FROM BULK TO BUDGET: BEST PRACTICES TO COM PRESS MULTIMODAL LARGE LANGUAGE MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

Multimodal large language models (MLLMs) are increasingly developed to meet diverse deployment needs, varying in scale and computational demand. While recent research has focused on building MLLMs from Small Language Models (SLMs), these efforts remain limited in flexibility and are still data- and computeintensive. In this paper, we present the first comprehensive study on flexibly compressing and recovering existing MLLMs in a data-efficient manner. Hence, we address a critical gap in the literature by empirically analyzing best practices for adapting to specific hardware or resource limitations. Our study investigates pruning and knowledge distillation techniques, examining their impact on downstream performance across various model compression strategies, including pruning paradigms, recovery training schemes, and data requirements. Key findings reveal that widthwise pruning is particularly effective in resource-constrained scenarios. For smaller compression ratios, finetuning the multimodal projector alone can restore most performance, while combining finetuning with hidden state knowledge distillation proves most effective across all compression levels. Notably, we demonstrate efficient model downsizing using as little as 5% of the original dataset for moderate compression. Our analysis suggests best practices for compressing MLLMs for resource-efficient deployment. With our best practices, Bunny-v1.0-3B retains over 95% of its original performance, while LLaVA-v1.5-7B maintains more than 97%, with compression ratios below 30%.

029 030

004

010 011

012

013

014

015

016

017

018

019

021

023

025

026

027

028

031 032

033

1 INTRODUCTION

State-of-the-art multimodal large language models (MLLMs) (Liu et al., 2023; Chu et al., 2023;
Chen et al., 2024b) based on Large Language Models (LLMs) require substantial resources. For
instance, the LLaVA family (Liu et al., 2023) includes models with parameter counts ranging from 7
to 34 billion. Even those designed to be more memory-efficient, such as Bunny-v1.0-3B (He et al., 2024), still require significant storage, with 3 billion parameters. Reducing the size of these models
without compromising performance is crucial for adapting them to diverse deployment scenarios with varying resource constraints.

Despite the growing need for efficient MLLMs, most existing research has focused on building MLLMs on SLMs (Zhu et al., 2024a; He et al., 2024; Chu et al., 2023). While these approaches successfully reduce the overall model size, their flexibility is constrained by the fixed size of the underlying SLM. Furthermore, training an SLM from scratch to meet desired specifications is computationally expensive (Chu et al., 2023). Meanwhile, efforts to compress multimodal models have largely focused on task-specific tuning (Wang et al., 2023; Shi et al., 2023). To the best of our knowledge, no previous work has investigated general-purpose model compression for MLLMs.

In this work, we aim to uncover the key practices for obtaining effective and compressed MLLMs.
 We evaluate several techniques designed for compressing LLMs on MLLMs and investigate how
 different design choices for performance recovery affect the downstream performance of the compressed MLLMs. Our comprehensive empirical study explores several key dimensions: the pruning
 paradigms applied, the objectives used to restore initial performance, and the amount of data required
 for effective recovery. Specifically, we examine two MLLMs: a large-scale model (LLaVA-v1.5-7B
 (Liu et al., 2024)) and a model already optimized for efficiency (Bunny-v1.0-3B (He et al., 2024)).



061

063

064

067

068

069

071



Figure 1: Overview of our best practices for MLLM compression. After evaluating two pruning strategies—widthwise and layerwise pruning—and multiple recovery strategies, we propose distinct compression approaches for MLLMs. The decision flow highlights the scenarios for applying each method, depending on resource availability and compression ratio requirements. For each approach, 073 we display the resulting compressed model performance across a range of compression ratios (0-074 60%) on a set of multimodal benchmarks. The spider plots illustrate the retained performance across 075 different tasks, demonstrating each strategy's effectiveness at various compression levels. 076

077

090

092

093

095

096

097 098

099

102

103

078 We assess their performance on visual question answering (Hudson & Manning, 2019; Lu et al., 079 2022) and instruction-following tasks (Li et al., 2023; Yin et al., 2023).

Given that the LLM contains the majority of the parameters (95% for LLaVA-v1.5-7B, 86% for 081 Bunny-v1.0-3B), we compress it by applying two distinct pruning strategies: layerwise pruning, 082 which removes entire transformer layers, and widthwise pruning, which reduces the number of attention heads and MLP hidden dimensions. We then evaluate different strategies to recover the 084 potential performance loss: supervised finetuning, knowledge distillation from the original model 085 using Kullback-Leibler divergence or reversed Kullback-Leibler divergence on the logits distribution or L2 loss on the intermediate features. Finally, we investigate how to combine these losses and their 087 effectiveness w.r.t. the amount of available training data.

Our systematic evaluation across different compression ratios led to different key findings:

- Widthwise pruning is more effective in low-resource scenarios as it produces an efficient model even without recovery training.
- With recovery training, layerwise pruning is better for small ratios, while widthwise pruning usually outperforms it at larger ones (greater than 40%).
- Finetuning only the multimodal projector is sufficient at small compression ratios, where pruning has a minimal impact on the language model itself but destroys the multimodal alignment.
- The best recovery strategy is supervised finetuning coupled with intermediate representation distillation, consistently achieving the the highest performance across all compression ratios.
- The higher the pruning ratio, the higher the amount of data needed for recovering the **performance**. While with small ratios (less than 50%) even 5% of the data might suffice, this quantity increases for larger ones.
- We highlight our key findings in Figure 1. These findings consitute a set of best practices that 105 practitioners can follow when compressing MLLMs and researchers can consider for developing effective pruning strategies. To ease their exploitation and future studies, we will release our codebase, 107 benchmark as well as the compressed model checkpoints upon acceptance.

108 2 RELATED WORK

109 110

Pruning. Unstructured pruning (Dong et al., 2017; Frankle & Carbin, 2019; Lee et al., 2020; Park 111 et al., 2020; Sanh et al., 2020b; Farina et al., 2024) removes individual weights or neurons. While 112 such approaches can achieve strong compression rates with minimal accuracy trade-offs, they usu-113 ally require specialized hardware or software for effective acceleration. Structured pruning (Ding 114 et al., 2019; Li et al., 2017; Liu et al., 2021; You et al., 2019) eliminates entire groups of parameters 115 to reduce both the model's size and its computational overhead. Semi-structured pruning offers a 116 middle ground between structured and unstructured methods by selectively pruning certain model 117 parts. In the context of LLMs, Fang et al. (2023) and Ma et al. (2023) have successfully applied structured pruning, achieving significant sparsity with minimal performance degradation. Xia et al. 118 (2024) targets transformer layers and demonstrates that some layers can be pruned more aggres-119 sively without compromising accuracy. Meanwhile, Dery et al. (2024) propose a dynamic pruning 120 strategy that adjusts pruning throughout the training process. In this work, we focus on structure 121 pruning as well as the recovery strategies for MLLMs. 122

123 Further Compression Methods. In addition to pruning, techniques like quantization and low-rank factorization are also widely used for model compression. Quantization (Bai et al., 2021; Yao et al., 124 2022; Zafrir et al., 2019) reduces model size and computational cost by lowering the precision of 125 model parameters, enabling efficient inference with minimal performance loss. Low-rank factoriza-126 tion (Hsu et al., 2022; Hu et al., 2021b; Lan et al., 2020; Ashkboos et al., 2024) compresses models 127 by approximating large weight matrices through the product of smaller matrices, effectively reduc-128 ing the number of parameters while maintaining most of the model's capacity. While these methods 129 can offer significant compression, we focus on pruning techniques, which allow for more granular 130 control over the architecture by directly targeting and removing redundant components. 131

Knowledge distillation (KD) (Hinton et al., 2015) is a standard method for compressing LLMs by 132 transferring knowledge from a large teacher model to a smaller student model (Gou et al., 2021; 133 Sanh et al., 2020a). In NLP classification settings, KD is often applied by having the student model 134 replicate the teacher's output distribution (Liang et al., 2021; Song et al., 2020; Zhang et al., 2023), 135 hidden states Jiao et al. (2020); Sun et al. (2019), or attention patterns (Wang et al., 2020; 2021), 136 allowing the student to learn from the teacher's internal representations effectively. For text gener-137 ation tasks, Xu et al. (2024) provides a comprehensive survey of the role of knowledge distillation 138 in language models. Hsieh et al. (2023) introduce multi-stage distillation, transferring intermediate 139 representations to help the student model capture more detailed features. Gu et al. (2023) propose to 140 replace the forward Kullback-Leibler divergence with a reverse Kullback-Leibler divergence to prevent the student model from overestimating the low-probability regions of the teacher distribution. 141

142 Efficient MLLMs. Recent studies (Jin et al., 2024; Zhu et al., 2024b; Lin et al., 2024; Wei et al., 143 2024) have explored Multimodal Small Language Models (MSLMs). Models such as LLaVA-Phi 144 (Zhu et al., 2024b) utilize pretrained small language models to lower computational costs, while 145 MobileVLM (Chu et al., 2023) concentrates on projector designs to enhance MSLM performance. The Bunny model (He et al., 2024) explores the effects of training data size on performace. Al-146 though these approaches reduce model size, they are constrained by the fixed dimansions of the 147 base SLM. Our study specifically addresses methods for customizing the size of existing MLLMs 148 through structured pruning and recovery strategies. 149

While most multimodal structured compression efforts, such as EfficientVLM (Wang et al., 2023) and UPOP (Shi et al., 2023), focus on task-specific tuning for tasks like visual question answering and image captioning, general-purpose model compression for MLLMs remains underexplored.
Our work addresses this gap by investigating structured pruning techniques and recovery strategies applicable across a variety of multimodal tasks. Unlike previous approaches optimizing models for specific tasks, our study provides general-purpose compression guidelines for MLLMs.

