

# COAP: Memory-Efficient Training with Correlation-Aware Gradient Projection

## Supplementary Material

### 1. Detailed Proposed Method

**Notation.** In this paper, we use the following notation conventions: Matrices and vectors are indicated with boldface capital and lowercase letters, *e.g.*,  $\mathbf{X}$  and  $\mathbf{x}$ , respectively. Tensors are represented using boldface calligraphic script, denoted as  $\mathcal{X}$ .

#### 1.1. Inter-projection Correlation-aware $P_t$ Update.

Considering the importance of incorporating the information from the previous projection into the current update, we propose calculating  $P_t$  via solving the optimization problem:

$$\min_{P_t} \underbrace{\text{MSE}(\hat{\mathbf{G}}_t, \mathbf{G}_t)}_{\text{reconstruction term}} \underbrace{(1 - \text{CosSim}(\hat{\mathbf{M}}_{t-1}, \mathbf{G}_t))}_{\text{direction term}}, \quad (1)$$

where  $\text{CosSim}(\cdot, \cdot)$  and  $\text{MSE}(\cdot, \cdot)$  return the cosine similarity and mean squared error, respectively. To simplify notation, the notation without the subscript  $t$  represents a general form of the optimization problem, *i.e.*,

$$\min_P \underbrace{\text{MSE}(\hat{\mathbf{G}}, \mathbf{G})}_{\text{reconstruction term}} \underbrace{(1 - \text{CosSim}(\hat{\mathbf{M}}, \mathbf{G}))}_{\text{direction term}}, \quad (2)$$

where  $\mathbf{P} \in \mathbb{R}^{n \times r}$ ,  $\hat{\mathbf{G}} \in \mathbb{R}^{m \times n} = \mathbf{G}\mathbf{P}\mathbf{P}^\top$  and  $\hat{\mathbf{M}} \in \mathbb{R}^{m \times n} = \mathbf{M}^{\text{proj}}\mathbf{P}^\top$  are the full-rank estimates of gradient and first-order moment projected back from the low-rank subspace, respectively. Notably, the *reconstruction term* is introduced to minimize the reconstruction error for gradients incurred by projection – achieving the similar goal that SVD essentially aims for, and *direction term* encourages the consistency of optimization direction after restoring from low-rank subspace.

To solve Eqn. 2 as a non-convex optimization problem, we propose using stochastic gradient descent to iteratively update  $P_t$  as follows:

$$\mathbf{P} := \mathbf{P} - \eta \left( \frac{\partial \text{MSE}(\hat{\mathbf{G}}, \mathbf{G})}{\partial \mathbf{P}} (1 - \text{CosSim}(\hat{\mathbf{M}}, \mathbf{G})) + \frac{\partial \text{CosSim}(\hat{\mathbf{M}}, \mathbf{G})}{\partial \mathbf{P}} \text{MSE}(\hat{\mathbf{G}}, \mathbf{G}) \right). \quad (3)$$

Here,  $\eta$  represents learning rate, set to 0.1 by default. The gradient expressions for the reconstruction term and the direction term are derived as follows:

#### Gradient of Reconstruction term (MSE).

$$\begin{aligned} \frac{\partial \text{MSE}(\hat{\mathbf{G}}, \mathbf{G})}{\partial \mathbf{P}} &= \frac{\partial}{\partial \mathbf{P}} \left( \frac{1}{mn} \text{tr}((\hat{\mathbf{G}} - \mathbf{G})^\top (\hat{\mathbf{G}} - \mathbf{G})) \right) \\ &= \frac{2}{mn} (\hat{\mathbf{G}}^\top \mathbf{G}\mathbf{P} - 2\mathbf{G}^\top \mathbf{G}\mathbf{P} + \mathbf{G}^\top \hat{\mathbf{G}}\mathbf{P}), \end{aligned} \quad (4)$$

where  $\text{tr}(\cdot)$  represents the trace of a matrix.

#### Gradient of Direction Term (CosSim).

Given the cosine similarity between matrices  $\hat{\mathbf{M}}$  and  $\mathbf{G}$ , defined as:

$$\begin{aligned} \text{CosSim}(\hat{\mathbf{M}}, \mathbf{G}) &= \frac{1}{m} \sum_{i=1}^m \text{CosSim}(\hat{\mathbf{M}}_i, \mathbf{G}_i) \\ &= \frac{1}{m} \sum_{i=1}^m \frac{\langle \hat{\mathbf{M}}_i, \mathbf{G}_i \rangle}{\|\hat{\mathbf{M}}_i\| \|\mathbf{G}_i\|}. \end{aligned} \quad (5)$$

Applying the chain rule to compute the gradient of  $\text{CosSim}(\hat{\mathbf{M}}, \mathbf{G})$  with respect to  $\mathbf{P}$ :

$$\begin{aligned} \frac{\partial \text{CosSim}(\hat{\mathbf{M}}, \mathbf{G})}{\partial \mathbf{P}} &= \left( \frac{\partial \text{CosSim}(\hat{\mathbf{M}}, \mathbf{G})}{\partial \hat{\mathbf{M}}} \right)^\top \frac{\partial \hat{\mathbf{M}}}{\partial \mathbf{P}} \\ &= \frac{1}{m} \sum_{i=1}^m \left( \frac{\partial \text{CosSim}(\hat{\mathbf{M}}_i, \mathbf{G}_i)}{\partial \hat{\mathbf{M}}_i} \right)^\top \mathbf{M}_i^{\text{proj}} \\ &= \frac{1}{m} \sum_{i=1}^m \left( \frac{\mathbf{G}_i}{\|\hat{\mathbf{M}}_i\| \|\mathbf{G}_i\|} - \frac{\hat{\mathbf{M}}_i \langle \hat{\mathbf{M}}_i, \mathbf{G}_i \rangle}{\|\hat{\mathbf{M}}_i\|^3 \|\mathbf{G}_i\|} \right)^\top \mathbf{M}_i^{\text{proj}}, \end{aligned} \quad (6)$$

where  $\|\cdot\|$  represents the Euclidean norm,  $\langle \cdot, \cdot \rangle$  denotes the inner product,  $\hat{\mathbf{M}}_i, \mathbf{G}_i, \mathbf{M}_i^{\text{proj}}$  denotes the  $i$ -th row of  $\hat{\mathbf{M}}, \mathbf{G}, \mathbf{M}^{\text{proj}}$ .

