

# AFORA: ACTIVATION-AWARE FACTORIZATION WITH OPTIMAL RANK ALLOCATION FOR TRAINING-FREE LLM COMPRESSION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Large language models are challenging to deploy because of their extreme size and compute demands. In this work, we propose *AFORA* (Activation-aware Factorization with Optimal Rank Allocation for Training-free LLM Compression), a simple and hardware-friendly framework that directly reduces the number of parameters through low-rank factorization of weight matrices. *AFORA* consists of two core components: (1) **Activation-aware Weight Factorization (AWF)**, a closed-form low-rank approximation that accounts for the input activation distribution to preserve task-relevant directions and ensure numerical stability; and (2) **Optimal Rank Allocation (ORA)**, a global rank allocation strategy that assigns heterogeneous ranks across layers to minimize activation distortion under a given budget. Evaluations across multiple large-scale language models show that our framework consistently outperforms existing approaches at the same compression ratios, while additionally reducing model size, saving memory, and decreasing computation with hardware-friendly layer dimensions. It also requires only a short runtime to perform compression, and offers a principled mathematical interpretation. These results demonstrate that activation-aware, globally optimized low-rank compression offers a practical and theoretically grounded path to efficient LLM deployment.

## 1 INTRODUCTION

Large language models (LLMs) have achieved remarkable success (Achiam et al., 2023; Brown et al., 2020; Touvron et al., 2023), but their rapidly growing parameter counts pose severe challenges for memory, compute, and deployment efficiency. To mitigate this, many efforts have attempted to reduce the number of parameters through various compression techniques. Among existing approaches, **quantization** reduces precision to save memory. However, even after quantization, model execution typically requires dequantization, meaning that inference-time acceleration gains are minimal while the primary benefit lies in storage reduction. To fully exploit quantization for runtime acceleration, one must rely on specialized hardware kernels, which are not always available or portable across devices (Frantar et al., 2022; Xiao et al., 2023; Lin et al., 2024; Yao et al., 2022). **Pruning** reduces parameter counts by removing weights or blocks, yet width pruning produces sparse matrices that rarely translate to proportional speedups and still require full memory loads. Depth pruning has been proposed as a more structured variant that removes entire blocks without changing the overall architecture. However, naively discarding less important blocks can overly simplify the model and compromise robustness (Han et al., 2015; Hoefler et al., 2021; Song et al., 2024; Dettmers et al., 2023).

These observations highlight the need for an alternative that reduces parameter counts in a hardware-friendly manner, avoids architectural modifications, and still preserves fine-grained modeling capacity. In other words, instead of simply eliminating less important components, we argue for mathematically principled methods that replace computations with efficient low-rank structures, thereby reducing operations while approximating the original behavior to minimize performance loss.

**We propose *AFORA* (Activation-aware Factorization with Optimal Rank Allocation for Training-free LLM Compression)**, a simple yet principled framework that directly reduces pa-

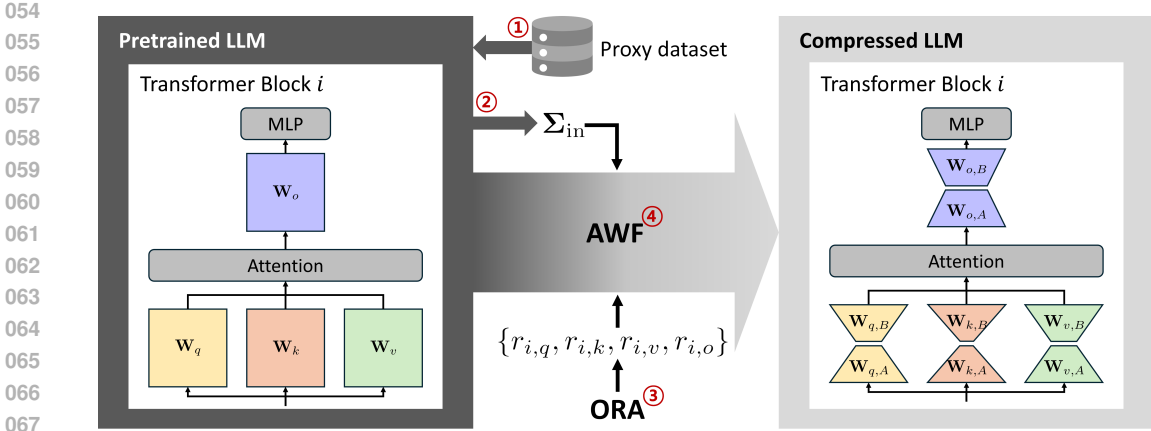


Figure 1: Overview of AFORA

rameters through low-rank factorization of less important weight matrices. Unlike quantization and pruning, *AFORA* is training-free, hardware-efficient, and supported by a theoretically grounded analysis that accounts for activation geometry. In practice, it also keeps the compression stage lightweight and leads to inference-time reductions that are consistent with the nominal FLOPs savings.

Our method is built upon two core developments: (i) **Activation-aware Weight Factorization (AWF)**, a closed-form low-rank approximation guided by the input activation distribution, preserving task-relevant directions and ensuring numerical stability. (ii) **Optimal Rank Allocation (ORA)**, a global allocation strategy that assigns rank across layers optimally under a given budget.

**Our contributions are threefold:**

- We propose a new weight factorization method, referred to as **AWF**. We provide a rigorous justification from both a statistical perspective and an optimization viewpoint, showing that it minimizes reconstruction loss in a principled manner.
- We develop the optimal rank allocation strategy that assigns heterogeneous ranks across layers. We theoretically prove that it yields the mathematically optimal rank assignment under a global compression budget.
- We demonstrate the effectiveness of the proposed **AFORA** framework through extensive experiments on LLaMA-2-7B, showing that at a compression ratio of 0.15, **AFORA** reduces perplexity on WikiText-2 by 10.27% compared to the best existing method. It also reduces compression time by up to 9.60× relative to the same baseline, achieves an inference speedup of 1.21×, and matches the zero-shot accuracy of the dense model.

Together, these results show that activation-aware, globally optimized low-rank distribution provides a practical and theoretically grounded path to efficient LLM deployment.

## 2 BACKGROUND AND RELATED WORK

### 2.1 LOW-RANK APPROXIMATION AS A COMPRESSION TOOL

Low-rank approximation has long been used to reduce the complexity of large linear operators. Given a weight matrix  $W \in \mathbb{R}^{d_{out} \times d_{in}}$ , it can be factorized as  $W \approx BA$  with  $B \in \mathbb{R}^{d_{out} \times r}$  and  $A \in \mathbb{R}^{r \times d_{in}}$ , where  $r \ll \min(d_{out}, d_{in})$ . This reduces both parameter count and FLOPs approximately as  $d_{out}d_{in} \rightarrow r(d_{out} + d_{in})$ . Because Transformers are dominated by large linear layers (attention projections and feed-forward blocks) Vaswani et al. (2017), low-rank approximation offers a hardware-friendly pathway to reduce model size and compute (FLOPs). The most common variant is truncated singular value decomposition (SVD), which selects the top- $r$  singular components of  $W$  under the Frobenius norm. By the **Eckart–Young–Mirsky theorem**, this truncated

SVD is the unique optimal rank- $r$  approximation in weight space (Eckart & Young, 1936; Mirsky, 1960). A formal statement and proof sketch are provided in Appendix A.

**Limitations of naive SVD.** Truncated SVD is simple and training-free but *data-agnostic*: it minimizes the reconstruction error of  $\mathbf{W}$ , not the induced error on outputs. Directions that seem unimportant in  $\mathbf{W}$  may be heavily amplified by input activations, leading to large output errors.

## 2.2 ACTIVATION-AWARE REDUCTION: VIEWPOINT AND PRIOR ATTEMPTS

**Activation-aware viewpoint.** Compression quality should be measured in the output space, under the input activation distribution. This motivates a whitened operator formulation (formalized later in Sec. 3). Concretely, the objective is

$$\min \|\mathbf{W}\mathbf{x} - \mathbf{B}\mathbf{A}\mathbf{x}\|_F^2, \quad (1)$$

which emphasizes minimizing the output error induced by the low-rank factorization.

**Prior heuristics.** Several works attempt to improve upon naive SVD by incorporating task-related statistics. **FWSVD** (Hsu et al., 2022) modifies the spectrum using Fisher information so that directions aligned with high curvature are preserved. **ASVD** (Yuan et al., 2023) instead rescales components according to the empirical activation distribution. Beyond such modifications of the spectrum, **LoRAP** (Li et al., 2024) highlights the heterogeneity across submodules: attention projections can tolerate low-rank designs better, while MLP layers are more sensitive to compression. **SVD-LLM** (Wang et al., 2024) further shows that even activation-aware variants still suffer significant performance loss, which must be compensated by parameter update. These findings suggest that while heuristics can partially mitigate the shortcomings of naive SVD, they are still largely heuristic in nature and lack a unified theoretical foundation. This indicates room for deeper mathematical analysis of activation-aware low-rank approximations—a direction we pursue in this work.

## 2.3 OTHER COMPRESSION PARADIGMS (TRAINING-FREE, HARDWARE-FRIENDLY)

We now contrast our approach with other paradigms that, like ours, can be applied **without training**. Our focus is on post-training compression methods that can reduce parameters or computations without requiring additional fine-tuning.

**Quantization.** Quantization reduces model size by lowering numerical precision and is especially attractive in post-training settings. Different designs vary in whether quantization is applied uniformly or in a group-wise manner, and whether activations are calibrated to control variance. At moderate bit-widths (e.g., 8–4 bits), quantization can achieve strong compression with minimal accuracy loss. However, quantized inference typically requires dequantization before matrix multiplication, which means that real efficiency gains often rely on specialized kernels. This introduces hardware-dependence and limits portability. For [ultra-low precision \(e.g., 2–3 bits\)](#), the [theoretical gains often depend on specialized kernels, and practical speedups are not consistently realized across hardware platforms \(Frantar et al., 2022; Xiao et al., 2023; Lin et al., 2024; Yao et al., 2022; Dettmers et al., 2023; Egiazarian et al., 2025\)](#).