156 157

3 Methodology

158 159

This section outlines our approach to compressing MLLMs. We first introduce two pruning strate gies: layerwise and widthwise pruning. We then describe methods to recover model performance through supervised finetuning and knowledge distillation.

Notation. Given a triplet $\mathbf{X} = {\mathbf{x}_v, \mathbf{x}_p, \mathbf{x}_r}$, the objective of an MLLM m_θ , parameterized by $\theta = {\psi, \phi, \mathbf{W}}$, is to generate a response \mathbf{x}_r based on an input image \mathbf{x}_v and a text prompt \mathbf{x}_p , such that $m_\theta(\mathbf{x}_v, \mathbf{x}_p) = \mathbf{x}_r$. The MLLM typically consists of a vision encoder $g_{\psi}(\cdot)$, an LLM $f_{\phi}(\cdot)$, and a multimodal projector \mathbf{W} aligning the two modalities. The prompt \mathbf{x}_p is tokenized into \mathbf{T}_p , while the vision encoder processes the image \mathbf{x}_v to extract visual features, which are then converted into language embedding tokens \mathbf{T}_v via the multimodal projector:

$$\mathbf{\Gamma}_{v} = \mathbf{W} \cdot g_{\psi}(\mathbf{x}_{v}) \text{ and } f_{\phi}(\mathbf{T}_{v} \odot \mathbf{T}_{p}) = \mathbf{x}_{r}.$$
 (1)

170 The concatenated visual tokens \mathbf{T}_v and prompt tokens \mathbf{T}_p are fed into the LLM's *M* layers, produc-171 ing hidden states $\{\mathbf{H}_i \in \mathbb{R}^{T \times d}\}_{i=1}^{M}$, where *T* is the number of tokens and *d* is the hidden dimension. 172 Finally, the probabilities $p_{m_\theta}(\mathbf{x}_r | \mathbf{x}_v, \mathbf{x}_p, \tau)$ are computed by passing the final hidden state through 173 the classification head with softmax temperature τ .

3.1 PRUNING

168 169

174 175

196

197

201

214

Pruning seeks to reduce the number of parameters in a model, thus decreasing its computational cost. In MLLMs, the majority of parameters θ are concentrated in the LLM f_{ϕ} , so downsizing it can significantly reduce the overall computational burden. The LLM is typically structured along two main dimensions: depth and width. Depth refers to the number of stacked transformer layers, while width pertains to the internal structure of each layer, including the multi-head attention mechanism and the multi-layer perceptron (MLP). In this paper, we explore two pruning strategies to reduce the parameter count in LLMs: layerwise pruning, which removes entire transformer layers, and widthwise pruning, which eliminates the least important components within each layer.

To determine which layers or components to prune, we draw a small subset of n samples from the the original visual instruct-tuning dataset as the calibration dataset $\mathcal{D} = {\{\mathbf{x}_v^j, \mathbf{x}_p^j, \mathbf{x}_r^j\}_{j=1}^n}$. The importance of each layer or component is assessed, and those with the lowest importance are pruned.

188 3.1.1 LAYERWISE PRUNING

Empirical research (Fan et al., 2019; Sajjad et al., 2023) has shown that large transformer models often contain redundant layers, allowing several to be removed with minimal impact on accuracy. To identify and remove these redundant layers, we use the Block Influence (BI) score (Men et al., 2024), which quantifies the importance of layer *i* through the cosine distance between input H_i and output hidden states H_{i+1} . The key assumption is that layers that cause larger changes in hidden states have a greater influence on model performance. The BI score of layer *i* is then calculated by

$$\mathbf{B}\mathbf{I}_{i}(\mathcal{D}) = 1 - \mathbb{E}_{\mathbf{X}\sim\mathcal{D},t} \left[\frac{\mathbf{H}_{i,t}^{\mathsf{T}}\mathbf{H}_{i+1,t}}{\|\mathbf{H}_{i,t}\|_{2}\|\mathbf{H}_{i+1,t}\|_{2}} \right],\tag{2}$$

where $\mathbf{H}_{i,t}$ represents the t^{th} row of \mathbf{H}_i . After calculating the BI scores, the layers are ranked by importance, and those with the lowest scores are pruned.

202 3.1.2 WIDTHWISE PRUNING

Previous research has shown that transformer layers also exhibit width redundancy, meaning that
only a subset of attention heads (Voita et al., 2019; Michel et al., 2019) or MLP dimensions (McCarley et al., 2019; Hudson & Manning, 2019) are critical for model performance. To address this,
we apply dependency-based structural pruning, which removes redundant widthwise components
while minimizing the impact on the model's performance. Specifically, we first identify groups of
interdependent structures and then prune entire groups based on their collective importance.

Following the methods of Fang et al. (2023) and Ma et al. (2023), we begin by constructing dependency relationships within each LLM layer. Let N_i and N_j represent two neurons in the layer, where In(N_i) and Out(N_i) represent the neurons connected to N_i as inputs and outputs, respectively. The dependency of neuron N_j on N_i is defined as:

$$N_j \in \operatorname{Out}(N_i) \cap \operatorname{Num}_{\operatorname{In}(N_j)} = 1, \text{ or } N_j \in \operatorname{In}(N_i) \cap \operatorname{Num}_{\operatorname{Out}(N_j)} = 1,$$
(3)

where $\operatorname{Num}_{\operatorname{In}(N_i)}$ refers to the number of nodes connected to N_i as inputs and $\operatorname{Num}_{\operatorname{Out}(N_j)}$ is the number of nodes connected to N_j as outputs. If neuron N_i is pruned, all its dependent neurons N_j

must also be pruned. This process results in a set of dependency graphs $\mathcal{G} = \{w_i^k\}_{i=1}^M$, where M is the number of structures in the graph and w_i^k represents the k^{th} weight parameter within a structure.

Once the dependency graphs are constructed, we assess their importance at the group level, since all weights within a graph must be pruned together. Group importance is evaluated by comparing the vision-language modeling loss $\mathcal{L}_{CE}(m_{\theta}(\mathbf{x}_{v}, \mathbf{x}_{q}), \mathbf{x}_{r})$ on the calibration dataset, both with and without the weight. To efficiently approximate the importance, we apply a Taylor expansion using gradient information. The importance function is given by:

$$I_{w_i^k}(\mathbf{X}) = |\mathcal{L}_{CE}(\mathbf{X}, m_\theta) - \mathcal{L}_{CE}(\mathbf{X}, m_\theta^{w_i^k = 0})| \approx \left| \frac{\partial \mathcal{L}_{CE}(\mathbf{X}, m_\theta)}{\partial w_i^k} w_i^k \right| .$$
(4)

We then prune the graphs with the lowest group importance $I_{\mathcal{G}}$:

$$I_{\mathcal{G}}(\mathcal{D}) = \mathbb{E}_{\mathbf{X} \sim \mathcal{D}} \left[\sum_{i}^{M} \sum_{k} I_{w_{i}^{k}}(\mathbf{X}) \right].$$
(5)

3.2 RECOVERY TRAINING

232

233

234

235

236 237

238

246

247 248

249

255

258 259 Pruning a large multimodal language model results in performance degradation, affecting both language modeling and cross-modality alignment. To mitigate this, we investigate two recovery training methods: supervised finetuning (Sec. 3.2.1) and knowledge distillation (Sec. 3.2.2). We consider the orginal teacher model m_{θ}^{T} , the pruned student model $m_{\theta'}^{S}$, and a recovery dataset \mathcal{D} .

3.2.1 Recovery Training with supervised finetuning

A simple yet effective approach to recovery training is supervised finetuning on the original dataset. This method helps counteract performance degradation by allowing the model to adapt its parameters to the modified architecture while taking advantage of the detailed annotations in the original dataset. Here, we first focus on training only the multimodal projector to realign the vision and language spaces. Second, we jointly finetune both the projector and the language model while keeping the vision encoder fixed, as finetuning the vision encoder does not improve performance (Karamcheti et al., 2024). We use the cross-entropy loss for supervised finetuning, denoted as

$$\mathcal{L}_{sft}(m_{\theta'}^{\mathsf{S}}, \mathcal{D}) = \mathbb{E}_{\mathbf{X} \sim \mathcal{D}}[\mathcal{L}_{CE}(m_{\theta'}^{\mathsf{S}}(\mathbf{x}_{v}, \mathbf{x}_{p}), \mathbf{x}_{r})].$$

(6)

3.2.2 Recovery Training with knowledge distillation

Knowledge distillation (KD) is a method used to transfer knowledge from a large, well-trained model (the teacher) to a smaller or pruned model (the student) (Hinton et al., 2015). This approach allows
the pruned model to regain lost performance by mimicking the decision-making process of the more capable teacher. In our setup, the uncompressed model acts as the teacher, while the pruned model serves as the student. We explore two main strategies, logits-based KD and hidden state based KD, evaluating different loss functions and their trade-offs.

Logits-based KD focuses on aligning the output probability distributions of the pruned model with those of the teacher model. The logits-based KD loss is defined as

$$\mathcal{L}_{logits}(m_{\theta'}^{\mathbf{S}}, m_{\theta}^{\mathbf{T}}, \mathcal{D}) = \mathbb{E}_{\mathbf{X}\sim\mathcal{D}} \left[\mathcal{L}_{KD}(p_{m_{\theta}^{\mathbf{T}}}(\mathbf{x}_{r} | \mathbf{x}_{v}, \mathbf{x}_{p}, \tau), p_{m_{\theta'}^{\mathbf{S}}}(\mathbf{x}_{r} | \mathbf{x}_{v}, \mathbf{x}_{p}, \tau)) \right].$$
(7)

We leverage two distinct KD losses to evaluate the differences between these logit distributions. Given the teacher distribution p_{θ} and the student distribution $q_{\theta'}$, the standard KD objective minimizes the approximated forward Kullback–Leibler (KL) divergence between these two distributions, denoted as $\mathcal{L}_{KD}(p_{\theta}, q_{\theta'}) = D_{KL}(p_{\theta} || q_{\theta'})$. This approach encourages the student distribution to match all the modes of the teacher distribution.