Incorporating the gradient expressions above, we derive the final update formula for  $\mathbf{P}$  as follows:

$$\begin{aligned} \mathbf{P} := & \mathbf{P} - \eta \left( \frac{2}{mn} (\hat{\mathbf{G}}^\top \mathbf{G}\mathbf{P} - 2\mathbf{G}^\top \mathbf{G}\mathbf{P} + \mathbf{G}^\top \hat{\mathbf{G}}\mathbf{P}) \right. \\ & \left. \left( 1 - \frac{1}{m} \sum_{i=1}^m \frac{\langle \hat{\mathbf{M}}_i, \mathbf{G}_i \rangle}{\|\hat{\mathbf{M}}_i\| \|\mathbf{G}_i\|} \right) + \right. \\ & \left. \frac{1}{m} \sum_{i=1}^m \left( \frac{\mathbf{G}_i}{\|\hat{\mathbf{M}}_i\| \|\mathbf{G}_i\|} - \frac{\hat{\mathbf{M}}_i \langle \hat{\mathbf{M}}_i, \mathbf{G}_i \rangle}{\|\hat{\mathbf{M}}_i\|^3 \|\mathbf{G}_i\|} \right)^\top \mathbf{M}_i^{\text{proj}} \right. \\ & \left. \frac{1}{mn} \text{tr}((\hat{\mathbf{G}} - \mathbf{G})^\top (\hat{\mathbf{G}} - \mathbf{G})) \right). \end{aligned} \quad (7)$$

### 1.2. Theory Analysis.

COAP proposes to approximately keep the same low-rank subspace across iterations. Intuitively, this requires the low-rank subspace to remain stable across iterations, which can be bounded in terms of the learning rate and Lipschitz constant. Theoretical justification composed of the following steps:

- Bounded update will lead to bounded change of the subspace.
- Equation 2 finds a subspace not much worse than the true subspace.

We aim to demonstrate that this procedure guarantees an approximation to the underlying true subspace. Let  $\mathbf{P}_t^{\natural}$  be the true subspace. Let  $\mathbf{P}_t^{\text{MSE}}$  and  $\mathbf{P}_t^{\text{SIM}}$  be the solutions of the MSE loss and CosSim loss respectively, then

$$\begin{aligned}\mathbf{P}_t^{\text{MSE}} &= \mathbf{P}_{t-1} \\ \mathbf{P}_t^{\text{SIM}} &= \operatorname{argmax} \left\langle \mathbf{G}_t \mathbf{P}_t, \mathbf{M}_{t-1}^{\text{proj}} \right\rangle,\end{aligned}\quad (8)$$

where  $\mathbf{M}_{t-1}^{\text{proj}} = \beta \mathbf{M}_{t-2}^{\text{proj}} + (1-\beta) \mathbf{P}_t^T \mathbf{G}_t = \sum_{i=0}^{t-1} \beta^i (1-\beta) \mathbf{P}_{t-i}^T \mathbf{G}_{t-i}$ ,  $\operatorname{argmax} \langle \cdot, \cdot \rangle$  maximizes an inner product between two entities. Therefore,

$$\begin{aligned}\mathbf{P}_t^{\text{SIM}} &= \sum_{i=0}^{t-1} \beta^i (1-\beta) \operatorname{argmax} \left\langle \mathbf{P}_t^T \mathbf{G}_t, \mathbf{P}_{t-i}^T \mathbf{G}_{t-i} \right\rangle \\ &= \sum_{i=0}^{t-1} \beta^i (1-\beta) \mathbf{P}_{t-i}\end{aligned}\quad (9)$$

**Theorem 1.1.** Assume that the gradient is Lipchitz with respect to the weight with Lipchitz constant  $L$ . Assume  $\kappa_1 \doteq \frac{\sigma_1}{\sigma_{r+1}} > 1$  and  $\kappa_r \doteq \frac{\sigma_r}{\sigma_{r+1}} > 1$  to be the ratio between the 1-th and  $r$ -th to the  $r+1$ -th singular value of the gradient  $\mathbf{G}$  respectively. Then the proposed procedure in Equation 2 recovers a subspace well approximating the true subspace.

$$\begin{aligned}\left\| \mathbf{P}_t - \mathbf{P}_t^{\natural} \right\|_F &\leq c(t \bmod (\lambda \times T_u)) \frac{r\kappa_1 + \min(m, n) - r}{\kappa_r - 1} L\eta.\end{aligned}\quad (10)$$

With learning rate  $\eta$ , Lipchitz property of the gradient gives

$$\begin{aligned}\left\| \mathbf{G}_{t-1} - \mathbf{G}_t \right\|_2 &\leq L \left\| \mathbf{W}_{t-1} - \mathbf{W}_t \right\|_2 \\ &= L \left\| \eta \mathbf{G}_{t-1}^{\text{proj}} \right\|_2 \leq L\eta \left\| \mathbf{G}_{t-1} \right\|_2.\end{aligned}\quad (11)$$

Consequently, the optimal subspace  $\mathbf{P}_t^{\natural}$  is also close enough to the previous optimal subspace  $\mathbf{P}_{t-1}^{\natural}$ :

$$\left\| \mathbf{P}_{t-1}^{\natural} - \mathbf{P}_t^{\natural} \right\|_F \leq \frac{r\kappa_1 + \min(m, n) - r}{\kappa_r - 1} L\eta.\quad (12)$$

When  $t \bmod (\lambda \times T_u) = 0$ , the proposed algorithm updates  $\mathbf{P}_t = \mathbf{P}_t^{\natural}$ . Otherwise, we show that  $\mathbf{P}_t$  highly approx-

imates  $\mathbf{P}_t^{\natural}$ ,

$$\begin{aligned}\left\| \mathbf{P}_t - \mathbf{P}_t^{\natural} \right\|_F &\leq \left\| \mathbf{P}_{t-1}^{\natural} - \mathbf{P}_t^{\natural} \right\|_F + \left\| \mathbf{P}_t - \mathbf{P}_{t-1}^{\natural} \right\|_F \\ &\leq \left\| \mathbf{P}_{t-1}^{\natural} - \mathbf{P}_t^{\natural} \right\|_F + \left\| \mathbf{P}_t^{\text{MSE}} - \mathbf{P}_{t-1}^{\natural} \right\|_F + \left\| \mathbf{P}_t^{\text{SIM}} - \mathbf{P}_{t-1}^{\natural} \right\|_F \\ &\leq \left\| \mathbf{P}_{t-1}^{\natural} - \mathbf{P}_t^{\natural} \right\|_F + (2-\beta) \left\| \mathbf{P}_{t-1} - \mathbf{P}_{t-1}^{\natural} \right\|_F \\ &\quad + (1-\beta) \sum_{i=2} \beta^i \left\| \mathbf{P}_{t-i} - \mathbf{P}_{t-1}^{\natural} \right\|_F \\ &\leq \left\| \mathbf{P}_{t-1}^{\natural} - \mathbf{P}_t^{\natural} \right\|_F + 2 \left\| \mathbf{P}_{t-1} - \mathbf{P}_{t-1}^{\natural} \right\|_F \\ &\quad + (1-\beta) \sum_{i=2} \beta^i \left\| \mathbf{P}_{t-i} - \mathbf{P}_{t-1}^{\natural} \right\|_F \\ &\leq 2^{(t \bmod (\lambda \times T_u))} \frac{r\kappa_1 + \min(m, n) - r}{\kappa_r - 1} L\eta\end{aligned}\quad (13)$$