**Pruning.** Pruning has been an important line of research for compressing LLMs, as it can reduce parameter redundancy and lower computational cost. There are two main forms: **unstructured pruning**, which removes individual weights and can reach very high sparsity levels (often above 90%), but struggles to deliver real acceleration due to irregular memory access and the need to fully load parameters into memory; and **structured pruning**, which removes predefined groups (e.g., rows, columns, or 2:4 patterns) to form hardware-friendly patterns. While structured pruning is more practical for deployment, its actual speedup is often much smaller than the theoretical pruning ratio, especially at small batch sizes, and preserving linguistic performance typically requires substantial retraining. Recent studies also combine pruning with low-rank approximations, suggesting that pruning can complement other compression techniques. Overall, pruning remains a valuable tool, but its limitations in robustness and hardware efficiency highlight the need for simpler and more principled alternatives beyond pruning (Sanh et al., 2020; Michel et al., 2019; Lagunas et al., 2021; Song et al., 2024; Wee et al., 2025).

**Low-rank factorization.** Low-rank approximation factorizes large weight matrices into smaller components, reducing both parameter counts and multiply-accumulate operations while retaining standard dense matrix multiplications. Its use in model compression dates back to early work on convolutional networks, and its effectiveness in language models is supported by evidence that their intrinsic dimensionality is much smaller than the parameter space. Classic truncated SVD provides a closed-form solution but is weight-space and data-agnostic. Recent refinements introduce activation statistics or focus on submodule specialization, alleviating some limitations but still lacking a unified theoretical foundation or global optimization strategy. In contrast, our method explicitly formulates compression in whitened activation space, yielding a mathematically justified loss-aware factorization, and combines this with optimal global rank allocation to distribute budgets across layers in a principled way (Hsu et al., 2022; Yuan et al., 2023; Li et al., 2024; Wang et al., 2024).

### 3 ACTIVATION-AWARE WEIGHT FACTORIZATION

#### 3.1 FORMULATION FROM ACTIVATION DISTRIBUTION

We begin with the activation-aware objective:

$$\min_{\text{rank}(BA) \leq r} \mathbb{E}_x \left[ \|Wx - BAx\|_F^2 \right], \quad (2)$$

which measures compression error directly in the output space under the input distribution. As shown in Appendix B, this is equivalent to

$$\min_{\text{rank}(BA) \leq r} \|W\Sigma_{\text{in}}^{1/2} - BA\Sigma_{\text{in}}^{1/2}\|_F^2, \quad (3)$$

where  $\Sigma_{\text{in}} = \mathbb{E}[xx^\top]$  is the input covariance. Letting  $\Sigma_{\text{in}} = Q\Lambda Q^\top$  with  $\Lambda \succeq 0$ , we compute  $\Sigma_{\text{in}}^{\pm 1/2} = Q\Lambda^{\pm 1/2}Q^\top$  and define the whitened operator

$$T = W\Sigma_{\text{in}}^{1/2} = U \text{diag}(s) V^\top. \quad (4)$$

By the Eckart–Young–Mirsky theorem, the optimal rank- $r$  approximation is the truncated SVD

$$T_r^* = U_r S_r V_r^\top, \quad B^* A^* = U_r S_r V_r^\top \Sigma_{\text{in}}^{-1/2}. \quad (5)$$

We call this solution the **Activation-aware Weight Factorization (AWF)**.

#### 3.2 LOSS-BASED DERIVATION

**Second-order expansion of the loss.** Consider a weight matrix  $\mathbf{W}$  and perturbation  $\Delta\mathbf{W}$ . Expanding the training loss around  $\mathbf{W}$  yields

$$\begin{aligned} \Delta\mathcal{L} &= \mathcal{L}(\mathbf{W} + \Delta\mathbf{W}) - \mathcal{L}(\mathbf{W}) \\ &\approx \langle \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}), \Delta\mathbf{W} \rangle + \frac{1}{2} \text{vec}(\Delta\mathbf{W})^\top \mathbf{H}_{\mathbf{W}} \text{vec}(\Delta\mathbf{W}), \end{aligned} \quad (6)$$

where  $\mathbf{H}_{\mathbf{W}} = \nabla_{\mathbf{W}}^2 \mathcal{L}(\mathbf{W})$ . For a converged (near-stationary) pretrained model, the *first-order term* is typically small, so the loss increase is well approximated by

$$\Delta\mathcal{L} \approx \frac{1}{2} \text{vec}(\Delta\mathbf{W})^\top \mathbf{H}_{\mathbf{W}} \text{vec}(\Delta\mathbf{W}). \quad (7)$$

**K-FAC approximation.** For a linear layer  $h = \mathbf{W}x$ , the Kronecker-factored curvature (K-FAC) approximation (Martens & Grosse, 2015; Grosse & Martens, 2016; Botev et al., 2017) gives

$$\mathbf{H}_{\mathbf{W}} \approx \Sigma_{\text{out}} \otimes \Sigma_{\text{in}}, \quad \Sigma_{\text{in}} = \mathbb{E}[xx^\top], \quad \Sigma_{\text{out}} = \mathbb{E}[gg^\top], \quad g = \nabla_h \ell(h). \quad (8)$$

Using the identity  $\text{vec}(MXN) = (N^\top \otimes M)\text{vec}(X)$ , the quadratic form in equation 7 becomes

$$\Delta\mathcal{L} \approx \frac{1}{2} \|\Sigma_{\text{out}}^{1/2} \Delta\mathbf{W} \Sigma_{\text{in}}^{1/2}\|_F^2. \quad (9)$$

**Practical choice  $\Sigma_{\text{out}} = \mathbf{I}$ .** Estimating  $\Sigma_{\text{out}}$  requires backpropagated gradients  $g$ , which is costly at LLM scale. In practice we set  $\Sigma_{\text{out}} = \mathbf{I}$  to keep the method training-free, which yields

$$\Delta\mathcal{L} \approx \frac{1}{2} \|\Delta\mathbf{W} \Sigma_{\text{in}}^{1/2}\|_F^2, \quad \text{i.e.,} \quad \min_{\text{rank}(BA) \leq r} \|\mathbf{W} \Sigma_{\text{in}}^{1/2} - BA \Sigma_{\text{in}}^{1/2}\|_F^2, \quad (10)$$

which is the activation-aware covariance objective used throughout.

### 3.3 INFORMATION-BALANCED FACTORIZATION

While AWF gives the optimal approximation in expectation, its factors  $B$  and  $A$  can in practice exhibit scale imbalance across directions. Such imbalance may amplify numerical sensitivity and reduce stability, especially when singular values or activation statistics are highly anisotropic. To mitigate this issue, we introduce **Information-Balanced Factorization**, which balances per-direction scales. We formulate the problem as

$$\min_{\alpha_i > 0} \max_i \left\{ \|B(\alpha)_{:,i}\|_2, \|A(\alpha)_{i,:}\|_2 \cdot \|(\Sigma_{\text{in}}^{1/2} V_r)_{:,i}\|_2 \right\}, \quad (11)$$

where  $B(\alpha) = U_r \text{diag}(\alpha)$  and  $A(\alpha) = \text{diag}(s_i/\alpha_i) V_r^\top \Sigma_{\text{in}}^{-1/2}$ . This criterion enforces balanced scaling across the left factor, the right factor, and the whitened input directions, preventing any single component from dominating.

The closed-form solution is

$$\alpha_i^* = \sqrt{s_i \cdot \|(\Sigma_{\text{in}}^{1/2} V_r)_{:,i}\|_2}, \quad i = 1, \dots, r, \quad (12)$$

which equalizes scales across factors by taking the geometric mean of the singular value and the whitened input norm, thereby improving conditioning. See Appendix C for a detailed derivation.

## 4 ORA: OPTIMAL RANK ALLOCATION

### 4.1 WHY GLOBAL RANK ALLOCATION?

AWF (Sec. 3) specifies the best rank- $r$  approximation *within* a layer once  $r$  is fixed. In practice, however, we must decide *how much* rank each layer should be allocated under a single memory or FLOPs budget. Naive strategies such as uniform ranks or hand-crafted schedules ignore two key observations from AWF: (i) whitened spectra differ widely across layers, inducing different marginal utilities (see Appendix E); (ii) the per-rank cost depends on the layer’s shape. This motivates a global allocator that explicitly trades off *value* against *cost* across layers, akin to classical knapsack formulations (Kellerer et al., 2004).

### 4.2 GLOBAL OPTIMIZATION OBJECTIVE

Let targets be indexed by  $m \in \{1, \dots, M\}$  with weights  $\mathbf{W}_m \in \mathbb{R}^{d_{\text{out},m} \times d_{\text{in},m}}$ . With AWF (using  $\Sigma_{\text{out}} = \mathbf{I}$ ), define the whitened operator

$$\mathbf{T}_m = \mathbf{W}_m \Sigma_{\text{in},m}^{1/2} = \mathbf{U}_m \text{diag}(\mathbf{s}_m) \mathbf{V}_m^\top, \quad \mathbf{s}_m = (s_{m,1} \geq s_{m,2} \geq \dots). \quad (13)$$

Retaining the top  $r_m$  directions yields utility  $F_m(r_m) = \sum_{i=1}^{r_m} v_{m,i}$  with per-direction utilities  $v_{m,i} = \frac{1}{2} s_{m,i}^2$ . The per-rank cost is  $c_m = d_{\text{in},m} + d_{\text{out},m}$ , corresponding to parameter counts; FLOPs scale in the same order but differ by a constant factor depending on sequence length. Given a global budget  $\mathcal{B}$ , we solve

$$\max_{\{r_m \in \mathbb{Z}_{\geq 0}\}} \sum_{m=1}^M F_m(r_m) \quad \text{s.t.} \quad \sum_{m=1}^M c_m r_m \leq \mathcal{B}. \quad (14)$$

### 4.3 LAYER-WISE NORMALIZATION

To remove spurious scale differences across layers while preserving within-layer order, we normalize utilities by the leading singular value:

$$\tilde{v}_{m,i} = \frac{v_{m,i}}{v_{m,1}} = \frac{s_{m,i}^2}{s_{m,1}^2}, \quad \tilde{F}_m(r) = \sum_{i=1}^r \tilde{v}_{m,i}. \quad (15)$$

This preserves the relative importance of directions within a layer while placing all layers on a comparable scale, similar in spirit to normalization in PCA and related factor models (Jolliffe & Cadima, 2016).