However, minimizing forward KL can lead q_{θ} to assign excessively high probabilities to areas where p has little or no probability mass (Malinin & Gales, 2019). In contrast, Reversed Kullback–Leibler divergence (RKL) minimizes $\mathcal{L}_{KD}(p_{\theta}, q_{\theta'}) = D_{KL}(q_{\theta'} || p_{\theta})$, encouraging $q_{\theta'}$ to focus on the major modes of p_{θ} while assigning low probabilities to its less significant regions. This helps the student model avoid learning unnecessary long-tail variations of the teacher distribution and instead focus on generating more accurate responses (Gu et al., 2023; Holtzman et al., 2019). 270 Hidden State Matching involves aligning the pruned model's intermediate representations (hidden 271 states) $\mathbf{H}_{i}^{m_{\theta'}^{S}}$ with the teacher model's $\mathbf{H}_{i}^{m_{\theta}^{T}}$. The corresponding loss for a layer *i* can be defined as 272

$$\mathcal{L}_{match}(m_{\theta'}^{\mathsf{S}}, m_{\theta}^{\mathsf{T}}, \mathcal{D}) = \mathbb{E}_{\mathbf{X} \sim \mathcal{D}} \left[\mathcal{L}_{feat}(\mathbf{H}_{i}^{m_{\theta'}^{\mathsf{S}}}, \mathbf{H}_{i}^{m_{\theta}^{\mathsf{T}}}) \right],$$
(8)

where \mathcal{L}_{feat} refers to a feature matching loss. Both Yang et al. (2024) and Popp et al. (2024) 276 suggest that applying a feature-based L2 distillation loss improves the student model's performance, 277 particularly for pre-trained vision-language models. Consequently, we employ L2 loss as the feature 278 matching loss $\mathcal{L}_{feat} = \|\cdot - \cdot\|_2^2$.

The total loss for recovery training is computed as:

$$\mathcal{L}(m_{\theta}^{\mathsf{S}}, m_{\theta}^{\mathsf{T}}, \mathcal{D}) = \alpha \mathcal{L}_{sft}(m_{\theta'}^{\mathsf{S}}, \mathcal{D}) + \beta \mathcal{L}_{logits}(m_{\theta'}^{\mathsf{S}}, m_{\theta}^{\mathsf{T}}, \mathcal{D}) + \gamma \mathcal{L}_{match}(m_{\theta'}^{\mathsf{S}}, m_{\theta}^{\mathsf{T}}, \mathcal{D})$$
(9)

where α , β , and γ are the coefficients that balance the contributions of the different loss components.

4 **EXPERIMENTS**

In this section, we first introduce our experimental setup and then demonstrate the main findings on 288 model pruning (Sec. 4.1) and performance recovery (Sec. 4.2, Sec. 4.3). Finally, we highlight the findings on recovery training using only a small fraction of data (Sec. 4.4).

291 **Experimental setup.** We evaluate pruning and knowledge distillation strategies on both a large-292 scale MLLM model (LLaVA-v1.5-7B (LLaVA) (Liu et al., 2024)) and a smaller-scale MLLM model 293 (Bunny-v1.0-3B (Bunny) (He et al., 2024)). LLaVA is built upon Vicuna-v1.5 (Chiang et al., 2023) 294 with 6.7 billion parameters, and Bunny is based upon Phi-2 (Javaheripi et al., 2023) with 2.8 billion parameters. We provide a detailed overview of the model architectures in Appendix A. For both 295 models, we exclusively use their visual instruction tuning datasets: LLaVA-v1-5-mix665k (Liu et al., 296 2024) for LLaVA and Bunny-695K (He et al., 2024) for Bunny. During pruning, we randomly 297 select 10 samples as the calibration dataset to compute layer importance. For recovery training, 298 we experiment with various portions of the original dataset (i.e. 5%, 10%, 20%, and 100%) for 299 finetuning and knowledge distillation. We set the distillation temperature to 2.0 for logits-based 300 distillation and use the final layer representation for hidden state matching (see Appendix B). 301

We evaluate the pruned and recovery-trained models on visual question-answering tasks using GQA 302 (Hudson & Manning, 2019) and SQA-I (Lu et al., 2022), as well as instruction-following tasks 303 with POPE (Li et al., 2023), MME-Cognition, MME-Perception (Yin et al., 2023) and MMMU 304 (Yue et al., 2024). To ensure consistency, we use the lmms-eval suite (Bo et al., 2024) for all 305 evaluations. For clearer comparisons, we calculate the relative performance as a percentage of the 306 original (uncompressed) model's performance on each benchmark. 307

308 309

273 274 275

279

284 285

286 287

289

290

4.1 The effect of pruning on the model performance and resources usage

310 In this section, we begin by exploring the techniques to obtain the best pruned model. We then 311 examine the resulting reductions in memory usage and computational requirements. 312

313 **Comparison of Pruning Techniques.** As illustrated in Figure 2 (blue lines), for both the Bunny 314 and LLaVA model, widthwise pruning consistently outperforms layerwise pruning in terms of model performance after pruning without recovery training. Specifically, for small compression ratios, 315 such as 15%, the Bunny model retains 95% of its performance, while LLaVA retains 93% (see Ap-316 pendix C for full results). In resource-constrained scenarios, widthwise pruning without recovery 317 training offers an efficient strategy when a small compression ratio is required. However, as the com-318 pression ratio increases, both widthwise and layerwise pruned models show significant performance 319 degradation. Overall, widthwise pruning better preserves the model's structure and information flow, 320 allowing it to keep performance with minimal adjustments, especially at lower compression ratios. 321

The impact of the pruning method on model performance after finetuning both the projector and 322 the LLM is also illustrated in Figure 2 (green lines). For smaller compression ratios (less than 323 40%), layerwise pruning offers a slight advantage over widthwise pruning, while widthwise pruning



333

334

335

336

337

338

Figure 2: Comparison of pruning and finetuning strategies on two MLLMs, Bunny and LLaVA. The plot shows the average relative performance under three scenarios: pruning only, pruning followed by finetuning the projector, and pruning followed by finetuning both the projector and the LLM. For smaller compression ratios, finetuning only the projector effectively recovers performance. For larger compression ratios, finetuning the projector and the LLM leads to better recovery, indicating the need for broader adjustments as more parameters are pruned.

339 delivers better overall performance for larger compression ratios (greater than 40%). This suggests 340 that finetuning plays a crucial role in reconstructing inter-layer connections and reoptimizing layer 341 components. 342

Best Practice for MMLM Pruning. Widthwise pruning generally proves more effective than layer-343 wise pruning in obtaining the best pruned model. A widthwise pruned model can often be deployed 344 without recovery training when targeting a small compression ratio (less than 20%). Regarding 345 post-finetuning performance, layerwise pruning shows a slight advantage at compression ratios be-346 low 30%, whereas widthwise pruning performs marginally better at higher compression ratios. 347

348 From compression ratio to resource usage. 349 Table 1 provides an overview of how different 350 compression ratios impact memory usage and 351 FLOPS for both the Bunny and LLaVA mod-352 els compressed via widthwise pruning. Mem-353 ory consumption refers to the allocated GPU memory, while FLOPS are measured using the 354 Calflops codebase¹. The results demonstrate 355 that higher compression ratios consistently lead 356 to both memory and compute reductions. For 357 example, at a 30% compression ratio, we ob-358 serve a memory reduction of 25% for Bunny 359 and 28% for LLaVA, with a corresponding de-

Table 1: Memory requirements (Mem.) and FLOPS for the Bunny and LLaVA models at various compression ratios. The models are pruned widthwise. The evaluation is performed in inference mode, where each model is provided with an image and a prompt containing 50 tokens.

Datia	Bun	iny	LLaVA			
Ratio	Mem. (MiB)	FLOPS (T)	Mem. (MiB)	FLOPS (T)		
0%	6,167	4.77	13,546	9.57		
15%	5,380	4.14	11,530	8.21		
30%	4,597	3.50	9,548	6.89		
45%	3,770	2.84	7,470	5.49		
60%	2,992	2.20	5,435	4.17		

360 crease in FLOPS of 27% for both models. These reductions continue to scale with larger com-361 pression ratios; at a 60% compression ratio, memory usage and FLOPS decrease by 50-60%. We 362 observe similar results for layerwise pruning (see Appendix D). This indicates that the achieved 363 compressions directly translate into improvements in memory efficiency and computational cost.

