000 001

# **Enhancing Pruned Models by Input Compensation**

#### **Anonymous Authors**

#### Abstract

Although foundation models are powerful, they are large and require substantial memory and computation resources to serve. To address this issue, many pruning methods have been proposed to reduce model size, thereby achieving memory and computational efficiency. These methods ad*just the retained weights* to compensate for the removed weights. In this paper, we propose a novel approach called input compensation (IC) to improve the performance of pruned models, i.e., adjust the input to compensate for the removed weights. Unlike existing pruning methods, which are designed in the parameter space, the proposed IC is designed in the input space. Hence, IC is complementary to existing methods and can be integrated with them. Extensive experiments on various tasks, including image classification, language modeling, and image generation, demonstrate that IC effectively boosts the performance of pruned models.

#### 1. Introduction

Foundation models (Radford et al., 2021; Touvron et al., 2023b; Podell et al., 2024) have achieved great success in a variety of domains, such as computer vision and natural language processing. As the availability of data and computational resources expands, these models have scaled in both size and performance (Touvron et al., 2023a;b; Meta, 2024). However, the substantial number of parameters in these models require extensive computational resources to serve, making it challenging to deploy them on resourceconstrained devices such as smartphones and laptops. To reduce the costs, numerous model compression techniques have been proposed to reduce the model size, e.g., distillation (Polino et al., 2018; Wang et al., 2019; Liang et al., 2023), quantization (Lin et al., 2024; Dettmers et al., 2022; Shao et al., 2024; Xiao et al., 2023), and pruning (Han et al., 2015; Frantar & Alistarh, 2023; Zhang et al., 2024; Sun et al., 2024). Quantization requires specialized hardware support, while distillation requires extensive retraining. Hence, we focus on pruning, which is a simple and representative technique.

*Pruning* reduces the model size by removing individual weights or rows/columns according to their importance scores. A pruned model can achieve promising performance with fewer parameters, resulting in a noticeable reduction in memory and computational demands. A simple but effective pruning method is Magnitude Pruning (Han et al., 2015) which removes weights according to their magnitudes. The underlying assumption is that weights with smaller values contribute less to overall performance. However, this assumption does not always hold and many advanced methods (Sun et al., 2024; Frantar & Alistarh, 2023; Zhang et al., 2024) have been proposed recently.

Current state-of-the-art pruning methods (Frantar & Alistarh, 2023; Das et al., 2023; Zhang et al., 2024; Sun et al., 2024; Dong et al., 2024; An et al., 2024) focus on the parameter space to improve pruning efficacy and can be roughly categorized into two groups: (i) designing an effective score to measure the importance of weight, and (ii) adjusting the remaining unpruned weights to reduce the error caused by the pruned weights. For example, Wanda (Sun et al., 2024) designs an importance score to incorporate input activations with weight magnitude to take outlier features into consideration instead of only weight magnitudes in Magnitude Pruning; SparseGPT (Frantar & Alistarh, 2023) proposes to adjust the unpruned weights by minimizing reconstruction loss using the Optimal Brain Surgeon framework (Hassibi et al., 1993; Singh & Alistarh, 2020; Frantar et al., 2021). The pruned model can be formulated as  $\mathcal{F}(\mathbf{X}; \mathbf{W} \odot \mathbf{M} + \boldsymbol{\Delta}_{\mathbf{w}})$ , where  $\mathcal{F}$  is the model, **X** is the input, **W** is the weight matrix, M is the weight mask determined by the importance score,  $\odot$  is element-wise multiplication, and  $\Delta_w$  (called weight compensation) is an update matrix for the retained weights.

In this paper, we propose a novel method, called input compensation (IC), for enhancing pruned models by adjusting the input to compensate for the removed weights. Specifically, we learn an input compensation  $\Delta_x$  to adjust the input **X**, therefore the output of the pruned model is determined by  $\mathcal{F}(\mathbf{X} + \Delta_x; \hat{\mathbf{W}})$ , where  $\hat{\mathbf{W}}$  is a sparse weight matrix corresponding to the pruned model. We learn a compensation pool consists of multiple candidate compensations and  $\Delta_x$ is a weighted combination of the candidate compensations via the attention mechanism (Vaswani et al., 2017). Different from existing pruning methods, the proposed IC is
designed in the *input space*. Hence, IC complements existing methods that operate in the parameter space and can be
integrated with them to boost their performance. Extensive
experiments on computer vision and natural language processing show that IC brings a large improvement to existing
pruning methods.