#### 4.4 CONVEX RELAXATION AND WATER-FILLING

Relax equation 14 to  $r_m \in \mathbb{R}_{\geq 0}$  with Lagrange multiplier  $\lambda \geq 0$ . KKT stationarity yields a *global density cutoff*:

$$\tilde{v}_{m,r_m^*}/c_m \geq \lambda^* \text{ and } \tilde{v}_{m,r_m^*+1}/c_m < \lambda^*, \quad \forall m, \quad (16)$$

i.e., keep all directions whose density (utility per cost) exceeds  $\lambda^*$ . In discrete form this is the familiar water-filling rule (Tse & Viswanath, 2005):

$$\text{keep } (m, i) \text{ iff } \tilde{v}_{m,i}/c_m \geq \lambda^*, \quad \sum_m c_m r_m(\lambda^*) = \mathcal{B}. \quad (17)$$

Because  $\tilde{v}_{m,i}$  is nonincreasing in  $i$ , selections are per-layer prefixes. Under this prefix property, a standard exchange argument shows that the relaxation admits an integral optimum; see Appendix D for details.

#### 4.5 RANK FLOORS AND CONSTRAINED GLOBAL ALLOCATION

We allow per-block floors  $r_m^{\min}$  for blocks  $m \in \{1, \dots, M\}$ . Let  $k_m := \text{rank}(\mathbf{W}_m)$ , per-rank cost  $c_m > 0$ , global budget  $\mathcal{B}$ , values  $v_{m,i}$  (e.g.,  $v_{m,i} = \frac{1}{2}s_{m,i}^2$ ), and densities  $\rho_{m,i} := v_{m,i}/c_m$ .

**Feasibility.** The floor-constrained problem

$$\max_{\{r_m \in \mathbb{Z}_{\geq 0}\}} \sum_{m=1}^M \sum_{i=1}^{r_m} v_{m,i} \quad \text{s.t.} \quad r_m \geq r_m^{\min} \quad (\forall m), \quad \sum_{m=1}^M c_m r_m \leq \mathcal{B}$$

is feasible iff

$$0 \leq r_m^{\min} \leq k_m \quad (\forall m), \quad \sum_{m=1}^M c_m r_m^{\min} \leq \mathcal{B}. \quad (18)$$

**Allocation when feasible.** Pre-assign floors and water-fill the remainder.

$$\mathcal{B}' \leftarrow \mathcal{B} - \sum_{m=1}^M c_m r_m^{\min}, \quad (19)$$

$$r_m^* = r_m^{\min} + r_m^+(\lambda^*), \quad (20)$$

$$r_m^+(\lambda) := \max \left\{ r \leq k_m - r_m^{\min} : \frac{v_{m,r_m^{\min}+r}}{c_m} \geq \lambda \right\}, \quad (21)$$

where  $\lambda^*$  is chosen (by bisection) so that  $\sum_m c_m r_m^+(\lambda^*) = \mathcal{B}'$ . Equivalently, keep all remaining directions with  $\rho_{m,i} \geq \lambda^*$ . [Algorithm 1](#) summarizes the full procedure. The final tie-fill step resolves ties that occur at the density cutoff: if multiple layers have the same marginal density at  $\lambda^*$ , we increment the one with the largest next density by one rank, repeating this until the budget is fully consumed.

**Result and connection to AWF.** Applying AWF with the ORA ranks  $\{r_m^*\}$  yields  $\mathbf{B}_m^* \mathbf{A}_m^* = \arg \min_{\text{rank} \leq r_m^*} \|\mathbf{W}_m \Sigma_{\text{in},m}^{1/2} - \tilde{\mathbf{W}}_m \Sigma_{\text{in},m}^{1/2}\|_F^2$ . By construction (additivity of the whitened quadratic loss), ORA maximizes  $\sum_m \tilde{F}_m(r_m)$  under  $\mathcal{B}$ , giving the smallest second-order loss increase among all rank assignments with the same cost. A full proof of optimality under the prefix property is provided in Appendix D.

## 5 EXPERIMENTS

We evaluate our framework *AFORA* (AWF + ORA) under strict **training-free** conditions. All results are obtained without gradient updates, using only a small calibration stream to estimate input covariances. We compare against standard training-free baselines and conduct ablations to isolate the effects of each component.

**Algorithm 1** ORA: Water-filling with Normalized Utilities and Floors

**Require:** Targets  $\{W_m\}_{m=1}^M$ , costs  $\{c_m\}$ , budget  $\mathcal{B}$ , floors  $\{f_m\}$   
**Ensure:** Assigned ranks  $\{r_m^*\}$

- 1: **Per-layer spectra:** For each  $m$ , estimate  $\Sigma_{\text{in},m}$  on calibration data; form  $T_m = W_m \Sigma_{\text{in},m}^{1/2}$ ; compute singular values  $\mathbf{s}_m = (s_{m,1}, s_{m,2}, \dots)$ .
- 2: **Normalized utilities:**  $\tilde{v}_{m,i} \leftarrow s_{m,i}^2 / s_{m,1}^2$  for all  $i$ .
- 3: **Floors & feasibility:**  $B_{\min} \leftarrow \sum_m c_m f_m$ ; if  $\mathcal{B} < B_{\min}$ , abort (infeasible).
- 4: **Initialize bisection:**  $\lambda_{\min} \leftarrow 0$ ,  $\lambda_{\max} \leftarrow \max_{m,i} \tilde{v}_{m,i} / c_m$ .
- 5: **while**  $\lambda_{\max} - \lambda_{\min} > \epsilon$  **do**
- 6:    $\lambda \leftarrow (\lambda_{\min} + \lambda_{\max}) / 2$ .
- 7:    $r_m(\lambda) \leftarrow f_m + \#\{i > f_m : \tilde{v}_{m,i} / c_m \geq \lambda\}$  for all  $m$ .
- 8:   **if**  $\sum_m c_m r_m(\lambda) > \mathcal{B}$  **then**
- 9:      $\lambda_{\min} \leftarrow \lambda$
- 10:   **else**
- 11:      $\lambda_{\max} \leftarrow \lambda$
- 12:   **end if**
- 13: **end while**
- 14: **Tie-fill:** Set  $r_m^* \leftarrow r_m(\lambda^*)$ . While  $\sum_m c_m r_m^* < \mathcal{B}$ , increment the  $m$ .
- 15: **return**  $\{r_m^*\}$ .

## 5.1 EXPERIMENTAL SETUP

**Model.** We compress LLaMA-2-7B, LLaMA-3-8B (Touvron et al., 2023), and QWEN-3-8B (Bai et al., 2023), each evaluated in the bfloat16 configuration for consistency across methods. Each model consists of 32 decoder layers with multi-head attention and feed-forward blocks. Our compression is applied to the attention projections (query, key, value, and output matrices). Unless noted otherwise, the same compression ratio is applied across the full model using our global allocator.

**Datasets and metrics.** For language modeling we report perplexity on WikiText-2 (Merity et al., 2016) and Penn Treebank (PTB) (Marcinkiewicz, 1994). For zero-shot reasoning we evaluate accuracy on ARC-Easy (Clark et al., 2018), HellaSwag (Zellers et al., 2019), PIQA (Bisk et al., 2020), WinoGrande (Sakaguchi et al., 2020), CB (De Marneffe et al., 2019), OpenbookQA (Mihaylov et al., 2018), and GSM8K (Cobbe et al., 2021). Evaluation follows the LM Evaluation Harness protocol (Gao et al., 2024). Perplexity is measured with standard left-to-right generation, while reasoning accuracy is reported as multiple-choice accuracy without any in-context examples.

**Compression ratio and budgets.** We report compression ratios as the fraction of parameters removed from the dense model. Since only the attention projection matrices are compressed, ranks are determined using our ORA algorithm, which constructs the whitened spectra of each projection matrix and identifies the retained directions via a density cutoff obtained through bisection search. Costs are measured in parameter count and therefore scale directly with the sum of retained ranks across layers. For example, in the case of LLaMA-2, a 0.15 compression ratio corresponds to removing about 47.07% of the attention weights, which proportionally reduces the average ranks of the projection matrices.

**Computation cost.** The computational cost of estimating  $\Sigma_{\text{in}}$  is proportional to a single forward pass over the calibration set. This step scales linearly with the number of calibration sequences. The subsequent SVD is performed once per linear projection layer, matching the cost structure of other SVD-based compression methods.

**Hardware and software environment.** All experiments are run on a server equipped with 4 NVIDIA RTX 6000 Ada Generation GPUs (48GB memory each), CUDA 12.4, and driver version 550.144.03. Models are evaluated in mixed-precision (bfloat16) with PyTorch 2.3 and HuggingFace Transformers 4.44. Multi-GPU evaluation is handled with the Accelerate library. Reported perplexities and accuracies are averaged over three independent runs to reduce variance.

## 5.2 MAIN RESULTS

**Perplexity under compression.** Table 1 reports perplexity at 15% compression. At this budget, our method exhibits the least performance degradation, maintaining results close to the dense model while all SVD-based baselines degrade substantially. Compared with ASVD, our approach achieves the same compression ratio with  $9.6\times$  faster compression time, highlighting its practical efficiency in large-scale settings. All compression times were measured on a single RTX 6000 Ada GPU under the experimental setup described in Section 5.1. The large gap in compression time stems from how ranks are assigned: ASVD determines ranks by exhaustively testing multiple compression ratios for every layer and measuring the resulting perplexity drop, which is computationally expensive. In contrast, AFORA performs a single global optimization step—constructing layerwise utilities and identifying the density cutoff through an  $O(\log N)$  bisection search—allowing it to match ASVD’s performance while reducing compression time by an order of magnitude.

Table 1: Language modeling performance at 15% compression ratio on LLaMA-2-7B.

Method	Compression Ratio in MHA $\uparrow$	PPL(WikiText-2) $\downarrow$	PPL(PTB) $\downarrow$	Compression time $\downarrow$
Dense	0	10.46	128.01	
SVD	0.47	851.48	2512.70	435 <i>sec</i>
FWSVD	0.47	940.68	1896.48	322 <i>sec</i>
ASVD		13.15	147.69	14604 <i>sec</i>
SVD-LLM		15.12	174.37	963 <i>sec</i>
SLEB (prune)		13.58	<b>141.54</b>	
<b>AWF</b>	0.47	<u>12.88</u> $\pm 0.10$	162.26 $\pm 2.37$	1127 <i>sec</i>
<b>AFORA</b>	0.47	<b>11.80</b> $\pm 0.06$	<u>148.70</u> $\pm 2.88$	1522 <i>sec</i>

Repeat for compression ratio 5% and 10% in Appendix F.1.

**Zero-shot reasoning.** Table 2 shows that AFORA preserves accuracy across multiple reasoning benchmarks. Across most tasks, our method consistently ranks first or second in accuracy, confirming that it preserves generalization capability. At the same time, we observe a small but consistent pattern: AFORA tends to perform better on truth-judgment tasks such as CB, whereas ASVD sometimes holds an edge on more specialized or multi-step reasoning benchmarks like OpenbookQA and GSM8K.