364 365 366

4.2 SUPERVISED FINETUNING FOR PERFORMANCE RECOVERY AFTER PRUNING

Compressing the LLM can impair its language modeling capabilities. Additionally, whether pruning 367 the LLM decoder would disrupt the alignment between vision and language remains underexplored. 368 To explore these issues, we experiment with two approaches: finetuning only the multimodal projec-369 tor and jointly finetuning both the projector and the LLM. Following the previous research (Karam-370 cheti et al., 2024), which shows that training the vision encoder may degrade overall model perfor-371 mance, we keep the vision encoder frozen in both setups. To facilitate fast recovery, we employ the 372 low-rank approximation, LoRA (Hu et al., 2021a), while finetuning the LLM. 373

374 **Finetuning the multimodal projector.** As shown in Figure 2 (orange lines), finetuning the mul-375 timodal projector significantly restores performance. At lower compression ratios (less than 20%), 376 finetuning only the projector achieves results comparable to jointly finetuning the LLM. For both 377

¹Caflops codebase: https://github.com/MrYxJ/calculate-flops.pytorch



Figure 3: Comparison of different distillation recovery strategies (KL loss, RKL loss, L2 loss, and their combinations) for Bunny and LlaVA models pruned with widthwise pruning. The plot shows the relative performance improvement of each strategy over standard finetuning across various compression ratios. The results demonstrate that distillation helps recover more performance than fine-tuning alone, with the L2 loss component consistently leading to the largest performance gains.

Bunny and LLaVA, finetuning the projector retains at least 95% of the performance at a compression ratio of 15%. As the compression ratio increases, the loss of language modeling ability becomes more pronounced, making projector-only finetuning insufficient to recover the model's performance fully. Nevertheless, even at a compression ratio of 60%, only finetuning the multimodal projector can still recover 60 to 80% of the performance by realigning the vision and language inputs. This shows that pruning specific LLM structures in the MLLM can both damage the language modeling ability and introduce modality misalignment, making the model incapable of comprehending vision.

400 401

388

389

390

391

392 393

394

395

396

397

398

399

402 Finetuning both the projector and the LLM. While a significant portion of the recovered per-403 formance is due to realigning the visual and textual inputs, we observe consistent gains from addi-404 tionally finetuning the pruned LLM (green lines in Figure 2), especially at higher compression ratios (greater than 40%). This indicates that the pruned model not only suffers from modality misalign-405 ment but also experiences a decline in its language modeling capabilities. We can partly restore these 406 lost capabilities by finetuning the LLM. At a compression ratio of 40%, finetuning both the projector 407 and the LLM restores more than 80% of the original model's performance. Even at a compression 408 ratio of 60%, finetuning recovers close to 80% of the model's original performance. 409

Best Practice for Supervised Finetuning. When a small compression ratio of around 15% is
 required, finetuning the multimodal projector alone is typically sufficient to recover most of the
 model's performance. For higher compression ratios (greater than 40%), incorporating finetuning of
 the LLM yields additional performance improvements.

414

4.3 KNOWLEDGE DISTILLATION FOR PERFORMANCE RECOVERY AFTER PRUNING

This section investigates the impact of combining knowledge distillation with finetuning.
We perform ablation studies by adjusting the weights of each loss component to evaluate their individual contributions (see 3.2.2). Additional ablation results are in the Appendix E.

As shown in Table 2 for the Bunny model compressed with layerwise pruning, we compare a logit-based approach (RKL) and a hidden state matching strategy (L2), both with and without a finetuning loss component. The results demon-

Table 2: Comparison of distillation strategies with and without finetuning for the Bunny model compressed via layerwise pruning. Finetuning helps stabilize performance and prevents *model collapse*, especially at higher compression ratios.

Ratio	FT	L2	Bunny L2+FT	RKL	RKL+FT
15%	96.30%	95.51%	99.59%	96.88%	98.70%
30%	94.33%	88.13%	95.03%	92.21%	93.81%
45%	86.70%	56.96%	90.19%	82.57%	88.50%
60%	69.38%	47.61%	72.62%	12.61%	69.85%

strate that incorporating a supervised finetuning loss significantly enhances and stabilizes distillation performance. For example, when applying only the distillation loss, the L2 and RKL methods
can recover 85% of the original performance at a compression ratio of 30%. However, for higher
compression ratios, adding the finetuning loss becomes critical in preventing model collapse. At
a compression ratio of 60%, combining the finetuning loss with distillation dramatically improves
performance—L2 distillation increases from 47.61% to 72.62% and RKL distillation from 12.61%



Figure 4: Comparison of recovery performance using different percentages of training data (100%, 20%, 10%, and 5%) for finetuning and distillation after pruning across Bunny and Llava models. For smaller compression ratios, even a small percentage of the training data (as low as 5%) is sufficient to recover most of the original performance. However, as the compression ratio increases, more training data is required to achieve higher recovery performance.

to 69.85%. This pattern is consistently observed across all models and compression techniques
 evaluated. While knowledge distillation alone can partially recover performance after pruning, its
 effectiveness is limited without the integration of finetuning.

451 Figure 3 compares various distillation strategies based on their relative improvement over finetuning 452 alone when widthiwse pruning is applied (see Appendix F for further results on layerwise pruning). Our results indicate that applying the L2 loss to align the hidden states of the student and teacher 453 in the final layer yields the best performance, or at least matches other methods. Unlike logit-454 based approaches, which require the student to replicate the teacher's output distribution, the L2 loss 455 method enables the student to better capture the teacher's feature representations directly, leading to 456 enhanced performance. Additionally, we observe that RKL generally outperforms KL across most 457 compression ratios, a result consistent with the findings of Gu et al. (2024). 458

Best Practice for Knowledge Distillation. Knowledge distillation, particularly when combined with
 finetuning and using L2 loss to map the intermediate states, delivers the most effective performance
 recovery after pruning across all compression ratios.

463 4.4 DATA EFFICIENT RECOVERY

442

443

444

445

446 447

461 462

In this section, we investigate the feasibility of performing recovery training using only a small fraction of the available data. Figure 4 shows the models' performance after recovery training with different portions of the original dataset relative to training with the full 100%. Both models undergo widthwise pruning and recovery training incorporating RKL and L2 loss functions. Remarkably, for

Table 3: Performance of the best compressed models. The size is the number of total parameters of the model, while the ratio, short for compression ratio, indicates the proportion of remaining LLM parameters compared to the pre-pruning state. When the compression ratio (Ratio) is below 40%, we apply depthwise pruning. For ratios above 40%, we use widthwise pruning. During the recovery phase, we employ supervised finetuning combined with L2 loss to match the hidden states. For both Bunny and LLaVA, 95% performance is retained if the compression ratio is smaller than 40%.

Method	Size	Ratio	MMMU	GQA	SQA	MME-C	MME-P	POPE	AVG	AVG-%
				Bur	ny-v1.0	-3B				
	3.2B	0%	34.10	54.72	70.70	289.30	1487.71	87.82	59.65	100.00%
Depth+FT+L2	2.8B	15%	33.00	54.56	70.00	304.29	1457.06	87.97	59.40	99.59%
Depth+FT+L2	2.5B	30%	32.30	53.08	68.12	252.50	1349.91	87.53	56.68	95.03%
Width+FT+L2	2.0B	45%	29.10	52.31	63.06	244.64	1281.66	87.09	54.37	91.15%
Width+FT+L2	1.6B	60%	28.10	48.72	53.20	216.07	1115.33	86.73	49.92	83.69%
				LLa	VA-v1.5	-7B				
	7.0B	0%	35.10	61.98	68.67	363.21	1511.33	86.99	62.28	100.00%
Depth+FT+L2	6.3B	15%	36.40	61.20	68.42	337.86	1442.35	86.94	61.22	98.29%
Depth+FT+L2	5.5B	30%	36.00	60.34	68.82	318.57	1496.60	85.98	60.96	97.88%
Width+FT+L2	3.8B	45%	30.80	57.74	52.90	215.00	1191.17	85.74	52.27	83.92%
Width+FT+L2	2.8B	60%	27.70	52.32	46.26	211.79	1085.97	84.06	48.52	77.90%

486 compression ratios below 50%, using just 5% of the original data is sufficient to achieve over 95% 487 of the performance compared to using the full dataset. However, as the compression ratio increases, 488 the amount of data required for effective recovery training also grows. For a compression ratio of 489 60%, the relative performance drops below 90% for LLaVA and diminishes even further to below 490 70% for Bunny. Nevertheless, using only a small portion of the training data appears to be a valid option, significantly lowering the required time and cost for compressing and finetuning MLLMs. 491

492 Best Practice for Data Efficient Recovery. At small to medium compression ratios less than 50%, 493 using just 5% of the dataset is enough to achieve performance comparable to full data training. 494 However, for compression ratios greater than 50%, full data training becomes necessary to recover 495 performance effectively.

496 497

498

499

500

501 502

507

511

513

514

515 516

517

518

4.5 MODEL COMPRESSION RESULTS FOLLOWING OUR BEST PRACTICES

In this section, we summarize our key findings as a set of best practices and highlight model performance achieved by following them. Based on the empirical results from the previous section, we outline the following best practices for compressing MLLMs:

Best Practices for MLLM Compression and Recovery Widthwise pruning is more effective in low-resource settings, yielding an efficient model 504 even without the need for recovery training. 505 506 • With recovery training, layerwise pruning excels for smaller compression ratios (below 40%), while widthwise pruning performs better at higher ratios (above 40%). • For small compression ratios, fine-tuning just the multimodal projector is often sufficient 509 to restore performance, with minimal impact from pruning. 510 • For recovery training, combining finetuning with knowledge distillation of the intermediate representations using L2 loss consistently achieves the highest performance across all 512

compression ratios. Data efficiency can be significantly boosted, requiring only 5% of the original data to match full-data training results, though full datasets are still needed for high compression ratios.

These guidelines provide insights for researchers aiming to develop new techniques for deploying MLLMs, enabling more effective model customization for specific deployment needs.