Our contributions are summarized as follows: (i) We propose IC which is a novel direction to enhance pruned models; (ii) IC is designed in the input space and, thus, is orthogonal to existing pruning methods designed in the parameter space. Hence, IC can be combined with existing pruning methods; (iii) Experimental results on various tasks demonstrate that IC is beneficial to existing pruning methods.

### 2. Related Work

062

063

064

065

066

067

068

069

070

Model Compression. Foundation Models (Touvron et al., 074 2023a;b; Meta, 2024; Ho et al., 2020; Radford et al., 2021) 075 are large pre-trained models designed to serve as base mod-076 els for various downstream tasks. Though foundation mod-077 els are powerful, their massive parameters usually require 078 extensive computational and memory resources. Many re-079 cent efforts have been devoted to reducing the cost via model compression (Frantar & Alistarh, 2022; Xu et al., 081 2024; Wang et al., 2024). The most popular methods for 082 model compression are pruning, quantization, and distilla-083 tion. Pruning (Han et al., 2015; Zhang et al., 2024; Sun et al., 2024; Dong et al., 2024; Das et al., 2023; An et al., 085 2024; Frantar & Alistarh, 2023) discards parts of the model 086 that are less important or redundant. Quantization (Lin et al., 087 2024; Dettmers et al., 2022; Shao et al., 2024; Xiao et al., 088 2023; Yao et al., 2022; Kim et al., 2024) is a technique to 089 reduce the computational complexity and memory footprint 090 of a neural network by converting the model's parameters 091 (weights and activations) from higher-precision representa-092 tions (such as 32-bit floating-point) to lower-precision ones 093 (such as 8-bit integers). The primary goal of quantization 094 and pruning is to make the model more compressed without 095 significantly sacrificing its performance. Distillation (Polino 096 et al., 2018; Wang et al., 2019; Liang et al., 2023) trains 097 a smaller and more efficient model to replicate the behav-098 ior of a larger and more complex model, thereby retaining 099 much of its performance while significantly reducing com-100 putational resources. Quantization demands specialized hardware (e.g., NVIDIA TensorRT) that supports lower precision arithmetic, while distillation requires an expensive training phase to transfer knowledge from a large teacher 104 model to a small student model. In this paper, we focus on 105 pruning, which is a simple and widely used approach. 106

**Pruning** aims to remove less important weights withoutsignificant performance degradation. Several important met-

rics have been designed recently. The simplest one is based on the parameter magnitude, i.e., Magnitude Pruning (Han et al., 2015). Wanda (Sun et al., 2024) further incorporates weight magnitude with their input activations to consider outlier features when calculating importance scores, while RIA (Zhang et al., 2024) uses relative importance as a pruning metric. Taylor pruning (Molchanov et al., 2022) designs a score based on the weight multiplied by its gradient, while Diff-Pruning (Fang et al., 2023) further uses Taylor expansion over pruned timesteps to identify and discard unimportant parameters. In addition to designing importance scores to find less useful parameters, one can update the unpruned weights to compensate for the error caused by the pruned weights. For example, SparseGPT (Frantar & Alistarh, 2023) and OBC (Frantar & Alistarh, 2022) propose to update the unpruned weights by minimizing a reconstruction loss by the Optimal Brain Surgeon framework (Hassibi et al., 1993; Singh & Alistarh, 2020; Frantar et al., 2021). Different from SparseGPT and OBC, we propose input compensation by *adjusting the inputs* to reduce the error caused by pruning.