Table 2: Zero-shot evaluation at 15% compression ratio on LLaMA-2-7B.

Method	ARC-E	HellaSwag	PIQA	WinoGrande	CB	OpenbookQA	GSM8K	Avg.
Dense	0.71	0.70	0.79	0.68	0.48	0.28	0.13	0.54
SVD	0.29	0.35	0.49	0.50	0.41	0.12	0.00	0.31
FWSVD	0.37	0.56	0.67	0.62	0.41	0.22	0.00	0.41
ASVD	<b>0.62</b>	<b>0.65</b>	<u>0.71</u>	0.69	<u>0.45</u>	<b>0.28</b>	<b>0.03</b>	<b>0.49</b>
SVD-LLM	0.57	0.60	0.66	<b>0.71</b>	0.41	0.22	<u>0.02</u>	0.46
SLEB (depth)	0.54	0.53	0.67	0.58	0.41	0.22	0.00	0.42
<b>AWF</b>	0.52	0.59	0.68	<u>0.70</u>	0.38	0.17	0.01	0.44
	$\pm 0.05$	$\pm 0.03$	$\pm 0.01$	$\pm 0.02$	$\pm 0.04$	$\pm 0.01$	$\pm 0.01$	$\pm 0.00$
<b>AFORA</b>	<u>0.61</u>	<u>0.64</u>	<b>0.74</b>	0.66	<b>0.52</b>	<u>0.23</u>	0.01	<b>0.49</b>
	$\pm 0.01$	$\pm 0.02$	$\pm 0.01$	$\pm 0.01$	$\pm 0.05$	$\pm 0.02$	$\pm 0.00$	$\pm 0.00$

Repeat for compression ratio 5% and 10% in Appendix F.2.

Table 3 shows that **AWF** is the dominant contributor, giving the largest perplexity reduction compared to SVD. **ORA** degrades performance when applied on top of SVD, indicating that naive rank allocation cannot compensate for poor factorization. This is because the optimization objective used in ORA is derived from the activation-aware formulation of AWF. As a result, ORA becomes meaningful only when paired with AWF, where it successfully refines rank assignments and yields further

improvements. Table 4 further confirms that replacing dense linear matrices with low-rank counterparts provides consistent inference-time gains in terms of wall-clock latency, and the measured acceleration closely follows the theoretical reduction in FLOPs for the multi-head attention modules. For clarity, the FLOP counts follow the standard cost of the attention projections. In particular, the dense model requires  $O(4Ld^2)$  operations for the  $q, k, v, o$  projections, whereas the low-rank variant uses  $O(2LdR)$  for rank- $R$  factors. All reported values in Table 4 use sequence length  $L = 1024$  and batch size 1. Importantly, this trend highlights that our method achieves computation savings that scale in proportion to, or even beyond, the nominal compression ratio, whereas many alternative approaches only offer sublinear reductions in FLOPs and wall-clock time relative to their compression rates (Li et al., 2023). This proportionality between compression and actual computational gain underscores a key strength of our approach over competing methods. Additional ablations highlighting the specific role of ORA—both its behavior on top of ASVD and its layerwise rank patterns—are provided in Appendix F.2 and Appendix J.

Table 3: Ablation study on LLaMA-2-7B.

Ratio	Metric	SVD	AWF	SVD + ORA	AFORA
5%	PPL ↓	94.67	10.97 ± 0.05	862.40	<b>10.95</b> ± 0.04
10%	PPL ↓	310.16	11.58 ± 0.06	1624.27	<b>11.36</b> ± 0.05
15%	PPL ↓	851.48	12.88 ± 0.10	3409.46	<b>11.80</b> ± 0.06

Table 4: Wall-clock speedup.

Compression Ratio	FLOPs <sub>MHA</sub>	Speedup
dense	1.00× ↓	1.00× ↑
0.05	0.84×	1.09×
0.10	0.69×	1.15×
0.15	0.53×	1.21×

### 5.3 SENSITIVITY ANALYSES

We report two supplementary analyses related to calibration size and cross-architecture behavior.

**Effect of calibration set size.** All experimental results use 200 single-sequence batches (batch size = 1, sequence length = 1024) to estimate  $\Sigma_{in}$ . To assess sensitivity to this choice, we vary the number of calibration batches while keeping all other settings fixed. Figure 2 shows that performance stabilizes once the batch count exceeds roughly 100. This suggests that our default setting (200 batches) is conservative but provides a reliably stable estimate of  $\Sigma_{in}$ .

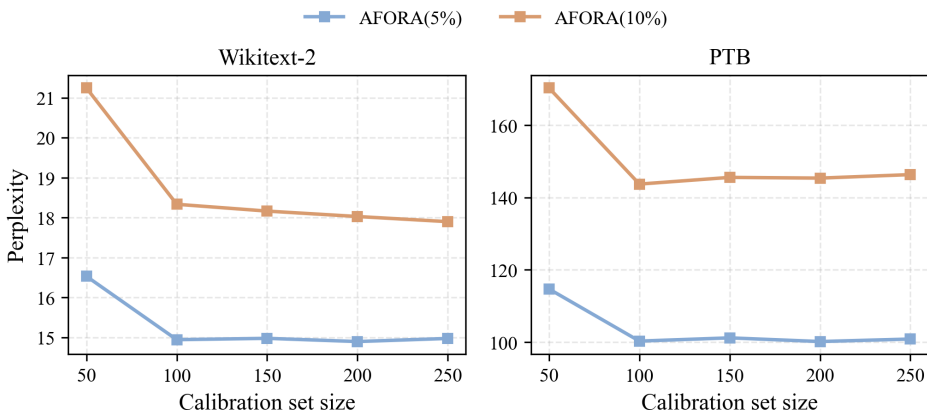


Figure 2: The effect of calibration set size.

**Cross-model generality.** We additionally apply AFORA to LLaMA-3-8B and Qwen-3-8B under the same training-free pipeline and comparable compression ratios. Both models show the same qualitative trend as LLaMA-2: AFORA consistently outperforms weight-space SVD and achieves stable perplexity under a fixed budget. All detailed results—including full perplexity tables, zero-shot accuracy—are provided in Appendix G, H.

## 6 FUTURE WORK

Our study opens several directions for further improvement and broader application of activation-aware, optimally allocated low-rank compression.

### 6.1 TOWARDS FULLY LOSS-AWARE FORMULATIONS

For tractability, we assumed  $\Sigma_{\text{out}} = I$  in the quadratic loss approximation. A natural extension is to incorporate nontrivial output covariances, either estimated from calibration data or approximated via Fisher information. This would yield a more refined objective

$$\|\Sigma_{\text{out}}^{1/2}(W - BA)\Sigma_{\text{in}}^{1/2}\|_F^2 \quad (22)$$

capturing gradient sensitivities in addition to input statistics.

### 6.2 INTEGRATION WITH OTHER COMPRESSION PARADIGMS

While AFORA is training-free and hardware-friendly, it can be combined with pruning and quantization. **Structured pruning can remove inactive blocks first, after which AFORA allocates ranks only to the remaining projections, avoiding unnecessary factorization.** For quantization, the ordering matters: applying quantization before AFORA introduces noise into the weight matrices and distorts the whitened operator, leading to suboptimal utilities and rank decisions. In contrast, applying quantization after AFORA is often more stable, as the information-balanced factorization reduces outlier directions in the factors  $(B, A)$  and can therefore lower quantization error.

### 6.3 ALTERNATIVE ALLOCATION STRATEGIES

Beyond the water-filling allocator, we also plan to evaluate the **incremental allocation** algorithm as an alternative (Fox, 1966). Unlike water-filling, which admits an efficient closed-form solution, the incremental allocation proceeds in discrete increments and is known to incur higher computational complexity. By comparing the two approaches, we aim to characterize the complexity–performance tradeoff and assess whether the incremental scheme offers any practical benefits in the context of rank allocation.

## 7 CONCLUSION

We introduced **AFORA**, a training-free and activation-aware framework for LLM compression. Our approach is built on two key components:

1. **Activation-aware reduction (AWF)**: a per-layer low-rank factorization that directly minimizes input-aware reconstruction error, providing a principled alternative to weight-space SVD.
2. **Optimal rank allocation (ORA)**: a global allocator that formulates rank assignment as a budget-constrained optimization problem, solved efficiently via water-filling and yielding provably optimal assignments.

We showed that AFORA outperforms truncated SVD and pruning-based baselines. At a 15% compression ratio, AFORA reduces perplexity on WikiText-2 by 8.36% compared to the best existing method, shortens compression algorithm execution time by up to 9.6 $\times$ , and preserves zero-shot accuracy on reasoning benchmarks at a level comparable to the dense model.

Beyond empirical performance, our analysis emphasizes that activation statistics and optimization-theoretic allocation can transform compression from heuristic design into a systematic optimization problem. AFORA thus provides both a practical tool for efficient LLM deployment and a theoretical foundation for further refinements, such as incorporating richer loss-aware objectives or exploring tighter integration with pruning-based approaches.

In summary, AFORA demonstrates that combining activation-aware factorization with principled global rank allocation offers a viable and theoretically grounded pathway toward efficient, training-free deployment of large language models.