To illustrate the performance at different compression ratios, Table 3 offers a detailed comparison 519 of results for both Bunny and LLaVA across various multimodal benchmarks. The results show 520 that, with compression ratios below 30%, Bunny retains over 95% of its original performance, while 521 LLaVA maintains more than 97%. Even at higher compression ratios, up to 60%, our best practices 522 preserve an average performance of 83% for Bunny and 78% for LLaVA. These findings underscore 523 the feasibility of compressing MLLMs without incurring significant performance degradation. We 524 provide qualitative results for the compressed models in Appendix G. Additionally, other compres-525 sion techniques, such as quantization, could be layered on top of the current framework to further 526 reduce inference time and memory usage, see Appendix H. We also include InternVL (Chen et al., 527 2024a) in Appendix I to show our best practices generalize to other model.

528 529

530

5 CONCLUSION

531 In this work, we investigated efficient compression techniques for MLLMs, focusing on two key 532 pruning strategies: width and layerwise. Our study assessed the impact of these strategies on model 533 performance, both before and after recovery training, across various compression ratios. We further 534 explored recovery methods such as supervised finetuning and knowledge distillation to address per-535 formance degradation caused by pruning. We formulate our findings as best practices, which offer practical guidelines for optimizing MLLMs, enabling a balance between model size, performance, 536 and data efficiency to meet specific deployment resource constraints. Due to computational con-537 straints, in this work we focused on two representative pruning techniques on two different models. 538 Future work could extend these best practices to include a broader range of pruning techniques and models to further refine these strategies.

540 REPRODUCIBILITY STATEMENT

We have made extensive efforts to ensure the reproducibility of our results. All pruning and recovery strategies are thoroughly detailed in Section 3. Section 4 provides a comprehensive overview of the experimental setup, with additional details included in the Supplementary Material. Upon publication, we will release the source code for pruning, recovery training, and the compressed model checkpoints, allowing researchers to replicate and extend our work. For evaluation, we used a widely adopted evaluation suite, which is publicly available and cited in our paper, ensuring a consistent and transparent benchmarking process.

References

549 550

551

555

556

558

563

564

565

- Saleh Ashkboos, Maximilian L Croci, Marcelo Gennari do Nascimento, Torsten Hoefler, and James
 Hensman. Slicegpt: Compress large language models by deleting rows and columns. *arXiv preprint arXiv:2401.15024*, 2024.
 - Haoli Bai, Wei Zhang, Lu Hou, Lifeng Shang, Jing Jin, Xin Jiang, Qun Liu, Michael Lyu, and Irwin King. Binarybert: Pushing the limit of bert quantization, 2021. URL https://arxiv.org/abs/2012.15701.
- Li* Bo, Zhang* Peiyuan, Zhang* Kaichen, Pu* Fanyi, Du Xinrun, Dong Yuhao, Liu Haotian, Zhang Yuanhan, Zhang Ge, Li Chunyuan, and Liu Ziwei. Lmms-eval: Accelerating
 the development of large multimoal models, March 2024. URL https://github.com/
 EvolvingLMMs-Lab/lmms-eval.
 - Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *arXiv preprint arXiv:2404.16821*, 2024a.
- Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong
 Zhang, Xizhou Zhu, Lewei Lu, et al. Internvl: Scaling up vision foundation models and aligning
 for generic visual-linguistic tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 24185–24198, 2024b.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality, March 2023. URL https: //lmsys.org/blog/2023-03-30-vicuna/.
- Xiangxiang Chu, Limeng Qiao, Xinyang Lin, Shuang Xu, Yang Yang, Yiming Hu, Fei Wei, Xinyu
 Zhang, Bo Zhang, Xiaolin Wei, et al. Mobilevlm: A fast, reproducible and strong vision language
 assistant for mobile devices. *arXiv preprint arXiv:2312.16886*, 2023.
- Lucio Dery, Steven Kolawole, Jean-François Kagy, Virginia Smith, Graham Neubig, and Ameet Talwalkar. Everybody prune now: Structured pruning of llms with only forward passes, 2024. URL https://arxiv.org/abs/2402.05406.
- 583 Xiaohan Ding, Guiguang Ding, Yuchen Guo, and Jungong Han. Centripetal sgd for pruning very
 584 deep convolutional networks with complicated structure, 2019. URL https://arxiv.org/
 585 abs/1904.03837.
- Xin Dong, Shangyu Chen, and Sinno Jialin Pan. Learning to prune deep neural networks via layerwise optimal brain surgeon, 2017. URL https://arxiv.org/abs/1705.07565.
- Angela Fan, Edouard Grave, and Armand Joulin. Reducing transformer depth on demand with structured dropout. *arXiv preprint arXiv:1909.11556*, 2019.
- Gongfan Fang, Xinyin Ma, Mingli Song, Michael Bi Mi, and Xinchao Wang. Depgraph: Towards any structural pruning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 16091–16101, 2023.

594 595 596	Matteo Farina, Massimiliano Mancini, Elia Cunegatti, Gaowen Liu, Giovanni Iacca, and Elisa Ricci. Multiflow: Shifting towards task-agnostic vision-language pruning. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 16185–16195, 2024.
597 598	Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. 2019. URL https://arxiv.org/abs/1803.03635.
599	····· ·, · · · · · · · · · · · · · · ·
601	Jianping Gou, Baosheng Yu, Stephen J. Maybank, and Dacheng Tao. Knowledge distillation:
602	A survey. International Journal of Computer Vision, 129(6):1789–1819, March 2021. ISSN 1572 1405 deite 10.1007/s11262.021.01452 = URL https://doi.org/10.1007/s12007
603	13/3-1403. doi: 10.100//S11203-021-01453-Z. UKL http://dx.doi.org/10.100//
604	511205-021-01455-2.
605	Yuxian Gu, Li Dong, Furu Wei, and Minlie Huang. Knowledge distillation of large language models. <i>arXiv preprint arXiv:2306.08543</i> , 2023.
607	Yuxian Gu, Li Dong, Furu Wei, and Minlie Huang. Minillm: Knowledge distillation of large lan-
608 609	guage models, 2024. URL https://arxiv.org/abs/2306.08543.
610 611	Muyang He, Yexin Liu, Boya Wu, Jianhao Yuan, Yueze Wang, Tiejun Huang, and Bo Zhao. Efficient multimodal learning from data-centric perspective. <i>arXiv preprint arXiv:2402.11530</i> , 2024.
612	Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network, 2015.
613	URL https://arxiv.org/abs/1503.02531.
614	Ari Holtzman, Jan Buye, Li Du, Maxwell Forbes, and Veiin Choi. The curious case of neural text
615	degeneration. arXiv preprint arXiv:1904.09751, 2019.
617	Chang Ve Heich Chan Lings Li Chik Kasa Veh Hester Nelkest Venkin Fuili Alexander
618	Ratner Raniay Krishna Chen-Yu Lee and Tomas Pfister Distilling step-by-step! outper-
619	forming larger language models with less training data and smaller model sizes. 2023. URL
620	https://arxiv.org/abs/2305.02301.
622	Yen-Chang Hsu, Ting Hua, Sungen Chang, Oian Lou, Yilin Shen, and Hongxia Jin, Language model
623	compression with weighted low-rank factorization, 2022. URL https://arxiv.org/abs/
624	
625 626 627	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. <i>arXiv preprint arXiv:2106.09685</i> , 2021a.
628	
629 630	Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models, 2021b. URL https:
631	//arxiv.org/abs/2106.09685.
632	Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning
633	and compositional question answering. In Proceedings of the IEEE/CVF conference on computer
634	vision and pattern recognition, pp. 6700–6709, 2019.
635	Majan Javaharini Sébastian Ruback Marsh Abdin Justi Anaja Sabastian Ruback Caja
636	César Teodoro Mendes Weizhu Chen Allie Del Giorno Ronen Eldan Siyakanth Goni et al
637	Phi-2: The surprising power of small language models. <i>Microsoft Research Blog</i> , 2023.
638	
639	Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and Qun Liu.
640 641	org/abs/1909.10351.
642	Vizhang Jin Jian Li Vavin Liu Tianiun Cu. Kai Wu Zhangkai Bang Muuang Ha Da Zhan Viz
643	Tan Zhenye Gan et al. Efficient multimodal large language models: A survey arXiv preprint
644	arXiv:2405.10739, 2024.
645	
646 647	Siddharth Karamcheti, Suraj Nair, Ashwin Balakrishna, Percy Liang, Thomas Kollar, and Dorsa Sadigh. Prismatic vlms: Investigating the design space of visually-conditioned language models. <i>arXiv preprint arXiv:2402.07865</i> , 2024.