### 1.3. Hyper-parameter Settings.

Table 1 presents the detailed hyper-parameters used in all experiments of COAP. Except for the Update Interval ( $T_u$ ) and Re-project Factor ( $\lambda$ ) specific to COAP, all other parameters remain consistent with the corresponding baselines. The optimizer, rank settings, learning rate, GPU configuration, and mixed precision are listed for each experiment. Hyper-parameters for different models, including LDM, SiT-XL/2, ControlNet SDXL, and LLaMA series, are provided to ensure reproducibility.

### 1.4. Adafactor with COAP.

Adafactor is an adaptive gradient optimization algorithm designed to reduce memory usage while maintaining training efficiency. Unlike Adam, which stores full-matrix second-moment estimates, Adafactor factorizes the second-moment accumulation to save memory, making it particularly suitable for training large-scale models. However, first-moment remains crucial for stabilizing training and accelerating convergence. Without it, optimization can become unstable, especially in deep networks where gradients vary significantly. The first moment helps smooth updates and improves the overall training dynamics. Algorithm 1 describes the Adafactor-based training procedure, by projecting the gradient and its first-moment estimate into a lower-dimensional subspace, we effectively compress the momentum while retaining its benefits. This approach allows us to maintain training stability with reduced memory overhead, making it highly efficient for large-scale models.

### 1.5. Extension to CONV Layer.

To enhance the generality and applicability of our algorithm, it is essential to extend support to higher-dimensional weight tensors, which are prevalent in architectures such as



Table 1. Table 1: Hyper-parameter settings for all experiments.

Optimizer	Exprement	Rank $r$	Rank Ratio $\alpha$	Update Interval $T_u$	Re-project Factor $\lambda$	Learning Rate $\eta$	GPU	Mixed-precision
Adafactor (COAP)	Table 1. Pre-training LDM	-	2	16	10	$2 \times 10^{-5}$	8×V100	FP32
	Table 2. Pre-training SiT-XL/2	512	-	200	5	$1 \times 10^{-4}$	8×H100	FP16
	Table 3. ControlNet SDXL	-	2	8	10	$1 \times 10^{-5}$	8×H100	BF16
		-	4	8	10	$1 \times 10^{-5}$	8×H100	BF16
AdamW (COAP)	Table 1. Pre-training LDM	-	2	16	10	$2 \times 10^{-5}$	8×V100	FP32
	Table 2. Pre-training SiT-XL/2	512	-	30	10	$1 \times 10^{-4}$	8×H100	FP16
	Table 5. Pre-training LLaMA-1B	512	-	40	5	$1 \times 10^{-2}$	8×H100	BF16
	Table 5. Pre-training LLaMA-7B	1024	-	100	1	$5 \times 10^{-3}$	8×H100	BF16
	Table 6. Fine-tuning LLaVA-v1.5-7B	-	4	32	1	$2 \times 10^{-5}$	1×A100	BF16

### Algorithm 1: Adafactor with COAP

**Input:** Weight matrix  $\mathbf{W} \in \mathbb{R}^{m \times n}$ , Learning rate  $\eta$ , Rank  $r$ , Betas  $[\beta_1, \beta_2]$ , Update interval  $[\lambda, T_u]$ , Decay rate  $\gamma$ .  
**Initialize:**  $\mathbf{M}_0^{\text{proj}} \in \mathbb{R}^{m \times r} \leftarrow 0$ ,  $\mathbf{V}_0^{\text{proj}} \in \mathbb{R}^{m \times r} \leftarrow 0$ ,  $t \leftarrow 0$ ,  $\mathbf{R}_0 \in \mathbb{R}^{m \times 1} \leftarrow 0$ ,  $\mathbf{C}_0 \in \mathbb{R}^{1 \times r} \leftarrow 0$ .  
**Randomly Initialize:**  $\mathbf{P}_0 \in \mathbb{R}^{n \times r}$   
**Compute:**  $\mathbf{P}_0 \leftarrow (\mathbf{P}_0, \mathbf{G}_0)$  ▷ Occasional Low-cost SVD  
**for**  $t$  in  $[1, 2, \dots]$  **do**  
    **Compute:** gradient  $\mathbf{G}_t$  of  $\mathbf{W}_t$  in the loss function.  
    **if**  $t \bmod T_u = 0$  **then**  
        **if**  $t \bmod (\lambda \times T_u) = 0$  **then**  
            ▷ Occasional Low-cost SVD  
            **Compute:**  $\mathbf{P}_t \leftarrow (\mathbf{P}_{t-1}, \mathbf{G}_t)$   
        **else**  
            **Update:**  $\mathbf{P}_t \leftarrow (\mathbf{P}_{t-1}, \mathbf{G}_t, \mathbf{M}_{t-1})$  ▷ Eqn. 2  
    **else**  
         $\mathbf{P}_t \leftarrow \mathbf{P}_{t-1}$   
        ▷ Project gradient and moments into low-rank space.  
         $\beta_2 \leftarrow 1 - t^\gamma$   
         $\mathbf{G}_t^{\text{proj}} \leftarrow \mathbf{G}_t \mathbf{P}_t$   
         $\mathbf{M}_t^{\text{proj}} \leftarrow \beta_1 \mathbf{M}_{t-1}^{\text{proj}} + (1 - \beta_1) \mathbf{G}_t^{\text{proj}}$   
         $\mathbf{R}_t = \beta_2 \mathbf{R}_{t-1} + (1 - \beta_2) \cdot \text{Sum}(\mathbf{G}_t^{\text{proj}^2}, -1)$   
         $\mathbf{C}_t = \beta_2 \mathbf{C}_{t-1} + (1 - \beta_2) \cdot \text{Sum}(\mathbf{G}_t^{\text{proj}^2}, -2)$   
         $\hat{\mathbf{V}}_t = \sqrt{\frac{\text{Mean}(\mathbf{R}_t, -1)}{\mathbf{R}_t \mathbf{C}_t}}$   
        ▷ Calculate the bias correction term in low-rank space.  
         $\Delta \mathbf{W}_t^{\text{proj}} \leftarrow \beta_1 \mathbf{M}_t^{\text{proj}} + (1 - \beta_1) \eta \hat{\mathbf{V}}_t \odot \mathbf{G}_t^{\text{proj}}$   
        ▷ Restore  $\Delta \mathbf{W}_t^{\text{proj}}$  to original space and update  $\mathbf{W}$ .  
         $\mathbf{W}_t \leftarrow \mathbf{W}_{t-1} - \Delta \mathbf{W}_t^{\text{proj}} \mathbf{P}_t^\top$   
**Return:** updated  $\mathbf{W}$