## REFERENCES

- 540  
541  
542 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-  
543 man, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical  
544 report. *arXiv preprint arXiv:2303.08774*, 2023.
- 545  
546 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge,  
547 Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
- 548  
549 Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. Piqa: Reasoning about physical com-  
550 monsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*,  
551 volume 34, pp. 7432–7439, 2020.
- 552  
553 Aleksandar Botev, Hippolyt Ritter, and David Barber. Practical gauss-newton optimisation for deep  
554 learning. In *International Conference on Machine Learning*, pp. 557–565. PMLR, 2017.
- 555  
556 Stephen P Boyd and Lieven Vandenberghe. *Convex optimization*. Cambridge university press, 2004.
- 557  
558 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,  
559 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are  
560 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 561  
562 Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and  
563 Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge.  
564 *arXiv preprint arXiv:1803.05457*, 2018.
- 565  
566 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,  
567 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to  
568 solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- 569  
570 Marie-Catherine De Marneffe, Mandy Simons, and Judith Tonhauser. The commitmentbank: In-  
571 vestigating projection in naturally occurring discourse. In *proceedings of Sinn und Bedeutung*,  
572 volume 23, pp. 107–124, 2019.
- 573  
574 Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning  
575 of quantized llms. *Advances in neural information processing systems*, 36:10088–10115, 2023.
- 576  
577 Carl Eckart and Gale Young. The approximation of one matrix by another of lower rank. *Psychome-*  
578 *trika*, 1(3):211–218, 1936.
- 579  
580 Vage Egiazarian, Roberto L Castro, Denis Kuznedelev, Andrei Panferov, Eldar Kurtic, Shubhra  
581 Pandit, Alexandre Marques, Mark Kurtz, Saleh Ashkboos, Torsten Hoefler, et al. Bridging  
582 the gap between promise and performance for microscaling fp4 quantization. *arXiv preprint*  
583 *arXiv:2509.23202*, 2025.
- 584  
585 Bennett Fox. Discrete optimization via marginal analysis. *Management science*, 13(3):210–216,  
586 1966.
- 587  
588 Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and Dan Alistarh. Gptq: Accurate post-training  
589 quantization for generative pre-trained transformers. *arXiv preprint arXiv:2210.17323*, 2022.
- 590  
591 Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Fos-  
592 ter, Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas Muen-  
593 nighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang  
Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. The language model  
evaluation harness, 07 2024. URL <https://zenodo.org/records/12608602>.
- Roger Grosse and James Martens. A kronecker-factored approximate fisher matrix for convolution  
layers. In *International Conference on Machine Learning*, pp. 573–582. PMLR, 2016.
- Song Han, Huizi Mao, and William J Dally. Deep compression: Compressing deep neural networks  
with pruning, trained quantization and huffman coding. *arXiv preprint arXiv:1510.00149*, 2015.

- 594 Torsten Hoefer, Dan Alistarh, Tal Ben-Nun, Nikoli Dryden, and Alexandra Peste. Sparsity in deep  
595 learning: Pruning and growth for efficient inference and training in neural networks. *Journal of*  
596 *Machine Learning Research*, 22(241):1–124, 2021.
- 597 Yen-Chang Hsu, Ting Hua, Sungen Chang, Qian Lou, Yilin Shen, and Hongxia Jin. Language model  
598 compression with weighted low-rank factorization. *arXiv preprint arXiv:2207.00112*, 2022.
- 600 Ian T Jolliffe and Jorge Cadima. Principal component analysis: a review and recent developments.  
601 *Philosophical transactions of the royal society A: Mathematical, Physical and Engineering Sci-*  
602 *ences*, 374(2065):20150202, 2016.
- 603 Hans Kellerer, Ulrich Pferschy, and David Pisinger. Multidimensional knapsack problems. In *Knap-*  
604 *sack problems*, pp. 235–283. Springer, 2004.
- 606 François Lagunas, Ella Charlaix, Victor Sanh, and Alexander M Rush. Block pruning for faster  
607 transformers. *arXiv preprint arXiv:2109.04838*, 2021.
- 608 Guangyan Li, Yongqiang Tang, and Wensheng Zhang. Lorap: Transformer sub-layers deserve dif-  
609 ferentiated structured compression for large language models. *arXiv preprint arXiv:2404.09695*,  
610 2024.
- 612 Zhuo Li, Hengyi Li, and Lin Meng. Model compression for deep neural networks: A survey.  
613 *Computers*, 12(3):60, 2023.
- 614 Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Wei-Ming Chen, Wei-Chen Wang, Guangxuan  
615 Xiao, Xingyu Dang, Chuang Gan, and Song Han. Awq: Activation-aware weight quantization  
616 for on-device llm compression and acceleration. *Proceedings of machine learning and systems*,  
617 6:87–100, 2024.
- 618 Mary Ann Marcinkiewicz. Building a large annotated corpus of english: The penn treebank. *Using*  
619 *Large Corpora*, 273:31, 1994.
- 621 James Martens and Roger Grosse. Optimizing neural networks with kronecker-factored approximate  
622 curvature. In *International conference on machine learning*, pp. 2408–2417. PMLR, 2015.
- 623 Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture  
624 models. *arXiv preprint arXiv:1609.07843*, 2016.
- 625 Paul Michel, Omer Levy, and Graham Neubig. Are sixteen heads really better than one? *Advances*  
626 *in neural information processing systems*, 32, 2019.
- 628 Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct  
629 electricity? a new dataset for open book question answering. *arXiv preprint arXiv:1809.02789*,  
630 2018.
- 631 Leon Mirsky. Symmetric gauge functions and unitarily invariant norms. *The quarterly journal of*  
632 *mathematics*, 11(1):50–59, 1960.
- 634 Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adver-  
635 sarial winograd schema challenge at scale. In *Proceedings of the AAAI Conference on Artificial*  
636 *Intelligence*, volume 34, pp. 8732–8740, 2020.
- 637 Victor Sanh, Thomas Wolf, and Alexander Rush. Movement pruning: Adaptive sparsity by fine-  
638 tuning. *Advances in neural information processing systems*, 33:20378–20389, 2020.
- 640 Jiwon Song, Kyungseok Oh, Taesu Kim, Hyungjun Kim, Yulhwa Kim, and Jae-Joon Kim. Sleb:  
641 Streamlining llms through redundancy verification and elimination of transformer blocks. *arXiv*  
642 *preprint arXiv:2402.09025*, 2024.
- 643 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-  
644 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda-  
645 tion and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 646 David Tse and Pramod Viswanath. *Fundamentals of wireless communication*. Cambridge university  
647 press, 2005.

648 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
649 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural informa-*  
650 *tion processing systems*, 30, 2017.

651 Xin Wang, Yu Zheng, Zhongwei Wan, and Mi Zhang. Svd-llm: Truncation-aware singular value  
652 decomposition for large language model compression. *arXiv preprint arXiv:2403.07378*, 2024.

653 Juyun Wee, Minjae Park, and Jaeho Lee. Prompt-based depth pruning of large language models.  
654 *arXiv preprint arXiv:2502.04348*, 2025.

655 Guangxuan Xiao, Ji Lin, Mickael Seznec, Hao Wu, Julien Demouth, and Song Han. Smoothquant:  
656 Accurate and efficient post-training quantization for large language models. In *International*  
657 *conference on machine learning*, pp. 38087–38099. PMLR, 2023.

658 Zhewei Yao, Reza Yazdani Aminabadi, Minjia Zhang, Xiaoxia Wu, Conglong Li, and Yuxiong  
659 He. Zeroquant: Efficient and affordable post-training quantization for large-scale transformers.  
660 *Advances in neural information processing systems*, 35:27168–27183, 2022.

661 Zhihang Yuan, Yuzhang Shang, Yue Song, Qiang Wu, Yan Yan, and Guangyu Sun. Asvd:  
662 Activation-aware singular value decomposition for compressing large language models. *arXiv*  
663 *preprint arXiv:2312.05821*, 2023.

664 Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a ma-  
665 chine really finish your sentence? *arXiv preprint arXiv:1905.07830*, 2019.

## 666 A OPTIMALITY OF TRUNCATED SVD

667 The following theorem is a classical result due to Eckart–Young–Mirsky. It formally establishes that  
668 truncated SVD gives the optimal rank- $r$  approximation under the Frobenius norm.

669 **Theorem A.1** (Eckart–Young–Mirsky). *For any matrix  $\mathbf{W} \in \mathbb{R}^{d_{out} \times d_{in}}$  with singular values  $\sigma_1 \geq$   
670  $\sigma_2 \geq \dots \geq \sigma_{\min(d_{out}, d_{in})}$ , the solution to*

$$671 \min_{\text{rank}(\mathbf{X}) \leq r} \|\mathbf{W} - \mathbf{X}\|_F \quad (23)$$

672 is given by

$$673 \mathbf{X}^* = \mathbf{U}_r \mathbf{S}_r \mathbf{V}_r^\top, \quad (24)$$

674 where  $\mathbf{U}_r, \mathbf{S}_r, \mathbf{V}_r$  are the top- $r$  singular components of  $\mathbf{W}$ . Moreover, the approximation error is

$$675 \|\mathbf{W} - \mathbf{X}^*\|_F^2 = \sum_{i=r+1}^{\min(d_{out}, d_{in})} \sigma_i^2. \quad (25)$$

676 This establishes truncated SVD as the optimal low-rank solution in the standard Frobenius sense.  
677 Our activation-aware formulation (Sec. 3) extends this principle by replacing the weight-space  
678 Frobenius norm with a whitened, loss-aware metric.

## 679 B FROM ACTIVATION-AWARE PROBLEM TO COVARIANCE FORMULATION

680 We show how the activation-aware objective in equation 2 is equivalent to the covariance formulation  
681 in equation 3.

682 Starting point:

$$683 \min_{\text{rank}(BA) \leq r} \mathbb{E}_x \left[ \|Wx - BAx\|_F^2 \right]. \quad (26)$$

684 Expand the Frobenius norm:

$$685 \mathbb{E}_x \left[ \|Wx - BAx\|_F^2 \right] = \mathbb{E}_x \left[ (Wx - BAx)^\top (Wx - BAx) \right] \quad (27)$$

$$686 = \mathbb{E}_x \left[ x^\top (W - BA)^\top (W - BA)x \right]. \quad (28)$$

Bring expectation inside:

$$= \text{Tr}\left((W - BA)^\top (W - BA) \mathbb{E}[xx^\top]\right). \quad (29)$$

Define the input covariance  $\Sigma_{\text{in}} = \mathbb{E}[xx^\top]$ , then

$$\mathbb{E}_x \left[ \|Wx - BAx\|_F^2 \right] = \text{Tr}\left((W - BA)^\top (W - BA) \Sigma_{\text{in}}\right). \quad (30)$$

Rewriting with the Frobenius norm identity:

$$= \|(W - BA)\Sigma_{\text{in}}^{1/2}\|_F^2. \quad (31)$$

Therefore the problem is equivalent to

$$\min_{\text{rank}(BA) \leq r} \|W\Sigma_{\text{in}}^{1/2} - BA\Sigma_{\text{in}}^{1/2}\|_F^2, \quad (32)$$

which is exactly equation 3.

## C APPENDIX: DERIVATION OF INFORMATION-BALANCED FACTORIZATION

We derive the closed-form scaling for information-balanced factorization. Starting from the truncated SVD factors

$$B_0 = U_r \text{diag}(\sqrt{s_{1:r}}), \quad (33)$$

$$A_0 = \text{diag}(\sqrt{s_{1:r}}) V_r^\top \Sigma_{\text{in}}^{-1/2}, \quad (34)$$

the column norms of  $B_0$  and row norms of  $A_0$  may be unbalanced.