651

657

658

659

- ⁶⁴⁸ Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. Albert: A lite bert for self-supervised learning of language representations, 2020. URL https://arxiv.org/abs/1909.11942.
- Namhoon Lee, Thalaiyasingam Ajanthan, Stephen Gould, and Philip H. S. Torr. A signal propagation perspective for pruning neural networks at initialization, 2020. URL https://arxiv.org/abs/1906.06307.
- Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning filters for
 efficient convnets, 2017. URL https://arxiv.org/abs/1608.08710.
 - Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. Evaluating object hallucination in large vision-language models. *arXiv preprint arXiv:2305.10355*, 2023.
- Kevin J Liang, Weituo Hao, Dinghan Shen, Yufan Zhou, Weizhu Chen, Changyou Chen, and
 Lawrence Carin. Mixkd: Towards efficient distillation of large-scale language models, 2021.
 URL https://arxiv.org/abs/2011.00593.
- Bin Lin, Zhenyu Tang, Yang Ye, Jiaxi Cui, Bin Zhu, Peng Jin, Junwu Zhang, Munan Ning, and
 Li Yuan. Moe-llava: Mixture of experts for large vision-language models. *arXiv preprint arXiv:2401.15947*, 2024.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning, 2023. URL https://arxiv.org/abs/2304.08485.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction
 tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recogni- tion*, pp. 26296–26306, 2024.
- Liyang Liu, Shilong Zhang, Zhanghui Kuang, Aojun Zhou, Jing-Hao Xue, Xinjiang Wang, Yimin Chen, Wenming Yang, Qingmin Liao, and Wayne Zhang. Group fisher pruning for practical network compression, 2021. URL https://arxiv.org/abs/2108.00708.
- Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord,
 Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal reasoning via thought chains for
 science question answering. *Advances in Neural Information Processing Systems*, 35:2507–2521, 2022.
- Xinyin Ma, Gongfan Fang, and Xinchao Wang. Llm-pruner: On the structural pruning of large
 language models. *Advances in neural information processing systems*, 36:21702–21720, 2023.
- Andrey Malinin and Mark Gales. Reverse kl-divergence training of prior networks: Improved uncertainty and adversarial robustness. *Advances in neural information processing systems*, 32, 2019.
- JS McCarley, Rishav Chakravarti, and Avirup Sil. Structured pruning of a bert-based question answering model. *arXiv preprint arXiv:1910.06360*, 2019.
- Kin Men, Mingyu Xu, Qingyu Zhang, Bingning Wang, Hongyu Lin, Yaojie Lu, Xianpei Han, and
 Weipeng Chen. Shortgpt: Layers in large language models are more redundant than you expect.
 arXiv preprint arXiv:2403.03853, 2024.
- Paul Michel, Omer Levy, and Graham Neubig. Are sixteen heads really better than one? *Advances in neural information processing systems*, 32, 2019.
- Sejun Park, Jaeho Lee, Sangwoo Mo, and Jinwoo Shin. Lookahead: A far-sighted alternative of
 magnitude-based pruning, 2020. URL https://arxiv.org/abs/2002.04809.
- Niclas Popp, Jan Hendrik Metzen, and Matthias Hein. Zero-shot distillation for image encoders: How to make effective use of synthetic data. *arXiv preprint arXiv:2404.16637*, 2024.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
 models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.

702 Hassan Sajjad, Fahim Dalvi, Nadir Durrani, and Preslav Nakov. On the effect of dropping layers of 703 pre-trained transformer models. Computer Speech & Language, 77:101429, 2023. 704 Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version 705 of bert: smaller, faster, cheaper and lighter, 2020a. URL https://arxiv.org/abs/1910. 706 01108. 707 708 Victor Sanh, Thomas Wolf, and Alexander M. Rush. Movement pruning: Adaptive sparsity by 709 fine-tuning, 2020b. URL https://arxiv.org/abs/2005.07683. 710 Dachuan Shi, Chaofan Tao, Ying Jin, Zhendong Yang, Chun Yuan, and Jiaqi Wang. Upop: Unified 711 and progressive pruning for compressing vision-language transformers. In International Confer-712 ence on Machine Learning, pp. 31292–31311. PMLR, 2023. 713 714 Kaitao Song, Hao Sun, Xu Tan, Tao Qin, Jianfeng Lu, Hongzhi Liu, and Tie-Yan Liu. Light-715 paff: A two-stage distillation framework for pre-training and fine-tuning, 2020. URL https: 716 //arxiv.org/abs/2004.12817. 717 Siqi Sun, Yu Cheng, Zhe Gan, and Jingjing Liu. Patient knowledge distillation for bert model 718 compression, 2019. URL https://arxiv.org/abs/1908.09355. 719 720 Elena Voita, David Talbot, Fedor Moiseev, Rico Sennrich, and Ivan Titov. Analyzing multi-head 721 self-attention: Specialized heads do the heavy lifting, the rest can be pruned. arXiv preprint 722 arXiv:1905.09418, 2019. 723 Tiannan Wang, Wangchunshu Zhou, Yan Zeng, and Xinsong Zhang. Efficientvlm: Fast and accurate 724 vision-language models via knowledge distillation and modal-adaptive pruning. In Findings of 725 the Association for Computational Linguistics: ACL 2023, pp. 13899–13913, 2023. 726 727 Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. Minilm: Deep 728 self-attention distillation for task-agnostic compression of pre-trained transformers, 2020. URL 729 https://arxiv.org/abs/2002.10957. 730 Wenhui Wang, Hangbo Bao, Shaohan Huang, Li Dong, and Furu Wei. Minilmv2: Multi-head 731 self-attention relation distillation for compressing pretrained transformers, 2021. URL https: 732 //arxiv.org/abs/2012.15828. 733 734 Haoran Wei, Lingyu Kong, Jinyue Chen, Liang Zhao, Zheng Ge, En Yu, Jianjian Sun, Chunrui 735 Han, and Xiangyu Zhang. Small language model meets with reinforced vision vocabulary. arXiv 736 preprint arXiv:2401.12503, 2024. 737 Mengzhou Xia, Tianyu Gao, Zhiyuan Zeng, and Danqi Chen. Sheared llama: Accelerating language 738 model pre-training via structured pruning, 2024. URL https://arxiv.org/abs/2310. 739 06694. 740 Xiaohan Xu, Ming Li, Chongyang Tao, Tao Shen, Reynold Cheng, Jinyang Li, Can Xu, Dacheng 741 Tao, and Tianyi Zhou. A survey on knowledge distillation of large language models, 2024. URL 742 https://arxiv.org/abs/2402.13116. 743 744 Chuanguang Yang, Zhulin An, Libo Huang, Junyu Bi, Xinqiang Yu, Han Yang, Boyu Diao, and 745 Yongjun Xu. Clip-kd: An empirical study of clip model distillation. In Proceedings of the 746 IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 15952–15962, 2024. 747 Zhewei Yao, Reza Yazdani Aminabadi, Minjia Zhang, Xiaoxia Wu, Conglong Li, and Yuxiong He. 748 Zeroquant: Efficient and affordable post-training quantization for large-scale transformers, 2022. 749 URL https://arxiv.org/abs/2206.01861. 750 751 Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. A survey on 752 multimodal large language models. arXiv preprint arXiv:2306.13549, 2023. 753 Zhonghui You, Kun Yan, Jinmian Ye, Meng Ma, and Ping Wang. Gate decorator: Global filter 754 pruning method for accelerating deep convolutional neural networks, 2019. URL https:// 755 arxiv.org/abs/1909.08174.

756 757 758 759 760	Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhu Chen. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In <i>Proceedings of CVPR</i> , 2024.
761 762 763 764 765	Ofir Zafrir, Guy Boudoukh, Peter Izsak, and Moshe Wasserblat. Q8bert: Quantized 8bit bert. In 2019 Fifth Workshop on Energy Efficient Machine Learning and Cognitive Computing - NeurIPS Edition (EMC2-NIPS). IEEE, December 2019. doi: 10.1109/emc2-nips53020.2019.00016. URL http://dx.doi.org/10.1109/EMC2-NIPS53020.2019.00016.
766 767 768	Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pp. 11975–11986, 2023.
769 770 771	Rongzhi Zhang, Jiaming Shen, Tianqi Liu, Jialu Liu, Michael Bendersky, Marc Najork, and Chao Zhang. Do not blindly imitate the teacher: Using perturbed loss for knowledge distillation, 2023. URL https://arxiv.org/abs/2305.05010.
772 773 774 775	Minjie Zhu, Yichen Zhu, Xin Liu, Ning Liu, Zhiyuan Xu, Chaomin Shen, Yaxin Peng, Zhicai Ou, Feifei Feng, and Jian Tang. A comprehensive overhaul of multimodal assistant with small language models. <i>arXiv preprint arXiv:2403.06199</i> , 2024a.
776 777 778 779	Yichen Zhu, Minjie Zhu, Ning Liu, Zhiyuan Xu, and Yaxin Peng. Llava-phi: Efficient multi-modal assistant with small language model. In <i>Proceedings of the 1st International Workshop on Efficient Multimedia Computing under Limited</i> , pp. 18–22, 2024b.
780 781 782	
783 784 785	
786 787	
788 789 790	
791 792	
793 794 795	
796 797 708	
798 799 800	
801 802	
803 804 805	
806 807	
808 809	

A MODEL ARCHITECTURE OF LLAVA AND BUNNY

Table 4 outlines the architectures of the Bunny and LLaVA models. LLaVA-v1.5-7B employs CLIP-ViT-L (Radford et al., 2021) as the vision encoder and Vicuna-v1.5 Chiang et al., 2023 as the language decoder, while Bunny-v1.0-3B utilizes SigLIP-SO (Zhai et al., 2023) as the vision encoder and Phi-2 (Javaheripi et al., 2023) as the language decoder. Both models leverage MLP layers to align the vision and language modalities.

Model	Parameters	Vision Encoder	Multimodal Projector	Language Decoder
LLaVA-v1.5-7B	7.0B	CLIP-ViT-L (0.3B)	mlp2x-gelu (0.01B)	Vicuna-v1.5 (6.7B)
Bunny-v1.0-3B	3.2B	SigLIP-SO (0.4B)	mlp2x-gelu (0.02B)	Phi-2 (2.8B)

Table 4: Architecture details of the uncompressed models. We present the number of parameters, along with the vision encoder, multimodal projector and the language decoder of the models included in our study.

B IMPLEMENTATION DETAILS OF HIDDEN STATE MATCHING

To determine which LLM layers' hidden states to map between the pruned and unpruned models, we explore three options: matching the last layer, the last two, and the last three layers. Table 5 shows that matching only the last layer's hidden state yields the best performance.