Convolutional Neural Networks (CNNs). While the most straightforward approach would be to reshape CNN weights into matrices and apply the same low-rank space construction method used for matrix operations, this naive strategy would inevitably result in the loss of intrinsic spatial characteristics inherent to CNNs [5].

For a convolutional layer with a weight tensor  $\mathcal{W} \in$

$\mathbb{R}^{O \times I \times K_1 \times K_2}$ , where  $I$  and  $O$  are the number of input channels and output channels, respectively, and  $K_1$  and  $K_2$  are the kernel sizes, we use Tucker-2 decomposition [7] as the factorization method. In this scenario,  $\mathcal{W}$  can be represented with a core tensor  $\mathcal{C}$  and two factor matrices ( $\mathbf{U}_1$  and  $\mathbf{U}_2$ ) along each mode as follows:

$$\mathcal{W} = \mathcal{C} \times_1 \mathbf{U}_1 \times_2 \mathbf{U}_2, \quad (14)$$

where “ $\times_n$ ” denotes the  $n$ -mode product. Specifically,  $\mathbf{U}_1 \in \mathbb{R}^{O \times r_O}$  represents the left singular vectors of the mode-1 unfolding of  $\mathcal{W}$ , denoted as  $\mathbf{W}_{(1)}$ , i.e.,  $\mathbf{W}_{(1)} = \mathbf{U}_1 \Sigma_1 \mathbf{V}_1^\top$ . Similarly,  $\mathbf{U}_2 \in \mathbb{R}^{I \times r_I}$  represents the left singular vectors of the mode-2 unfolding of  $\mathcal{W}$ , denoted as  $\mathbf{W}_{(2)}$ , i.e.,  $\mathbf{W}_{(2)} = \mathbf{U}_2 \Sigma_2 \mathbf{V}_2^\top$ . The core tensor  $\mathcal{C}$  is obtained by projecting  $\mathcal{W}$  onto the subspaces spanned by  $\mathbf{U}_1$  and  $\mathbf{U}_2$ , i.e.,

$$\mathcal{C} = \mathcal{W} \times_1 \mathbf{U}_1^\top \times_2 \mathbf{U}_2^\top, \quad (15)$$

where  $\mathcal{C} \in \mathbb{R}^{r_O \times r_I \times K_1 \times K_2}$ . Here,  $r_O$  and  $r_I$  are the Tucker-2 tensor ranks, determining the dimensionality of the factor matrices and the core tensor. Initially, the values for  $\mathbf{U}_1$ ,  $\mathbf{U}_2$ , and  $\mathcal{C}$  are typically obtained using Higher-Order Singular Value Decomposition (HOSVD) [1]. To refine these initial values and achieve the final decomposition, the Alternating Least Squares (ALS) [4] method is employed. This iterative optimization technique alternates between updating the core tensor  $\mathcal{C}$  and the factor matrices  $\mathbf{U}_1$  and  $\mathbf{U}_2$  to minimize the reconstruction error.

Algorithm 2 outlines the Adam-based training procedure, integrating the proposed  $\mathbf{P}_t$  update method for convolutional layers.

Typically, for a convolutional layer with a weight tensor  $\mathcal{W} \in \mathbb{R}^{O \times I \times K_1 \times K_2}$ ,  $I$  and  $O$  are significantly larger than the kernel sizes  $K_1$  and  $K_2$  (i.e.,  $I, O \gg K_1, K_2$ ). Therefore, we propose using the format of Tucker-2 decomposition to handle CNNs while employing our own decomposition method for the actual factorization. According to

Eq. 14 and Eq. 15, the low-rank projection space of  $\mathcal{G}_t$  becomes  $[P_{O_t} \in \mathbb{R}^{O \times r_O}, P_{I_t} \in \mathbb{R}^{I \times r_I}]$ . The gradient  $\mathcal{G}_t$  in the low-rank space is  $\mathcal{G}_t^{\text{proj}} = \mathcal{G}_t \times_1 P_{O_t}^\top \times_2 P_{I_t}^\top$ , and the restored tensor from the low-rank space is  $\hat{\mathcal{G}}_t = \mathcal{G}_t^{\text{proj}} \times_1 P_{O_t} \times_2 P_{I_t}$ . Here,  $P_{O_t}$  and  $P_{I_t}$  can be updated according to Eqn. 7 and Occasional Low-cost SVD on the mode-1 and mode-2 unfolding of tensor  $\mathcal{G}$ , respectively.

---

**Algorithm 2:** Adam with COAP (CONV)

---

**Input:** Weight tensor  $\mathcal{W} \in \mathbb{R}^{O \times I \times K_1 \times K_2}$ , Learning rate  $\eta$ , Rank ratio  $\alpha$ , Betas  $[\beta_1, \beta_2]$ , Update interval  $[\lambda, T_u]$ .