Introduce positive scalings  $\alpha_i > 0$ :

$$B(\alpha) = U_r \text{diag}(\alpha_1, \dots, \alpha_r), \quad (35)$$

$$A(\alpha) = \text{diag}\left(\frac{s_i}{\alpha_i}\right) V_r^\top \Sigma_{\text{in}}^{-1/2}. \quad (36)$$

For each component  $i$ , the relevant norms are

$$\|B(\alpha)_{:,i}\|_2 = \alpha_i, \quad (37)$$

$$\|A(\alpha)_{i,:}\|_2 \propto \frac{s_i}{\alpha_i}, \quad (38)$$

$$\|(\Sigma_{\text{in}}^{1/2} V_r)_{:,i}\|_2 = \text{fixed}. \quad (39)$$

Balancing the first two terms against the third gives

$$\alpha_i^* = \sqrt{s_i \cdot \|(\Sigma_{\text{in}}^{1/2} V_r)_{:,i}\|_2}, \quad i = 1, \dots, r. \quad (40)$$

## D OPTIMALITY OF ORA UNDER PREFIX STRUCTURE

We show that the convex relaxation of the ORA objective (Sec. 4) admits an integral optimum under the prefix structure of singular values. This establishes that the water-filling procedure yields an exact solution when no floors are imposed.

Consider the allocation problem

$$\max_{\{r_m \in \mathbb{Z}_{\geq 0}\}} \sum_{m=1}^M \sum_{i=1}^{r_m} v_{m,i} \quad \text{s.t.} \quad \sum_{m=1}^M c_m r_m \leq \mathcal{B}, \quad (41)$$

with per-rank utilities  $v_{m,i}$  that are nonincreasing in  $i$ . The problem is a variant of the multiple-choice knapsack problem (MCKP) (Boyd & Vandenberghe, 2004), where each direction  $(m, i)$  is an item with value  $v_{m,i}$  and cost  $c_m$ .

756 Because singular values are sorted, utilities satisfy

$$757 \quad v_{m,1} \geq v_{m,2} \geq \dots, \quad (42)$$

758  
759 so each block must be taken as a prefix  $\{1, \dots, r_m\}$ . Otherwise, replacing a later singular direction  
760 with an earlier one increases value at no higher cost.

761 Relaxing to  $r_m \in \mathbb{R}_{\geq 0}$  with multiplier  $\lambda \geq 0$  gives the stationarity condition

$$762 \quad \frac{v_{m,i}}{c_m} \geq \lambda \text{ for } i \leq r_m^*, \quad \frac{v_{m,i}}{c_m} < \lambda \text{ for } i > r_m^*. \quad (43)$$

763  
764 This is exactly the water-filling rule: keep all directions with density above the global cutoff  $\lambda^*$ .  
765 Since selections are per-layer prefixes, the solution  $r_m^*$  is automatically integer-valued. If ties occur  
766 at the cutoff, multiple integer solutions exist with equal objective value.

767 When minimum ranks  $r_m^{\min}$  are required, the same prefix argument applies after fixing the floors.  
768 Feasibility requires

$$769 \quad \sum_m c_m r_m^{\min} \leq \mathcal{B}. \quad (44)$$

770 If infeasible, some floors must be reduced. Our heuristic (Sec. 4.5) reduces the floor with the smallest  
771 boundary density until feasible. This is practical but not guaranteed globally optimal; more exact  
772 strategies (e.g., integer programming) could be applied.

773 Thus ORA admits an exact integral optimum under the prefix structure, and the water-filling proce-  
774 dure recovers it directly. With floors, optimality holds for the feasible case, while infeasible cases  
775 require heuristic handling.

## 776 E LAYER-WISE WHITENED SPECTRA

777 We whiten inputs with the estimated second moment and compute singular values per linear module:

$$778 \quad T_m = W_m \Sigma_{\text{in},m}^{1/2}, \quad (45)$$

$$779 \quad T_m = U_m \text{diag}(s_{m,1}, \dots, s_{m,d_m}) V_m^\top. \quad (46)$$

780 We visualize the squared singular values  $s_{m,i}^2$  together with the cumulative energy

$$801 \quad E_m(k) = \frac{\sum_{i=1}^k s_{m,i}^2}{\sum_{i=1}^{d_m} s_{m,i}^2}, \quad (47)$$

802 which indicates how quickly energy concentrates across ranks. Large early values of  $s_{m,i}^2$  (and  
803 rapidly rising  $E_m$ ) mean a few directions capture most of the input-aware effect, so low ranks suf-  
804 fice; flatter spectra imply higher ranks are needed. Across layers and modules, the shapes vary  
805 substantially, and we use these diagnostics to summarize compression difficulty and guide budgeted  
806 allocation (see Fig. 3).

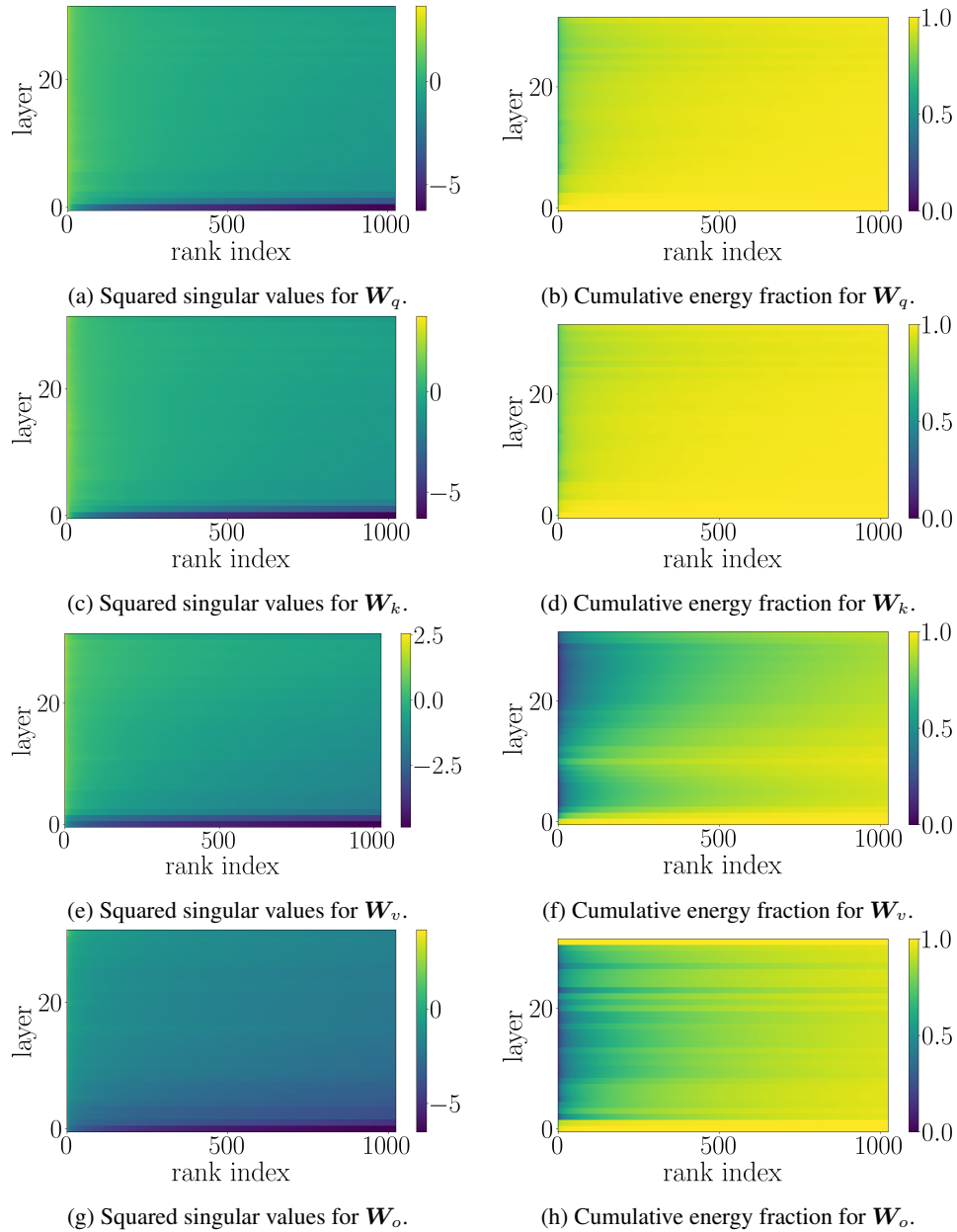


Figure 3: **Heatmaps of spectra across modules (attention projections  $q, k, v, o$ ).** Each row corresponds to one projection type ( $q, k, v, o$  from top to bottom). The left column shows squared singular values  $s_{m,i}^2$  across layers (horizontal axis: rank index; vertical axis: layer depth). The right column shows the corresponding cumulative energy curves  $E_m(k)$ .

## F ADDITIONAL EXPERIMENTS ON LLaMA-2-7B

## F.1 LANGUAGE MODELING (PPL) VS. BUDGET ON LLaMA-2-7B

Table 5: Language modeling performance at 5% compression ratio.

Method	Compression Ratio in MHA $\uparrow$	PPL(WikiText-2) $\downarrow$	PPL(PTB) $\downarrow$
Dense		10.46	128.01
SVD	0.16	94.67	400.01
FWSVD	0.16	136.91	489.16
ASVD		<b>10.35</b>	<b>127.45</b>
SVD-LLM		13.40	151.78
SLEB (prune)		11.48	130.94
<b>AWF</b>	0.16	10.97 $\pm$ 0.05	129.48 $\pm$ 0.31
<b>AFORA</b>	0.16	<u>10.95</u> $\pm$ 0.04	<u>131.47</u> $\pm$ 0.58

Table 6: Language modeling performance at 10% compression ratio.

Method	Compression Ratio in MHA $\uparrow$	PPL(WikiText-2) $\downarrow$	PPL(PTB) $\downarrow$
Dense		10.46	128.01
SVD	0.31	310.16	1082.86
FWSVD	0.31	396.18	969.46
ASVD		<b>11.12</b>	<b>127.52</b>
SVD-LLM		14.18	161.82
SLEB (prune)		11.99	132.45
<b>AWF</b>	0.31	11.58 $\pm$ 0.06	139.23 $\pm$ 1.72
<b>AFORA</b>	0.31	<u>11.36</u> $\pm$ 0.05	<u>137.39</u> $\pm$ 0.55

## F.2 ZERO-SHOT ACCURACY

Table 7: Zero-shot evaluation at 5% compression ratio on LLaMA-2-7B.