Ratio	Layer-1	Layer-1,2	Layer-1,2,3
12.8%	95.34%	95.17%	96.25%
25.5%	91.02%	90.48%	90.97%
39.0%	87.08%	86.12%	84.84%
51.8%	75.25%	72.56%	72.87%

Table 5: Results for recovering widthwise pruned Bunny with hidden state mapping. We compare the relative performance for mapping the last layer (layer-1), the last two layers (layer-1,2), and the last three layers (layer-1,2,3). By only mapping the last LLM layer the best performance is achieved.

C MORE RESULTS FOR PRUNING

Table 6 presents the model performance for widthwise and layerwise pruning. For both of the Bunny and LLaVA models, widthwise pruning consistently outperforms layerwise pruning. This performance gap widens as the compression ratio increases, with widthwise pruning showing a more significant advantage at higher compression ratios.

D MORE RESULTS ON MODEL EFFICIENCY

For the models pruned by layerwise method, we also assess their memory consumption as FLOPs. Memory consumption refers to the allocated GPU memory, while FLOPS are measured using the Calflops codebase². The results in Table 7 show the same trend as widthwise pruning, indicating that the achieved compressions directly translate into improvements in memory efficiency and computational cost for both widthwise and layerwise pruning.

E KNOWLEDGE DISTILLATION WITH AND WITHOUT FINETUNING

Figure 5 and Figure 6 compare logits-based knowledge distillation (represented by RKL) and hidden
state matching-based knowledge distillation (represented by L2 loss), with and without supervised
fine-tuning, following widthwise and layerwise pruning, respectively. While knowledge distillation

²Caflops codebase: https://github.com/MrYxJ/calculate-flops.pytorch

$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	Method	Size	PruneRatio	MMMU	GQA	SQA	MME-C	MME-P	POPE	AVG	AVG-%
Width-wise 6.3B 15% 32.40 59.34 63.21 268.93 1432.47 86.57 57.79 92.799 5.5B 30% 31.00 52.59 54.29 253.21 1174.93 86.29 52.43 84.179 4.8B 45% 27.60 20.86 12.10 70.00 347.45 45.96 22.11 35.499 4.0B 60% 23.30 0.43 0.40 2.14 19.24 3.94 4.88 7.849 Depth-wise 6.3B 15% 31.80 42.77 55.23 202.14 701.83 86.38 46.09 74.009 5.5B 30% 32.70 42.18 59.64 210.71 921.88 78.69 47.61 76.439 4.8B 45% 26.90 14.39 3.82 132.86 616.63 51.69 24.04 38.609 4.0B 60% 25.80 0.00 0.00 0.00 0.00 4.08 6.909 Bunny-v1_0-3B<	LLaVA-v1.5-7B	7.0B		35.10	61.98	68.67	363.21	1511.33	86.99	62.28	100.00%
5.5B 30% 31.00 52.59 54.29 253.21 1174.93 86.29 52.43 84.179 4.8B 45% 27.60 20.86 12.10 70.00 347.45 45.96 22.11 35.499 4.0B 60% 23.30 0.43 0.40 2.14 19.24 3.94 4.88 7.849 Depth-wise 6.3B 15% 31.80 42.77 55.23 202.14 701.83 86.38 46.09 74.009 5.5B 30% 32.70 42.18 59.64 210.71 921.88 78.69 47.61 76.439 4.8B 45% 26.90 14.39 3.82 132.86 616.63 51.69 24.04 38.609 4.0B 60% 25.80 0.00 0.00 0.00 0.00 4.08 6.995 100.09 Bunny-v1_0-3B 3.2B 34.10 54.72 70.70 289.30 1487.71 87.82 59.65 100.009 Width-w	Width-wise	6.3B	15%	32.40	59.34	63.21	268.93	1432.47	86.57	57.79	92.79%
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$		5.5B	30%	31.00	52.59	54.29	253.21	1174.93	86.29	52.43	84.17%
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$		4.8B	45%	27.60	20.86	12.10	70.00	347.45	45.96	22.11	35.49%
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$		4.0B	60%	23.30	0.43	0.40	2.14	19.24	3.94	4.88	7.84%
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	Depth-wise	6.3B	15%	31.80	42.77	55.23	202.14	701.83	86.38	46.09	74.00%
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$		5.5B	30%	32.70	42.18	59.64	210.71	921.88	78.69	47.61	76.43%
4.0B 60% 25.80 0.00 0.00 0.00 0.00 4.30 6.90% Bunny-v1_0-3B 3.2B 34.10 54.72 70.70 289.30 1487.71 87.82 59.65 100.00% Width-wise 2.8B 15% 30.90 51.83 65.64 242.50 1207.85 87.94 54.50 95.489 2.5B 30% 28.40 45.65 55.73 199.64 807.95 87.13 47.04 87.579 2.0B 45% 25.70 37.92 3.42 200.00 618.25 83.12 34.35 60.669 1.6B 60% 24.80 6.12 0.00 141.07 293.23 2.34 10.93 13.529 Depth-wise 2.8B 15% 33.80 29.42 69.66 271.43 1456.41 87.91 54.59 91.529 2.0B 45% 23.90 16.85 3.47 191.43 867.37 80.09 31.94 53.549 <		4.8B	45%	26.90	14.39	3.82	132.86	616.63	51.69	24.04	38.60%
Bunny-v1_0-3B 3.2B 34.10 54.72 70.70 289.30 1487.71 87.82 59.65 100.009 Width-wise 2.8B 15% 30.90 51.83 65.64 242.50 1207.85 87.94 54.50 95.489 2.5B 30% 28.40 45.65 55.73 199.64 807.95 87.13 47.04 87.579 2.0B 45% 25.70 37.92 3.42 200.00 618.25 83.12 34.35 60.669 1.6B 60% 24.80 6.12 0.00 141.07 293.23 2.34 10.93 13.529 Depth-wise 2.8B 15% 33.80 29.42 69.66 271.43 1456.41 87.91 54.59 91.529 2.5B 30% 29.00 24.77 28.76 272.86 1273.34 86.50 44.47 74.559 2.0B 45% 23.90 16.85 3.47 191.43 867.37 80.09 31.94 53.549 </td <td></td> <td>4.0B</td> <td>60%</td> <td>25.80</td> <td>0.00</td> <td>0.00</td> <td>0.00</td> <td>0.00</td> <td>0.00</td> <td>4.30</td> <td>6.90%</td>		4.0B	60%	25.80	0.00	0.00	0.00	0.00	0.00	4.30	6.90%
Width-wise 2.8B 15% 30.90 51.83 65.64 242.50 1207.85 87.94 54.50 95.489 2.5B 30% 28.40 45.65 55.73 199.64 807.95 87.13 47.04 87.579 2.0B 45% 25.70 37.92 3.42 200.00 618.25 83.12 34.35 60.669 1.6B 60% 24.80 6.12 0.00 141.07 293.23 2.34 10.93 13.529 Depth-wise 2.8B 15% 33.80 29.42 69.66 271.43 1456.41 87.91 54.59 91.529 2.5B 30% 29.00 24.77 28.76 272.86 1273.34 86.50 44.47 74.559 2.0B 45% 23.90 16.85 3.47 191.43 867.37 80.09 31.94 53.549 1.6B 60% 26.60 0.02 17.15 0.71 55.92 0.02 7.78 13.049 <td>Bunny-v1_0-3B</td> <td>3.2B</td> <td></td> <td>34.10</td> <td>54.72</td> <td>70.70</td> <td>289.30</td> <td>1487.71</td> <td>87.82</td> <td>59.65</td> <td>100.00%</td>	Bunny-v1_0-3B	3.2B		34.10	54.72	70.70	289.30	1487.71	87.82	59.65	100.00%
2.5B 30% 28.40 45.65 55.73 199.64 807.95 87.13 47.04 87.579 2.0B 45% 25.70 37.92 3.42 200.00 618.25 83.12 34.35 60.669 1.6B 60% 24.80 6.12 0.00 141.07 293.23 2.34 10.93 13.529 Depth-wise 2.8B 15% 33.80 29.42 69.66 271.43 1456.41 87.91 54.59 91.529 2.5B 30% 29.00 24.77 28.76 272.86 1273.34 86.50 44.47 74.559 2.0B 45% 23.90 16.85 3.47 191.43 867.37 80.09 31.94 53.549 1.6B 60% 26.60 0.02 17.15 0.71 55.92 0.02 7.78 13.049	Width-wise	2.8B	15%	30.90	51.83	65.64	242.50	1207.85	87.94	54.50	95.48%
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$		2.5B	30%	28.40	45.65	55.73	199.64	807.95	87.13	47.04	87.57%
1.6B 60% 24.80 6.12 0.00 141.07 293.23 2.34 10.93 13.529 Depth-wise 2.8B 15% 33.80 29.42 69.66 271.43 1456.41 87.91 54.59 91.529 2.5B 30% 29.00 24.77 28.76 272.86 1273.34 86.50 44.47 74.559 2.0B 45% 23.90 16.85 3.47 191.43 867.37 80.09 31.94 53.549 1.6B 60% 26.60 0.02 17.15 0.71 55.92 0.02 7.78 13.049		2.0B	45%	25.70	37.92	3.42	200.00	618.25	83.12	34.35	60.66%
Depth-wise 2.8B 15% 33.80 29.42 69.66 271.43 1456.41 87.91 54.59 91.529 2.5B 30% 29.00 24.77 28.76 272.86 1273.34 86.50 44.47 74.559 2.0B 45% 23.90 16.85 3.47 191.43 867.37 80.09 31.94 53.549 1.6B 60% 26.60 0.02 17.15 0.71 55.92 0.02 7.78 13.049		1.6B	60%	24.80	6.12	0.00	141.07	293.23	2.34	10.93	13.52%
2.5B 30% 29.00 24.77 28.76 272.86 1273.34 86.50 44.47 74.55% 2.0B 45% 23.90 16.85 3.47 191.43 867.37 80.09 31.94 53.54% 1.6B 60% 26.60 0.02 17.15 0.71 55.92 0.02 7.78 13.04%	Depth-wise	2.8B	15%	33.80	29.42	69.66	271.43	1456.41	87.91	54.59	91.52%
2.0B 45% 23.90 16.85 3.47 191.43 867.37 80.09 31.94 53.54% 1.6B 60% 26.60 0.02 17.15 0.71 55.92 0.02 7.78 13.04%	*	2.5B	30%	29.00	24.77	28.76	272.86	1273.34	86.50	44.47	74.55%
$1.6B \qquad 60\% \qquad 26.60 \qquad 0.02 \qquad 17.15 \qquad 0.71 \qquad 55.92 \qquad 0.02 \qquad 7.78 \qquad 13.04\%$		2.0B	45%	23.90	16.85	3.47	191.43	867.37	80.09	31.94	53.54%
		1.6B	60%	26.60	0.02	17.15	0.71	55.92	0.02	7.78	13.04%

Table 6: Pruning results for LLaVA-v1.5-7B and Bunny-v1-3B. Size is the number of total parameters of the model, while the compression ratio (Ratio) indicates the proportion of remaining language model parameters compared to the pre-pruning state. For both models, width-wise pruning results in better performance without finetuning compared to depth-wise pruning.