**Initialize:**  $r_O = O \sqrt{\alpha}$ ,  $r_I = I \sqrt{\alpha}$ ,  $t \leftarrow 0$ ,  
 $\mathcal{M}_0^{\text{proj}} \in \mathbb{R}^{r_O \times r_I \times K_1 \times K_2} \leftarrow 0$ ,  
 $\mathcal{V}_0^{\text{proj}} \in \mathbb{R}^{r_O \times r_I \times K_1 \times K_2} \leftarrow 0$

**Randomly Initialize:**  $P_O \in \mathbb{R}^{O \times r_O}$ ,  $P_I \in \mathbb{R}^{I \times r_I}$

**Define:**  $G_{O_t} \leftarrow \text{reshape}(\mathcal{G}_t, [O, IK_1K_2])$ ,  
 $G_{I_t} \leftarrow \text{reshape}(\mathcal{G}_t, [I, OK_1K_2])$

**Compute:**  $P_{O_0} \leftarrow (P_O, G_{O_0})$ ,  $P_{I_0} \leftarrow (P_I, G_{I_0})$   
 $\triangleright$  Occasional Low-cost SVD

**for**  $t$  in  $[1, 2, \dots]$  **do**

**Compute:** gradient  $G_t$  of  $W_t$  in the loss function.

**if**  $t \bmod T_u = 0$  **then**

**if**  $t \bmod (\lambda \times T_u) = 0$  **then**

**Compute:**  $P_{O_t} \leftarrow (P_{O_{t-1}}, G_{O_t})$ ,  
       $P_{I_t} \leftarrow (P_{I_{t-1}}, G_{I_t})$   
       $\triangleright$  Occasional Low-cost SVD

**else**

**Update:**  $P_{O_t} \leftarrow (P_{O_{t-1}}, G_{O_t}, M_{O_{t-1}})$   $\triangleright$   
      Eqn. 7

**Update:**  $P_{I_t} \leftarrow (P_{I_{t-1}}, G_{I_t}, M_{I_{t-1}})$   $\triangleright$   
      Eqn. 7

**else**

$P_{O_t} \leftarrow P_{O_{t-1}}$ ,  $P_{I_t} \leftarrow P_{I_{t-1}}$

$\triangleright$  Project gradient and moments into low-rank space.  
 $\mathcal{G}_t^{\text{proj}} \leftarrow \mathcal{G}_t \times_1 P_{O_t}^\top \times_2 P_{I_t}^\top$   
 $\mathcal{M}_t^{\text{proj}} \leftarrow \beta_1 \mathcal{M}_{t-1}^{\text{proj}} + (1 - \beta_1) \mathcal{G}_t^{\text{proj}}$   
 $\mathcal{V}_t^{\text{proj}} \leftarrow \beta_2 \mathcal{V}_{t-1}^{\text{proj}} + (1 - \beta_2) (\mathcal{G}_t^{\text{proj}})^2$

$\triangleright$  Calculate the bias correction term in low-rank space.  
 $\Delta \mathcal{W}_t^{\text{proj}} \leftarrow \frac{\mathcal{M}_t^{\text{proj}} / (1 - \beta_1^t)}{\sqrt{\mathcal{V}_t^{\text{proj}} / (1 - \beta_2^t) + \epsilon}}$

$\triangleright$  Restore  $\Delta \mathcal{W}_t^{\text{proj}}$  to original space and update  $\mathcal{W}$ .  
 $\mathcal{W}_t \leftarrow \mathcal{W}_{t-1} - \eta \Delta \mathcal{W}_t^{\text{proj}} \times_1 P_{O_t} \times_2 P_{I_t}$

**Return:** updated  $\mathcal{W}$

---

### 1.6. Impact of Low-rank Matrix Projection Formats on CNN Models Performance.

We compare different formats of Tucker decomposition, *i.e.*, Tucker-1, Tucker-2, and Tucker. In this context, the default Tucker format applies projections along all dimensions of a tensor. For instance, if the tensor is 4-dimensional, it requires 4 projection matrices. Tucker-2, on the other hand, uses only two projection matrices, while Tucker-1 requires

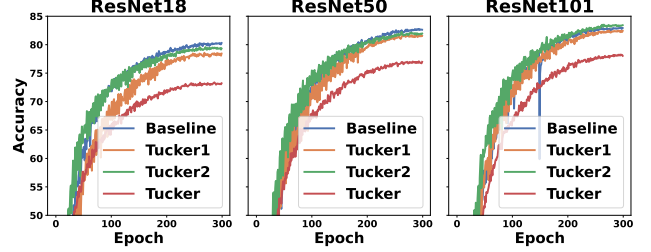


Figure 1. The comparison of Top-1 accuracy for ResNet under different low-rank projection formats is conducted with a rank ratio of 4. The models are trained for 100 epochs on the CIFAR-100 dataset.

Table 2. Pre-training DDPM on CIFAR-10 and CelebA-HQ datasets on  $8 \times V100$ . FID scores are reported, along with the GPU memory usage of optimizer states in FP32 format.

Dataset	Method	Rank Ratio	Optimizer Mem. (MB) $\downarrow$	FID $\downarrow$
CIFAR-10 ( $32 \times 32$ ) 800K steps	AdamW	-	272.72	5.42
	GaLore	1.5	302.43	6.09
	<b>COAP</b>	1.5	<b>214.66</b>	5.66
	Adafactor	-	222.71	5.43
	GaLore	1.5	196.11	7.14
	<b>COAP</b>	1.5	<b>180.87</b>	<b>5.41</b>
CelebA-HQ ( $256 \times 256$ ) 460K steps	AdamW	-	867.26	12.82
	GaLore	2	562.56	27.95
	<b>COAP</b>	2	<b>525.18</b>	17.37
	Adafactor	-	714.64	12.38
	GaLore	2	549.27	19.12
	<b>COAP</b>	2	<b>447.59</b>	<b>12.30</b>

just one, making it a variant of SVD.

As shown in the Fig. 1, Tucker-2 achieves performance closest to the baseline across different scales of ResNet models. Thus, we select Tucker-2 as the primary format for computing convolution projection matrices.

## 2. Experimental Results

The COAP optimizer is versatile and can be applied to various model training scenarios. We have provided training results for DDPM and ControlNet, and this method is also suitable for other application scenarios [8–11].

### 2.1. Pre-training DDPM

**Experimental Settings.** We implement DDPM based on the Diffusers from Hugging Face and conduct experiments on  $8 \times V100$  GPUs following the training and evaluation settings in [2]. For CIFAR-10 [3], the model is trained with a

[https://github.com/huggingface/diffusers/tree/main/examples/unconditional\\_image\\_generation](https://github.com/huggingface/diffusers/tree/main/examples/unconditional_image_generation)

batch size of 128 for 800K steps. For CelebA-HQ [6], the model is trained with a batch size of 64 for 460K steps. We generate 50K images to compute the FID (Frechet Inception Distance) with respect to the training dataset and the images generated with 1000 DDPM steps.

**Comparison Results.** Table 2 presents the performance of our method and GaLore when compressing AdamW and Adafactor optimizers on the DDPM model. Our approach consistently outperforms GaLore on CIFAR-10 and CelebA-HQ datasets. Specifically, using the Adafactor optimizer, our method reduces FID by 1.7 and 6.8 compared to GaLore. Additionally, at compression rates of  $1.2\times$  and  $1.6\times$ , our approach surpasses the baseline performance.