Prompting (Radford et al., 2019; Brown et al., 2020; Liu et al., 2022; Ding et al., 2022) is a popular method used in transformer-based models which inserts additional tokens that instruct the model to generate a specific kind of response. These tokens can be either discrete tokens (e.g., "The topic is" for topic classification (Zhang et al., 2022a; Hou et al., 2022; Jiang et al., 2023), "Let's think step by step" for reasoning tasks (Kojima et al., 2022)) or learnable continuous vectors (e.g., prompt tuning (Lester et al., 2021; Liu et al., 2021; Zhang et al., 2022b) or prefix learning (Li & Liang, 2021; Liu et al., 2023)). Unlike prompting that inserts extra tokens into the inputs, our input compensation edits the inputs directly. Furthermore, compensations are input-dependent, while prompts are usually input-independent (Ding et al., 2022; Lester et al., 2021; Liu et al., 2021; Zhang et al., 2022b; Bahng et al., 2022).

In control systems, the idea of **input compensation** (Kuo & Golnaraghi, 1995; Franklin et al., 2002) is practically used to adjust the control signal to reduce the influence of disturbance. The goal is to adjust the input so that the overall system achieves the desired behavior, such as better stability, faster response, or improved accuracy. For example, in feedforward compensation (Campos & Lewis, 1999; Krstic, 2009), if a disturbance is known ahead of time (e.g., wind gusts affecting an airplane), this information can be incorporated into the control signal so that the system compensates for it before it affects the output. In model pruning, the pruned weights can be viewed as disturbances and we use input compensation to enhance pruned models.

**Enhancing Pruned Models by Input Compensation** 



Figure 1: Input compensation for pruned models.

### 3. Preliminary on Pruning Models

124 125 126

127

128

129

130

131

132

133

134

135

136

137

138

139

Let  $\mathbf{W} \in \mathbb{R}^{d_i \times d_o}$  be a weight matrix of a model  $\mathcal{F}$  and  $\mathbf{S}$  be a scoring matrix whose  $\mathbf{S}_{i,j}$  measures the importance of  $\mathbf{W}_{i,j}$ . To prune p% parameters of  $\mathbf{W}$ , a threshold  $\beta$  is determined to satisfy  $\frac{\#\{\mathbf{S}_{i,j} | \mathbf{S}_{i,j} | < \beta\}}{\#\{\mathbf{S}_{i,j}\}} = p\%$ . Using the threshold, we construct a binary weight mask  $\mathbf{M}$  whose  $\mathbf{M}_{i,j} = 1$  if  $|\mathbf{S}_{i,j}| \geq \beta$  else 0 and prune the model as  $\mathbf{W} \odot \mathbf{M}$ . To improve the performance of the pruned model, one can adjust the unpruned weights to compensate for the removed weights. Generally, the pruned model is formulated as:

$$\mathcal{F}(\mathbf{X}; \mathbf{W} \odot \mathbf{M} + \boldsymbol{\Delta}_{\mathbf{w}}), \tag{1}$$

140 where  $\Delta_{\mathbf{w}}$  (called *weight compensation*) is an update matrix 141 for the retained weights. Various methods have been pro-142 posed to enhance pruning by designing an effective scoring 143 metric or learning an effective weight compensation  $\Delta_{\mathbf{w}}$ , 144 e.g., Han et al. (2015); Zhang et al. (2024); Sun et al. (2024); 145 Dong et al. (2024); Das et al. (2023); An et al. (2024) for the 146 former, and Frantar & Alistarh (2023; 2022) for the latter. 147 We briefly review three representative pruning approaches. 148

Magnitude Pruning (Han et al., 2015) is the simplest tech-149 nique whose score matrix is defined as  $\mathbf{S}_{i,j} = |\mathbf{W}_{i,j}|$ , i.e., 150 removing the weights whose magnitudes are below a prede-151 fined threshold. In practice, magnitude pruning is performed 152 layer-wise: for each layer, a layer-dependent threshold is de-153 termined based on the local distribution of weights. Though 154 155 Magnitude pruning has stood out as a strong baseline for pruning models (Blalock et al., 2020), it has a major limi-156 tation: it ignores the importance of input activation, which 157 plays an equally importance role as weight magnitudes in 158 159 determining the output of linear layers (e.g., fully connected 160 layers, attention layers).