Method	ARC-E	HellaSwag	PIQA	WinoGrande	CB	OpenbookQA	GSM8K	Avg.
Dense	0.71	0.70	0.79	0.68	0.48	0.28	0.13	0.54
SVD	0.46	0.45	0.60	0.58	0.46	0.21	0.00	0.39
FWSVD	<u>0.67</u>	0.67	0.74	0.67	0.39	0.25	0.01	0.49
ASVD	<u>0.67</u>	<u>0.69</u>	<u>0.76</u>	<b>0.70</b>	0.34	<b>0.28</b>	<u>0.09</u>	0.50
SVD-LLM	0.59	0.61	0.70	0.67	0.41	0.21	0.05	0.46
SLEB (depth)	0.65	<b>0.70</b>	0.72	0.65	0.30	<u>0.26</u>	0.02	0.47
<b>AWF</b>	0.66	0.63	0.73	<b>0.70</b>	<b>0.59</b>	<u>0.26</u>	<u>0.07</u>	<u>0.52</u>
	$\pm$ 0.01	$\pm$ 0.01	$\pm$ 0.02	$\pm$ 0.04	$\pm$ 0.08	$\pm$ 0.02	$\pm$ 0.02	$\pm$ 0.01
<b>AFORA</b>	<b>0.69</b>	<u>0.69</u>	<b>0.79</b>	0.67	<b>0.59</b>	<u>0.26</u>	<u>0.07</u>	<b>0.54</b>
	$\pm$ 0.01	$\pm$ 0.01	$\pm$ 0.00	$\pm$ 0.01	$\pm$ 0.02	$\pm$ 0.01	$\pm$ 0.02	$\pm$ 0.00

Table 8: Zero-shot evaluation at 10% compression ratio on LLaMA-2-7B.

Method	ARC-E	HellaSwag	PIQA	WinoGrande	CB	OpenbookQA	GSM8K	Avg.
Dense	0.71	0.70	0.79	0.68	0.48	0.28	0.13	0.54
SVD	0.28	0.46	0.52	0.52	0.36	0.12	0.01	0.32
FWSVD	0.51	0.60	0.69	0.65	0.44	0.25	0.00	0.45
ASVD	0.66	<b>0.68</b>	0.75	<b>0.74</b>	0.48	<b>0.31</b>	<b>0.09</b>	<b>0.53</b>
SVD-LLM	0.56	0.62	0.70	0.70	0.41	0.21	0.03	0.46
SLEB (depth)	0.59	0.60	0.69	0.59	0.45	0.23	0.01	0.45
<i>AWF</i>	0.63	0.61	<b>0.78</b>	0.68	<b>0.56</b>	0.24	0.03	0.52
	$\pm 0.02$	$\pm 0.02$	$\pm 0.01$	$\pm 0.02$	$\pm 0.02$	$\pm 0.02$	$\pm 0.02$	$\pm 0.01$
<i>AFORA</i>	<b>0.71</b>	0.65	<b>0.78</b>	0.68	<b>0.56</b>	0.24	0.03	0.52
	$\pm 0.03$	$\pm 0.02$	$\pm 0.02$	$\pm 0.01$	$\pm 0.05$	$\pm 0.01$	$\pm 0.00$	$\pm 0.01$

## G ADDITIONAL EXPERIMENTS ON LLaMA-3-8B

### G.1 LANGUAGE MODELING (PPL) VS. BUDGET ON LLaMA-3-8B

Table 9: Language modeling performance at 2.5% compression ratio.

Method	Compression Ratio in MHA $\uparrow$	PPL(WikiText-2) $\downarrow$	PPL(PTB) $\downarrow$
Dense		14.16	90.90
SVD	0.35	24.63	133.51
FWSVD	0.35	50425.62	21320.24
ASVD		<b>14.33</b>	<b>93.31</b>
<i>AFORA</i>	0.15	14.59 $\pm 0.02$	95.91 $\pm 0.44$

Table 10: Language modeling performance at 5% compression ratio.

Method	Compression Ratio in MHA $\uparrow$	PPL(WikiText-2) $\downarrow$	PPL(PTB) $\downarrow$
Dense		14.16	90.90
SVD	0.35	24.63	133.51
FWSVD	0.35	50425.62	21320.24
ASVD		15.43	<b>98.01</b>
<i>AFORA</i>	0.30	14.89 $\pm 0.04$	98.15 $\pm 0.40$

Table 11: Language modeling performance at 7.5% compression ratio.

Method	Compression Ratio in MHA $\uparrow$	PPL(WikiText-2) $\downarrow$	PPL(PTB) $\downarrow$
Dense		14.16	90.90
SVD	0.45	85.72	436.64
FWSVD	0.45	32150.11	24532.78
ASVD		18.17	107.56
<i>AFORA</i>	0.45	15.30 $\pm 0.04$	105.45 $\pm 0.64$

## G.2 ZERO-SHOT ACCURACY

Table 12: Zero-shot evaluation at 2.5% compression ratio on LLaMA-3-8B.

Method	ARC-E	HellaSwag	PIQA	WinoGrande	CB	OpenbookQA	GSM8K	Avg.
Dense	0.82	0.68	0.82	0.74	0.63	0.31	0.53	0.65
SVD	0.69	0.65	0.76	0.70	0.41	0.20	0.07	0.50
FWSVD	0.33	0.32	0.55	0.50	0.41	0.08	0.02	0.32
ASVD	<b>0.80</b>	<u>0.68</u>	<u>0.82</u>	<b>0.76</b>	<b>0.54</b>	<b>0.45</b>	<b>0.50</b>	<b>0.65</b>
<i>AFORA</i>	<u>0.76</u>	<u>0.70</u>	<b>0.83</b>	<b>0.76</b>	<u>0.52</u>	<u>0.30</u>	<u>0.37</u>	<u>0.61</u>
	± 0.01	± 0.01	± 0.01	± 0.01	± 0.02	± 0.01	± 0.04	± 0.01

Table 13: Zero-shot evaluation at 5% compression ratio on LLaMA-3-8B.

Method	ARC-E	HellaSwag	PIQA	WinoGrande	CB	OpenbookQA	GSM8K	Avg.
Dense	0.82	0.68	0.82	0.74	0.63	0.31	0.53	0.65
SVD	0.69	0.65	0.76	0.70	0.41	0.20	0.07	0.50
FWSVD	0.33	0.32	0.55	0.50	0.41	0.08	0.02	0.32
ASVD	<b>0.77</b>	<b>0.70</b>	<b>0.83</b>	<b>0.79</b>	<u>0.48</u>	<b>0.36</b>	<b>0.50</b>	<b>0.63</b>
<i>AFORA</i>	<u>0.74</u>	<u>0.66</u>	<u>0.82</u>	<u>0.76</u>	<b>0.49</b>	<u>0.30</u>	<u>0.25</u>	<u>0.58</u>
	± 0.02	± 0.02	± 0.02	± 0.01	± 0.02	± 0.01	± 0.04	± 0.00

Table 14: Zero-shot evaluation at 7.5% compression ratio on LLaMA-3-8B.

Method	ARC-E	HellaSwag	PIQA	WinoGrande	CB	OpenbookQA	GSM8K	Avg.
Dense	0.82	0.68	0.82	0.74	0.63	0.31	0.53	0.65
SVD	0.60	0.53	0.75	0.62	0.38	0.20	0.03	0.44
FWSVD	0.25	0.33	0.49	0.52	0.41	0.07	0.00	0.30
ASVD	<b>0.75</b>	<b>0.68</b>	<b>0.79</b>	<b>0.77</b>	<b>0.54</b>	<u>0.30</u>	<b>0.26</b>	<b>0.58</b>
<i>AFORA</i>	<u>0.73</u>	<u>0.66</u>	<u>0.78</u>	<u>0.74</u>	<u>0.41</u>	<b>0.31</b>	<u>0.14</u>	<u>0.54</u>
	± 0.02	± 0.02	± 0.01	± 0.02	± 0.00	± 0.02	± 0.03	± 0.01

## H ADDITIONAL EXPERIMENTS ON QWEN-3-8B

### H.1 LANGUAGE MODELING (PPL) VS. BUDGET ON QWEN-3-8B

Table 15: Language modeling performance at 2.5% compression ratio.

Method	Compression Ratio in MHA $\uparrow$	PPL(WikiText-2) $\downarrow$	PPL(PTB) $\downarrow$
Dense		18.47	116.95
SVD	0.35	33.76	273.60
FWSVD	0.35	166445.64	116773.84
ASVD		<b>18.58</b>	<b>109.15</b>
<i><b>AFORA</b></i>	0.14	<u>19.62</u> $\pm$ 0.03	<u>119.69</u> $\pm$ 1.19

Table 16: Language modeling performance at 5% compression ratio.

Method	Compression Ratio in MHA $\uparrow$	PPL(WikiText-2) $\downarrow$	PPL(PTB) $\downarrow$
Dense		18.47	116.95
SVD	0.35	33.76	273.60
FWSVD	0.35	166445.64	116773.84
ASVD		<b>20.00</b>	<b>118.02</b>
<i><b>AFORA</b></i>	0.28	<u>20.40</u> $\pm$ 0.07	<u>130.11</u> $\pm$ 4.81

Table 17: Language modeling performance at 7.5% compression ratio.

Method	Compression Ratio in MHA $\uparrow$	PPL(WikiText-2) $\downarrow$	PPL(PTB) $\downarrow$
Dense		18.47	116.95
SVD	0.42	114.51	6053.02
FWSVD	0.42	4740616.01	3600659.72
ASVD		<u>21.77</u>	<b>126.90</b>
<i><b>AFORA</b></i>	0.42	<b>20.78</b> $\pm$ 0.10	<u>137.17</u> $\pm$ 3.50

## H.2 ZERO-SHOT ACCURACY

Table 18: Zero-shot evaluation at 2.5% compression ratio on QWEN-3-8B.