## 2.2. Qualitative Comparisons

We present qualitative results, showcasing images generated by models trained with COAP and other optimizers. These results provide a visual comparison that complements the quantitative analysis in the main paper. Comparisons for DDPM (CIFAR-10), DDPM (CelebA-HQ), LDM, SiT-XL/2, and ControlNet-XL are detailed in Tables 3, 4, 5, 6, and 7, respectively.

Table 3. Comparison of images generated by DDPM trained on CIFAR-10 with different Adam-based optimizers.

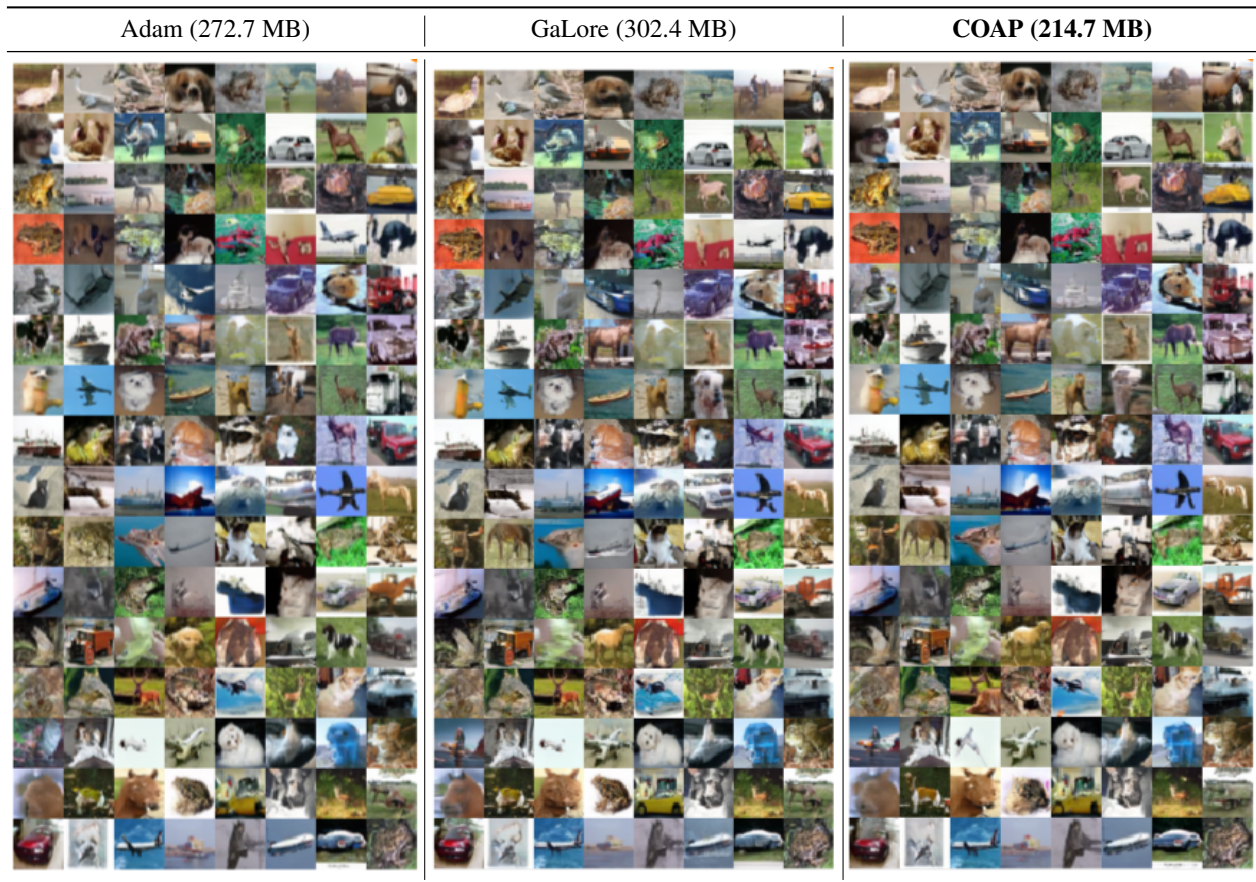




Table 4. Comparison of images generated by DDPM trained on CelebA-HQ with different Adafactor-based optimizers.



Table 5. Random class-conditional samples generated by LDM trained on the ImageNet dataset using the COAP optimizer.

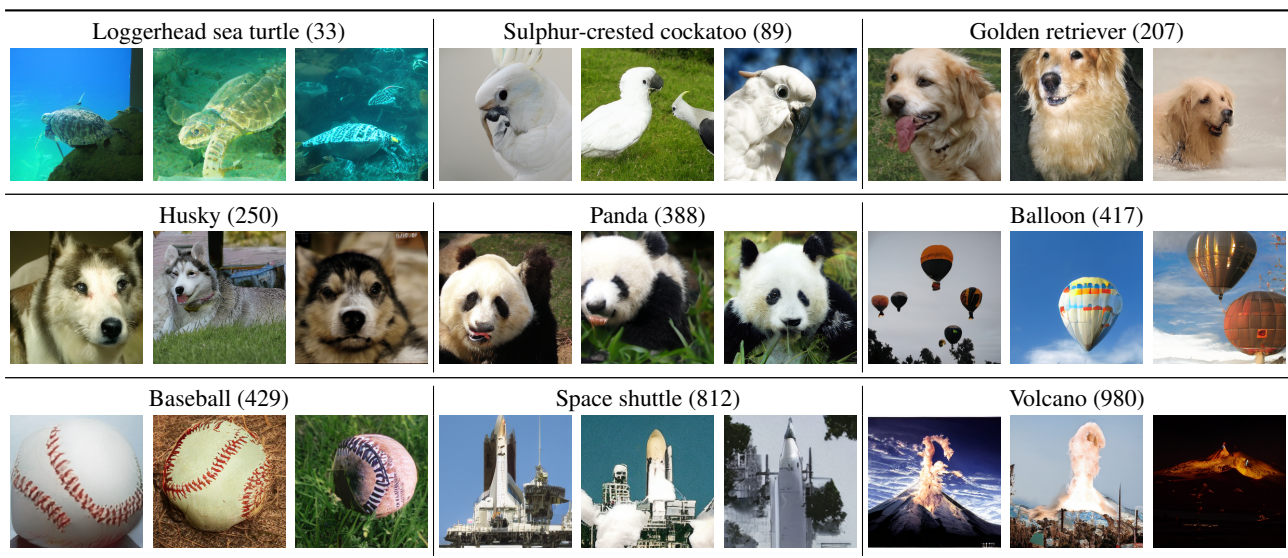




Table 6. Random class-conditional samples generated by SiT-XL/2 trained on the ImageNet dataset using the COAP optimizer.

