | 0.11 | 0.15 | 0.50 | 0.00 | 0.00 | 0.07 | | 1 0.07 | 0.01 |

Table 11: Reconstruction error In-Distribution for various tasks and methods

| 1605 | | | | | | | | | | | | | |
|------|--------------|--------------|---|---|---|---|---|---|---|---|---|---|---|
| 1606 | Model Type | Method Type | | | | | Ta | sks | | | | | Average |
| 1607 | | | task039 | task190 | task280 | task290 | task391 | task442 | task620 | task1342 | task1391 | task1598 | |
| 1007 | | 16 F | 0.46 ± 0.01 | 0.63 ± 0.00 | 0.50 ± 0.00 | 0.55 ± 0.01 | 0.50 ± 0.00 | 0.49 ± 0.01 | 0.50 ± 0.01 | 0.50 ± 0.01 | 0.61 ± 0.01 | 0.47 ± 0.01 | 0.52 ± 0.06 |
| 1608 | 10 full (F) | 32 F 64 F | 0.30 ± 0.01 0.15 ± 0.00 | $\begin{array}{c} 0.37 \pm 0.00 \\ 0.15 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.31 \pm 0.00 \\ 0.16 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.35 \pm 0.00 \\ 0.17 \pm 0.00 \end{array}$ | 0.34 ± 0.00 0.17 ± 0.00 | $\begin{array}{c} 0.31 \pm 0.00 \\ 0.16 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.33 \pm 0.00 \\ 0.16 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.31 \pm 0.00 \\ 0.16 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.38 \pm 0.00 \\ 0.17 \pm 0.00 \end{array}$ | $0.30 \pm 0.00 \\ 0.15 \pm 0.00$ | $\begin{array}{c} 0.33 \pm 0.03 \\ 0.16 \pm 0.01 \end{array}$ |
| 1609 | | 16 F | 0.80 ± 0.02 | 0.82 ± 0.01 | 0.85 ± 0.01 | 0.90 ± 0.02 | 0.78 ± 0.01 | 0.95 ± 0.01 | 0.76 ± 0.01 | 0.75 ± 0.00 | 0.79 ± 0.01 | 0.82 ± 0.01 | 0.82 ± 0.06 |
| 1610 | 50 full (F) | 32 F 64 F | $\begin{array}{c} 0.65 \pm 0.01 \\ 0.50 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.67 \pm 0.01 \\ 0.52 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.72 \pm 0.01 \\ 0.52 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.76 \pm 0.02 \\ 0.55 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.65 \pm 0.01 \\ 0.52 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.82 \pm 0.02 \\ 0.62 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.66 \pm 0.01 \\ 0.54 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.65 \pm 0.01 \\ 0.51 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.67 \pm 0.02 \\ 0.54 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.69 \pm 0.00 \\ 0.57 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.69 \pm 0.06 \\ 0.54 \pm 0.03 \end{array}$ |
| 1611 | | 16 F | 0.93 ± 0.02 | 0.90 ± 0.02 | 0.93 ± 0.01 | 0.91 ± 0.02 | 0.88 ± 0.03 | 0.98 ± 0.01 | 0.96 ± 0.01 | 0.78 ± 0.00 | 0.82 ± 0.00 | 0.93 ± 0.02 | 0.90 ± 0.06 |
| 1612 | 100 full (F) | 32 F 64 F | $\begin{array}{c} 0.87 \pm 0.01 \\ 0.65 \pm 0.04 \end{array}$ | $\begin{array}{c} 0.81 \pm 0.01 \\ 0.69 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.85 \pm 0.02 \\ 0.71 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.80 \pm 0.01 \\ 0.67 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.79 \pm 0.02 \\ 0.64 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.91 \pm 0.00 \\ 0.76 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.90 \pm 0.01 \\ 0.77 \pm 0.01 \end{array}$ | $\begin{array}{c} 0.74 \pm 0.01 \\ 0.67 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.70 \pm 0.02 \\ 0.61 \pm 0.00 \end{array}$ | $\begin{array}{c} 0.85 \pm 0.02 \\ 0.75 \pm 0.06 \end{array}$ | $\begin{array}{c} 0.82 \pm 0.07 \\ 0.69 \pm 0.06 \end{array}$ |
| 1613 | | 16 F | 0.98 ± 0.04 | 0.98 ± 0.01 | 0.99 ± 0.01 | 1.00 ± 0.00 | 0.99 ± 0.00 | 0.96 ± 0.05 | 0.93 ± 0.10 | 0.94 ± 0.09 | 1.00 ± 0.00 | 0.99 ± 0.00 | 0.98 ± 0.05 |
| 1614 | 500 full (F) | 52 F 64 F | 0.92 ± 0.07 0.80 ± 0.00 | 0.84 ± 0.20 0.67 ± 0.21 | 0.92 ± 0.10 0.78 ± 0.11 | 0.98 ± 0.02 0.90 ± 0.07 | 0.97 ± 0.02 0.86 ± 0.08 | 0.89 ± 0.08 0.76 ± 0.00 | $\begin{array}{c} 0.82 \pm 0.13 \\ 0.67 \pm 0.00 \end{array}$ | 0.84 ± 0.11 0.70 ± 0.00 | $\begin{array}{c} 0.99 \pm 0.00 \\ 0.96 \pm 0.03 \end{array}$ | 0.99 ± 0.02 0.99 ± 0.00 | 0.92 ± 0.10 0.81 ± 0.13 |
| 1615 | | | | | | | | | | | | | |