Method	ARC-E	HellaSwag	PIQA	WinoGrande	CB	OpenbookQA	GSM8K	Avg.
Dense	0.82	0.67	0.82	0.78	0.73	0.30	0.91	0.72
SVD	0.74	0.62	<u>0.77</u>	0.64	<u>0.57</u>	<u>0.30</u>	0.52	<u>0.59</u>
FWSVD	0.37	0.43	<u>0.59</u>	0.42	<u>0.27</u>	0.19	0.02	0.33
ASVD	<u>0.79</u>	<b>0.64</b>	<b>0.84</b>	<b>0.77</b>	0.16	<b>0.32</b>	<u>0.62</u>	<u>0.59</u>
<i>AFORA</i>	<b>0.80</b>	<u>0.63</u>	<u>0.79</u>	<u>0.75</u>	<b>0.76</b>	0.29	<b>0.81</b>	<b>0.69</b>
	± 0.02	± 0.01	± 0.02	± 0.02	± 0.04	± 0.01	± 0.01	± 0.01

Table 19: Zero-shot evaluation at 5% compression ratio on QWEN-3-8B.

Method	ARC-E	HellaSwag	PIQA	WinoGrande	CB	OpenbookQA	GSM8K	Avg.
Dense	0.82	0.67	0.82	0.78	0.73	0.30	0.91	0.72
SVD	0.74	<u>0.62</u>	<u>0.77</u>	0.64	<u>0.57</u>	<u>0.30</u>	0.52	0.59
FWSVD	0.37	0.43	<u>0.59</u>	0.42	<u>0.27</u>	0.19	0.02	0.33
ASVD	<b>0.79</b>	<u>0.62</u>	<b>0.83</b>	<b>0.73</b>	0.39	<b>0.31</b>	<u>0.58</u>	<u>0.61</u>
<i>AFORA</i>	<u>0.78</u>	<b>0.63</b>	<u>0.77</u>	<u>0.71</u>	<b>0.75</b>	0.27	<b>0.71</b>	<b>0.66</b>
	± 0.02	± 0.00	± 0.03	± 0.02	± 0.03	± 0.01	± 0.05	± 0.02

Table 20: Zero-shot evaluation at 7.5% compression ratio on QWEN-3-8B.

Method	ARC-E	HellaSwag	PIQA	WinoGrande	CB	OpenbookQA	GSM8K	Avg.
Dense	0.82	0.67	0.82	0.78	0.73	0.30	0.91	0.72
SVD	0.48	<u>0.61</u>	0.69	0.60	<u>0.57</u>	0.23	0.00	0.45
FWSVD	0.40	0.36	0.58	0.44	0.48	0.12	0.00	0.34
ASVD	<b>0.81</b>	<b>0.62</b>	<b>0.83</b>	<b>0.75</b>	0.48	<b>0.32</b>	<u>0.31</u>	<u>0.59</u>
<i>AFORA</i>	<u>0.78</u>	0.60	<u>0.77</u>	<u>0.68</u>	<b>0.73</b>	<u>0.27</u>	<b>0.59</b>	<b>0.63</b>
	± 0.02	± 0.01	± 0.03	± 0.02	± 0.05	± 0.01	± 0.13	± 0.01

## I ADDITIONAL ABLATION RESULTS: ASVD COMBINED WITH ORA

To further isolate the contribution of ORA, we also evaluate the combination **ASVD + ORA**. This ablation helps clarify whether ORA provides value even when paired with a different activation-aware factorization method. Table 21 and Table 22 report results on LLaMA-3-8B at a compression ratio of 0.15. ORA yields small improvements over ASVD in zero-shot accuracy, although the perplexity can be higher, especially when utilities are derived from whitened operators (ORA<sup>†</sup>). These observations are consistent with the design of ASVD, which selects ranks by directly evaluating every possible choice, whereas ORA chooses ranks by solving a simple optimization problem. As a result, ASVD’s allocation tends to achieve lower perplexity, while ORA is far more computationally efficient and can still provide modest gains in generalization metrics.

Table 21: Perplexity of ASVD with and without ORA on LLaMA-3-8B (compression ratio 0.05).

Method	PPL(Wikitext-2)	PPL(PTB)
ASVD	14.33	98.01
ASVD + ORA	21.36 $\pm$ 0.22	152.45 $\pm$ 1.39

Table 22: Zero-shot accuracy of ASVD with and without ORA on LLaMA-3-8B (compression ratio 0.05).

Method	ARC-E	HellaSwag	PIQA	WinoGrande	CB	OpenbookQA	GSM8K	Avg.
ASVD	0.77	0.70	0.83	0.79	0.48	0.36	0.50	0.63
<b>ASVD + ORA</b>	0.82	0.68	0.82	0.74	0.63	0.31	0.53	0.65
	$\pm$ 0.00	$\pm$ 0.00	$\pm$ 0.00	$\pm$ 0.00	$\pm$ 0.00	$\pm$ 0.00	$\pm$ 0.00	$\pm$ 0.00

Overall, these results indicate that ORA remains compatible with alternative activation-aware factorizations such as ASVD. While ASVD’s exhaustive per-layer search often yields lower perplexity, ORA provides a significantly cheaper mechanism for rank allocation and can still improve zero-shot performance relative to ASVD alone.

## J RANK-ALLOCATION PATTERNS: ORA VS. ASVD

Figure 4 visualizes the rank distributions produced by ORA and ASVD. Recall that AFORA applies compression only to the attention projections, whereas ASVD also includes the MLP blocks. Within the attention layers, both methods show broadly similar prioritization patterns, with small differences due to whitening.

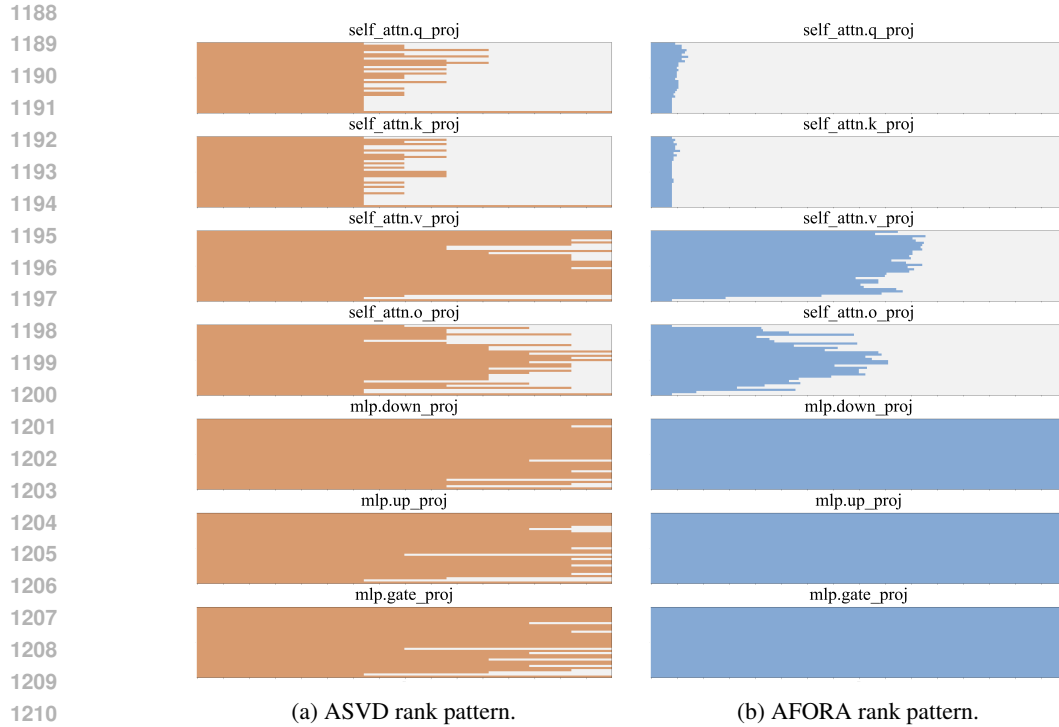


Figure 4: **Comparison of rank allocation: ASVD vs. AFORA.** ASVD uses exhaustive per-layer search, while AFORA derives rank assignments from activation-aware utilities and a global budget.

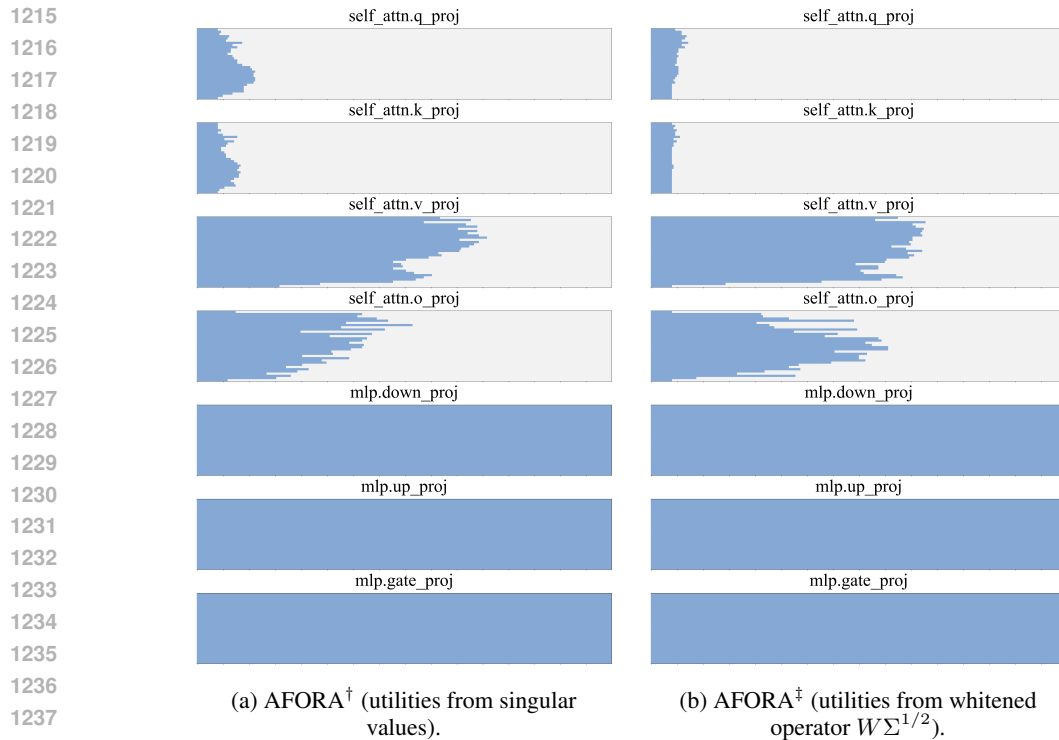


Figure 5: **Comparison of AFORA variants.** AFORA<sup>†</sup> computes utilities from the singular values of  $W$ , while AFORA<sup>‡</sup> uses the whitened operator  $W\Sigma^{1/2}$ , leading to slightly different rank distributions under the same compression budget.