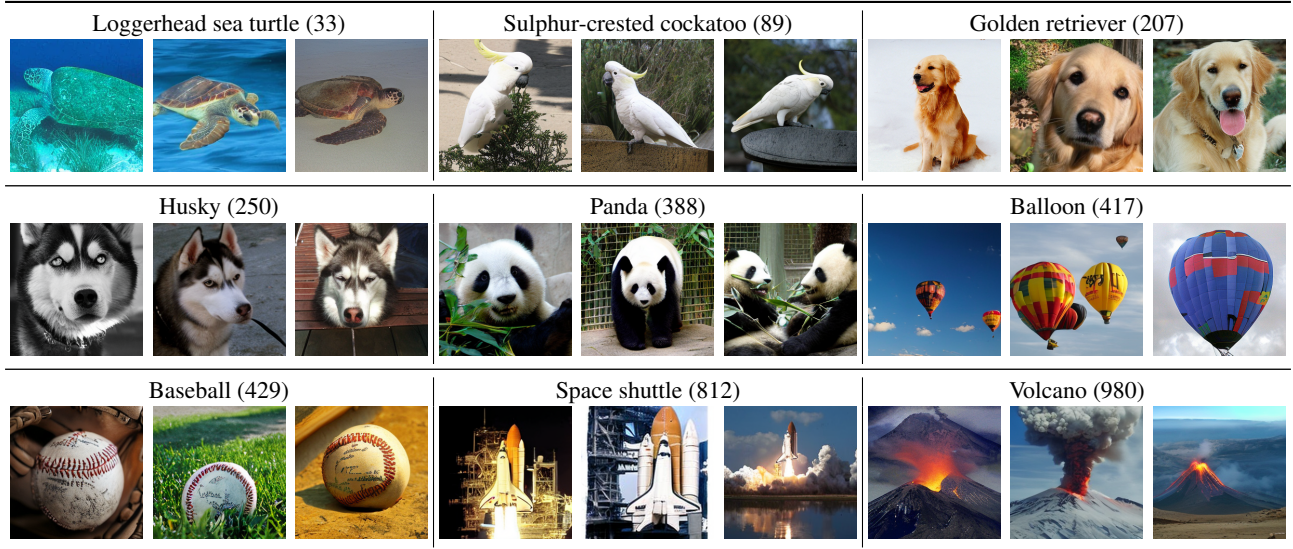
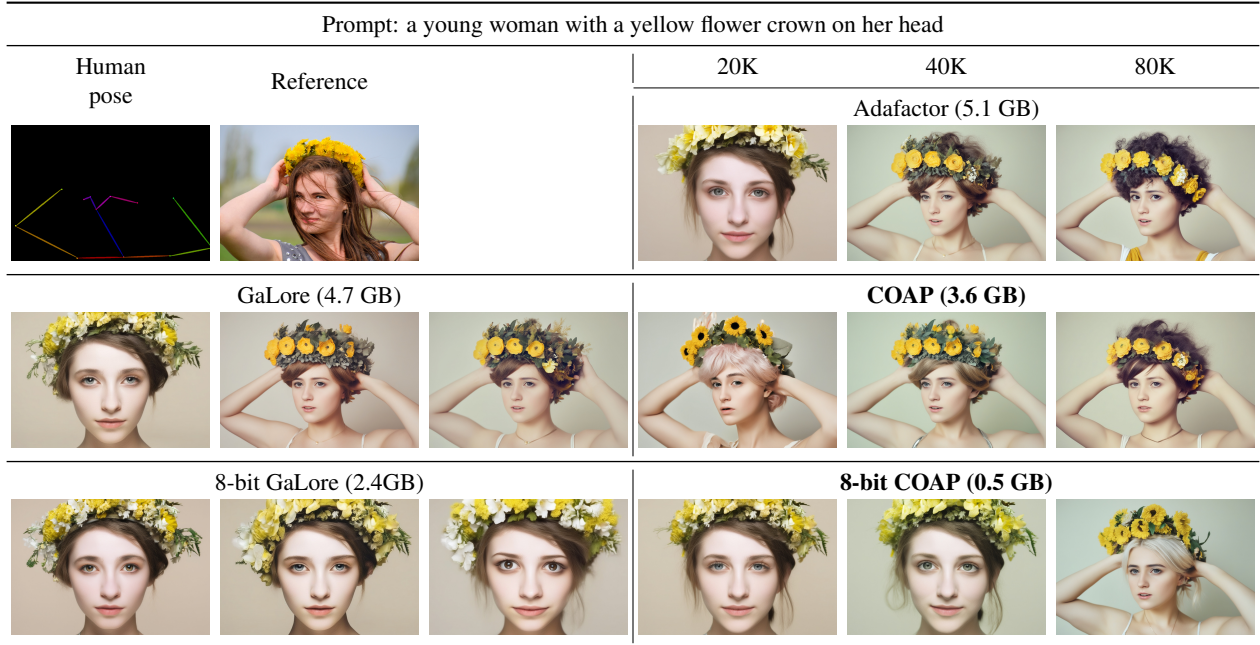
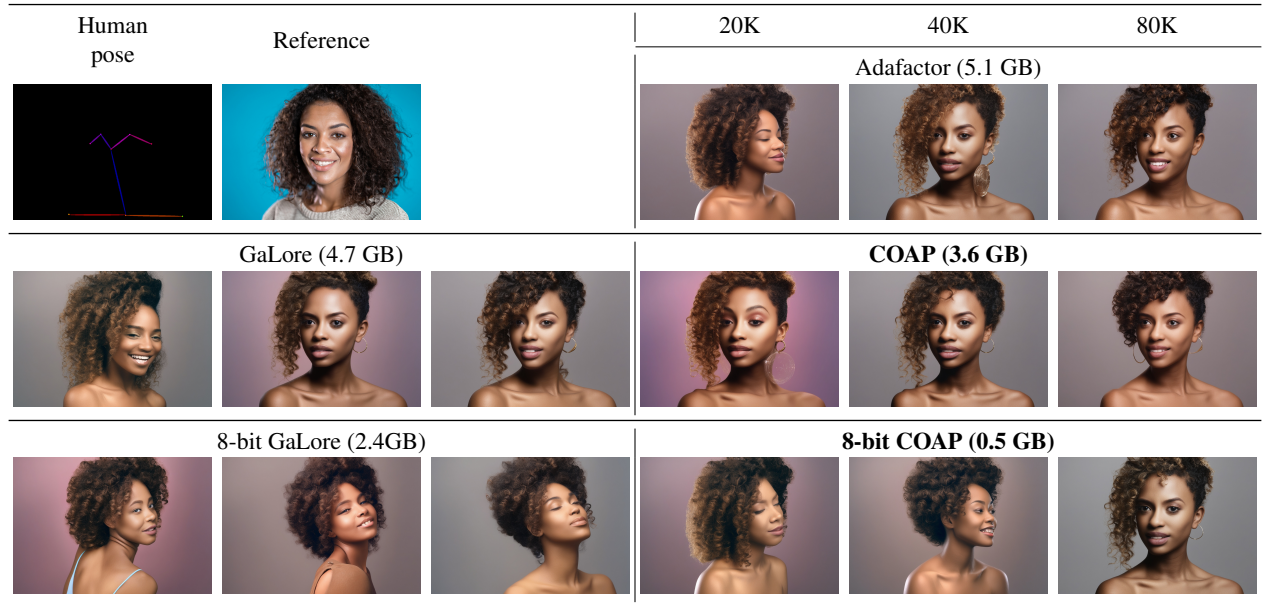