1616Table 12: Reconstruction error on random LoRAs The error is larger in comparison to recon-1617structing trained (i.e., non-random) LoRAs in Table 11 for the corresponding compression methods.

1632 1633

1634 1635

1649 1650

1620 Tasks 1621 Model Type Method Type Average task620 | task1342 | task1391 | task1598 task039 task190 task280 | task290 | task391 task442 1622 24.44 1.60 19.13 39.22 10.27 35.46 7.85 6.22 17.82 38.87 20.09 base 1623 78.93 lora 95.00 86.00 99.00 93.67 94.33 74.88 74.40 26.68 95.00 50.32 1624 94.00 96.00 97.00 90.00 93.33 94.67 94.67 74.09 74.29 72.13 74.80 27.83 26.63 50.71 51.04 79.28 32 F 99.00 79.41 10 full (F) 64 F 95.00 89.00 99.00 93.67 1625 92.33 32 F 96.00 88.00 99.00 93.67 72.30 75.97 29.89 94.00 45.68 78.68 1626 50 full (F) 64 F 98.00 89.00 99.00 93.67 93.33 72.74 76.50 29.33 96.00 45.71 79.33 1627 42.36 44.71 32 F 92.10 83.00 99.00 93.67 92.00 71.09 63.29 27.87 88.00 75.24 97.00 92.33 29.98 95.00 100 full (F) 64 F 99.00 74.69 1628 87.00 93.67 72.23 78.56 68.92 93.50 43.00 78.00 87.00 91.00 91.67 92.33 90.67 90.33 70.08 72.55 51.16 57.49 14.40 15.44 83.00 85.00 32 F 41.97 64.19 1629 64 F 42.31 500 full (F) 71.80 1630

Table 13: Performance with convergence In-Distribution Rouge-L

Table 14: Agreement Comparison. 100 LoRAs

| Configuration | | Agreement (%) |
|--------------------|---------|---------------|
| Base Model | | 83.015 |
| Uncompressed LoRAs | | 100.000 |
| Joint Compression | | |
| Diagonal | Rank 8 | 87.032 |
| 2 | Rank 16 | 88.908 |
| | Rank 32 | 91.545 |
| | Rank 64 | 94.659 |
| Full | Rank 8 | 87.686 |
| | Rank 16 | 90.163 |
| | Rank 32 | 94.018 |
| | Rank 64 | 96.918 |