Datia	Bun	ny	LLaVA			
Ratio	Mem. (MiB)	FLOPS (T)	Mem. (MiB)	FLOPS (T)		
0%	6,167	4.77	13,546	9.57		
15%	5,411	4.16	11,604	8.03		
30%	4,659	3.56	9,664	6.92		
45%	3,907	2.95	7,724	5.55		
60%	3,006	2.22	5,496	3.9		

Table 7: Memory requirements (Mem.) and FLOPS for the Bunny and LLaVA models at various compression ratios. The models are pruned layerwise. The evaluation is performed in inference mode, where each model is provided with an image and a prompt containing 50 tokens.

alone helps in recovering performance post-pruning, it remains less effective than supervised finetuning. However, when combined with supervised fine-tuning, it results in superior performance.



Figure 5: Comparison of L2 and RKL distillation strategies with and without additional fine-tuning loss for Bunny and Llava models compressed by widthwise pruning. The plot shows performance differences relative to standard fine-tuning across varying compression ratios.



Figure 6: Comparison of L2 and RKL distillation strategies with and without additional fine-tuning loss for Bunny and Llava models compressed by layerwise pruning. The plot shows performance differences relative to standard fine-tuning across varying compression ratios.

F KNOWLEDGE DISTILLATION STRATEGIES COMPARISON AFTER LAYERWISE PRUNING

Figure 7 compares various distillation strategies based on their relative improvement over finetuning alone after layerwise pruning. Similar to the results after widthwise pruning, applying hidden states matching yields the best performance, or at least matches other methods. The trend that RKL generally outperforms KL is also observed here.



Figure 7: Comparison of different distillation recovery strategies (KL loss, RKL loss, L2 loss, and their combinations) for Bunny and LlaVA models pruned with layerwise pruning. The plot shows the relative performance improvement of each strategy over standard fine-tuning across various compression ratios. The results demonstrate that distillation helps recover more performance than fine-tuning alone, with the L2 loss component consistently leading to the largest performance gains, particularly at higher compression ratios.

G OU

QUALITATIVE RESULTS OF THE COMPRESSED MODELS

Table 8 shows the qualitative results of the compressed LLaVA models. After compression, the models demonstrate a strong ability to comprehend visual inputs and generate detailed textual outputs.

H PRUNING VS. QUANTIZATION

967 The table includes results for LLaVA-6B and LLaVA-5B, where the former is layerwise pruned and
968 recovery trained, and the latter is widthwise pruned and recovery trained. Our findings indicate
969 that combining structured pruning with quantization provides a balanced approach: it significantly
970 reduces memory usage while mitigating the latency overhead associated with quantization alone.
971 This synergy leverages the strengths of both methods, offering an effective strategy for compressing large models while maintaining performance and efficiency.

Model	Response
	And a state of the second s
	Materia
	and the second
	Prompt : A chat between a curious user and an artificial intelligence
	assistant. The assistant gives helpful, detailed, and polite answers to the
	user's questions. USER: $< image >$ What is funny in the image?
LLaVA-6.3B	The image is funny because it features a dog wearing a Renaissance-
	style dress and a bandana, posing as if it's posing for a painting. This
	is humorous because it is a playful and unconventional representation
	of a dog, which is not something you would typically see in a painting
	or a Renaissance-style dress. The dog's pose and attire create a sense
	viewers.
LLaVA-5 5B	The image is funny because it features a dog dressed up in a human-like
ELU VII 5.5D	attire, wearing a long black dress and a hat. This unusual and amusing
	scene is not something you would typically see in real life, as dogs don't
	wear clothes or accessories like humans. The dog's expression and the
	way it is posing for the picture further contribute to the comic effect of
	the image. The image captures the viewer's attention and sparks a sense
	dressed like a human
LL WA 2 OD	The image is furner because it factures a day drawed in a larger 1'1.
LLavA-3.8B	the image is furning because it features a dog dressed in a human-like dress posing in a pose as if it were a person. The dog appears to be
	dressed in a dress, which is a humorous and unexpected scene as doos
	are not typically dressed like humans. The dog's pose and the overall
	scene create a sense of humor and playfulness, making the image a
	delightful and entertaining piece.

Table 8: Qualitative analysis of compressed LLaVA models with respect to the image. We feed the
 compressed LLaVA models the image paired with a prompt to generate a response. The compressed
 LLaVA models are able to understand the visual inputs and output reasonable texts.

1015 1016

1018

1017 I GENERALIZABILITY OF THE BEST PRACTICES

We also explored the impact of recovery training, specifically finetuning the multimodal projector and the language model, as well as the importance of incorporating knowledge distillation. These results, shown in Table 11, confirm that finetuning only the multimodal projector is sufficient at small compression ratios, as pruning minimally impacts the language model but disrupts multimodal alignment. For InternVL, finetuning the projector alone recovers 96.9% of the performance at a 15% compression ratio, compared to 97.8% when both the projector and the language model are finetuned. At a 30% compression ratio, projector-only finetuning recovers 75.1% while finetuning both components recovers 86.6%. Furthermore, combining supervised finetuning with intermediate

Model	Memory(MiB)	Ratio	MMMU	GQA	SQA	MME-C	MME-P	POPE	AVG	Latency
LLaVA-7B	13546		35.10	61.98	68.67	363.21	1511.33	86.99	62.28	105ms ± 1.5ms
LLaVA-7B.int8()	7518		35.2	61.87	68.22	350.71	1508.41	86.54	61.85	398ms ± 1.31ms
LLaVA-6B	11604	15%	35.40	61.17	68.07	328.57	1454.20	86.51	60.82	95 ms ± 8.1 ms
LLaVA-6B.int8()	6473	15%	35.40	61.17	68.07	328.57	1454.20	86.51	60.82	125 ms ± 937 µs
LLaVA-5B	9548	30%	31.80	60.71	60.54	252.50	1407.08	86.68	56.94	$80.7 \text{ ms} \pm 634 \mu$
LLaVA-5B.int8()	5389	30%	31.6	60.65	60.09	263.57	1410.28	86.78	57.10	141 ms ± 2.4 ms

Table 9: Pruning and Quantization comparison on LLaVA-v1.5-7B.

034	Model	Size	Ratio	MMMU	GQA	SQA	MME-C	MME-P	POPE	AVG	AVG-%
035	InternVL-Chat-4B-V1-5	4B		43.20	62.57	93.30	547.50	1,596.71	88.00	72.56	100%
36	Layerwise	3.5B	15%	42.70	54.43	92.96	527.86	1,534.37	88.09	70.15	96.68%
)37	Widthwise	3.5B	15%	43.60	56.35	93.12	510.10	1,588.30	87.96	70.70	97.44%

Table 10: Comparison of different pruning techniques on the model InternVL-Chat-4B.

representation distillation consistently achieves the highest performance across all compression ratios. For InternVL, this approach recovers 98.2% of performance at 15% compression and 87.2% at 30%, underscoring its effectiveness.

Overall, the results from InternVL demonstrate that our best practices generalize effectively to different model architectures, confirming their applicability and robustness.

Model	Size	Ratio	MMMU	GQA	SQA	MME-C	MME-P	POPE	AVG	AVG-
Mini-InternVL-Chat-4B-V1-5	4B	-	43.20	62.57	93.30	547.50	1,596.71	88.00	72.56	100%
SFT mm only	3.5B	15%	42.70	54.43	92.96	527.86	1,534.37	88.09	70.15	96.68
SFT mm + LLM	3.5B	15%	43.10	56.34	93.36	524.64	1,585.83	88.10	70.96	97.80
SFT+KD mm + LLM	3.5B	15%	43.30	56.17	93.41	539.64	1,582.58	87.89	71.23	98.16
SFT mm only	3B	30%	33.30	27.39	62.82	197.14	845.60	73.20	43.94	60.56
SFT mm + LLM	3B	30%	34.40	53.46	76.80	432.27	1,438.28	86.57	62.86	86.64
SFT+KD mm + LLM	3B	30%	36.20	53.77	76.60	448.21	1,410.88	86.60	63.29	87.23

Table 11: Comparison of recovery training only multimodal projector (mm) and large language model on Mini-InternVL-Chat-4B-V1-5 after layerwise pruning with different recovery strategies, i.e., supervised finetuning (SFT) and knowledge distillation (KD) on intermediate representations.