Table 7. Comparison of images generated at different training steps (20K, 40K, 80K) with ControlNet-XL trained under various optimizers. DDIM with a guidance scale of 5.0 is applied for image generation, with the number of inference steps set to 50.



Prompt: beautiful african american woman with curly hair on blue background

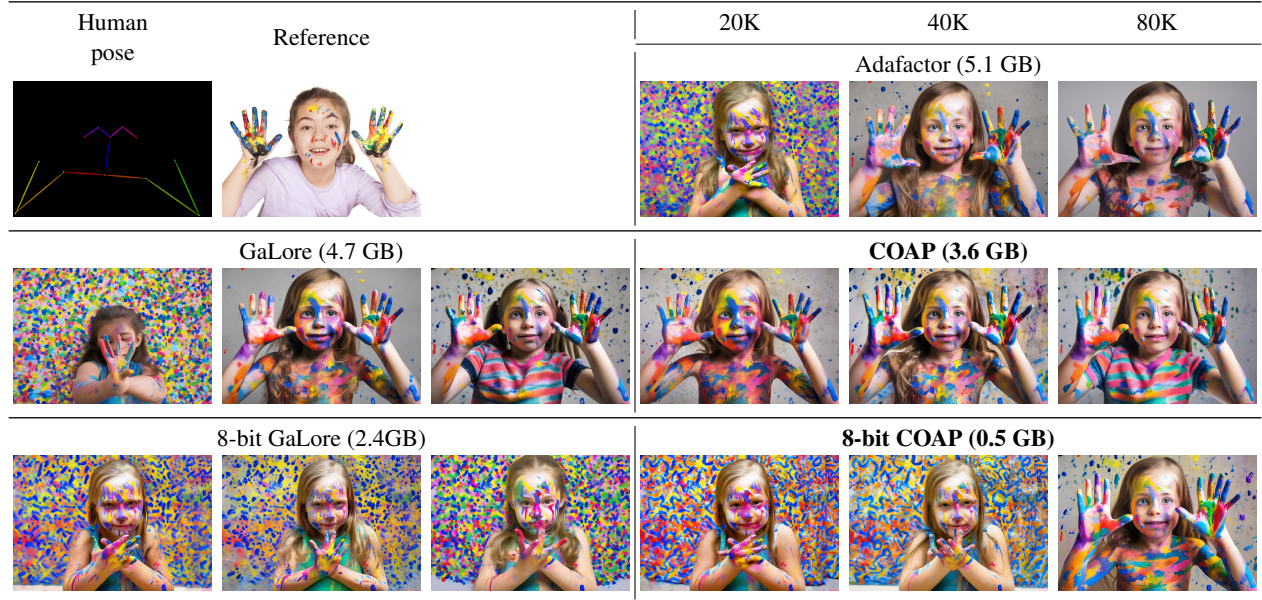


Prompt: a man in a scarf and sweater leaning against a brick wall





Prompt: a young girl with her hands painted with colorful paint



## References

- [1] Lieven De Lathauwer, Bart De Moor, and Joos Vandewalle. A multilinear singular value decomposition. *SIAM journal on Matrix Analysis and Applications*, 21(4):1253–1278, 2000. 3
- [2] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 4
- [3] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. 4
- [4] Pieter M Kroonenberg and Jan De Leeuw. Principal component analysis of three-mode data by means of alternating least squares algorithms. *Psychometrika*, 45:69–97, 1980. 3
- [5] Ji Liu, Przemyslaw Musialski, Peter Wonka, and Jieping Ye. Tensor completion for estimating missing values in visual data. *IEEE transactions on pattern analysis and machine intelligence*, 35(1):208–220, 2012. 3
- [6] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, 2015. 5
- [7] Ledyard R Tucker. Some mathematical notes on three-mode factor analysis. *Psychometrika*, 31(3):279–311, 1966. 3
- [8] Cheng Yang, Yang Sui, Jinqi Xiao, Lingyi Huang, Yu Gong, Yuanlin Duan, Wenqi Jia, Miao Yin, Yu Cheng, and Bo Yuan. Moe-i<sup>2</sup>: Compressing mixture of experts models through inter-expert pruning and intra-expert low-rank decomposition, 2024. 4
- [9] Yimeng Zhang, Xin Chen, Jinghan Jia, Yihua Zhang, Chongyu Fan, Jiancheng Liu, Mingyi Hong, Ke Ding, and Sijia Liu. Defensive unlearning with adversarial training for robust concept erasure in diffusion models. *Advances in Neural Information Processing Systems*, 37:36748–36776, 2024.
- [10] Yimeng Zhang, Jinghan Jia, Xin Chen, Aochuan Chen, Yihua Zhang, Jiancheng Liu, Ke Ding, and Sijia Liu. To generate or not? safety-driven unlearned diffusion models are still easy to generate unsafe images... for now. In *European Conference on Computer Vision*, pages 385–403. Springer, 2024.
- [11] Yimeng Zhang, Tiancheng Zhi, Jing Liu, Shen Sang, Liming Jiang, Qing Yan, Sijia Liu, and Linjie Luo. Id-patch: Robust id association for group photo personalization. *arXiv preprint arXiv:2411.13632*, 2024. 4