Table 15: Performance Comparison. 100 LoRAs

| Configuration | | Average Performance |
|--------------------|---------|---------------------|
| Base Model | | 32.28 |
| Uncompressed LoRAs | | 48.32 |
| Join Compression | | |
| Diagonal | Rank 8 | 41.90 |
| C | Rank 16 | 45.44 |
| | Rank 32 | 46.89 |
| | Rank 64 | 47.43 |
| Full | Rank 8 | 43.88 |
| | Rank 16 | 45.79 |
| | Rank 32 | 46.83 |
| | Rank 64 | 47.66 |

Table 16: Task-Based Performance Evaluation Across Different Models and Ranks

| 1663 | | | | | | | |
|------|-------------------------------|------------|-------|-------------|--------------|--------------|--------------|
| 1664 | Task | Base Model | LoRA | Diagonal R8 | Diagonal R16 | Diagonal R32 | Diagonal R64 |
| 1665 | Causal Judgement | 57.47 | 64.37 | 55.17 | 58.62 | 58.62 | 58.62 |
| 1666 | Date Understanding | 15.33 | 23.33 | 20.67 | 22.00 | 21.33 | 22.67 |
| 1000 | Formal Fallacies | 51.33 | 56.00 | 52.67 | 52.67 | 53.33 | 54.67 |
| 1667 | Hyperbaton | 6.67 | 68.00 | 57.33 | 63.33 | 67.33 | 68.00 |
| 1668 | Logical Deduction (5 Objects) | 21.33 | 37.33 | 32.00 | 36.67 | 37.33 | 37.33 |
| 1660 | Logical Deduction (7 Objects) | 12.67 | 44.00 | 31.33 | 42.67 | 44.67 | 45.33 |
| 1005 | Movie Recommendation | 62.67 | 67.33 | 62.00 | 64.67 | 66.67 | 67.33 |
| 1670 | Object Counting | 34.67 | 38.00 | 35.33 | 36.67 | 36.67 | 38.00 |
| 1671 | Snarks | 50.00 | 61.54 | 53.85 | 56.41 | 58.97 | 57.69 |
| 1672 | Temporal Sequences | 16.67 | 23.33 | 18.67 | 20.67 | 24.00 | 24.67 |
| 1673 | Average | 32.88 | 48.32 | 41.90 | 45.44 | 46.89 | 47.43 |

Table 17: Task-Based Performance Evaluation Across Different Models and Ranks

| Task | Base Model | LoRA | Full R8 | Full R16 | Full R32 | Full R64 |
|-------------------------------|------------|-------|---------|----------|----------|----------|
| Causal Judgement | 57.47 | 64.37 | 56.32 | 57.47 | 58.62 | 60.92 |
| Date Understanding | 15.33 | 23.33 | 19.33 | 22.00 | 22.67 | 22.67 |
| Formal Fallacies | 51.33 | 56.00 | 51.33 | 52.67 | 53.33 | 56.00 |
| Hyperbaton | 6.67 | 68.00 | 63.33 | 66.00 | 69.33 | 68.00 |
| Logical Deduction (5 Objects) | 21.33 | 37.33 | 35.33 | 36.00 | 35.33 | 37.33 |
| Logical Deduction (7 Objects) | 12.67 | 44.00 | 40.00 | 44.67 | 44.67 | 44.67 |
| Movie Recommendation | 62.67 | 67.33 | 63.33 | 65.33 | 67.33 | 67.33 |
| Object Counting | 34.67 | 38.00 | 35.33 | 36.67 | 37.33 | 37.33 |
| Snarks | 50.00 | 61.54 | 53.85 | 55.13 | 57.69 | 58.97 |
| Temporal Sequences | 16.67 | 23.33 | 20.67 | 22.00 | 22.00 | 23.33 |
| Average | 32.88 | 48.32 | 43.88 | 45.79 | 46.83 | 47.66 |