Wanda (Sun et al., 2024) addresses this limitation by incorporating both weights and inputs into defining the weight importance. Specifically, let  $\mathbf{X} \in \mathbb{R}^{N \times d_i}$  (where N is the sequence length) be the input activation of a training sample. Consider a linear layer  $\mathbf{Y} = \mathbf{X}\mathbf{W}$ , Wanda defines the importance of  $\mathbf{W}_{i,j}$  as  $\mathbf{S}_{i,j} = |\mathbf{W}_{i,j}| \cdot ||\mathbf{X}_{:,i}||_2$ .

**SparseGPT** (Frantar & Alistarh, 2023) introduces a more sophisticated pruning approach by incrementally pruning each column of **W**, followed by adjusting the remaining weights to compensate for those that have been pruned by the Optimal Brain Surgeon framework (Hassibi et al., 1993; Singh & Alistarh, 2020; Frantar et al., 2021). The score matrix is determined by  $\mathbf{S}_{i,j} = \frac{|\mathbf{W}_{i,j}|^2}{[\mathbf{H}^{-1}]_{i,i}}$  and  $\mathbf{H} = \mathbf{X}^{\top}\mathbf{X} + \lambda \mathbf{I}$  ( $\lambda$  is a small positive constant) is the Hessian matrix of the reconstruction loss.

## 4. Methodology

Different from existing methods, which primarily focus on enhancing pruning in the *parameter* space, e.g., learning a good scoring metric **S** or weight compensation  $\Delta_w$ , we propose enhancing pruning in the *input* space. Let  $\hat{W}$  be the weight of a pruned model. Our objective is to *determine* an input compensation  $\Delta_x$  for the input such that the output of the pruned model approximates that of the dense model, i.e.,

$$\mathcal{F}(\mathbf{X} + \mathbf{\Delta}_{\mathbf{x}}; \hat{\mathbf{W}}) \approx \mathcal{F}(\mathbf{X}; \mathbf{W}).$$
 (2)

The compensation  $\Delta_x$  depends on the input X. Obviously, learning  $\Delta_x$  from scratch for each input is inefficient. To deal with this issue, we begin by developing a framework of learning input compensation within the context of a simple linear layer and subsequently extend this approach to general models.

#### 4.1. Linear Models

Recent studies (Yu et al., 2017; Li et al., 2023b; Ding et al., 2023) have shown that the weight matrix **W** of neural networks can be approximated by a combination of a sparse

(3)

matrix  $\mathbf{S} \in \mathbb{R}^{d_i \times d_o}$  (assume rank $(\mathbf{S}) = d_o$ ) and a low-rank 165 matrix  $\mathbf{AB}^{\top}$  (where  $\mathbf{A} \in \mathbb{R}^{d_i \times r}$  and  $\mathbf{B} \in \mathbb{R}^{d_o \times r}$ , r is the 167 rank). Hence, for a linear layer, its output can be approxi-168 mated as:

169  
170 
$$\mathbf{Y} = \mathbf{X}\mathbf{W} \approx \mathbf{X}(\mathbf{S} + \mathbf{A}\mathbf{B}^{\top}) = \mathbf{X}\mathbf{S} + \mathbf{X}\mathbf{A}\underbrace{\mathbf{B}^{\top}(\mathbf{S}^{\top}\mathbf{S})^{-1}\mathbf{S}^{\top}}_{\hat{\mathbf{n}}}\mathbf{S}$$

 $= \left( \mathbf{X} + \underbrace{\mathbf{X} \mathbf{A} \hat{\mathbf{B}}^{\top}}_{i \in \mathbf{A}} \right) \mathbf{S}.$ 

172

- 173
- 174

175 176

177

178

179

180

181 182

183

184

185

186

187

188

211

212 213

214

215

216

217

218

219

Note that input compensation and weight compensation are dual for linear models, i.e., for an input X and its compensation  $\Delta_{\mathbf{x}}$ , one can design a weight compensation  $\Delta_{\mathbf{w}} = \mathbf{X}^{\top} (\mathbf{X} \mathbf{X}^{\top})^{-1} \Delta_{\mathbf{x}} \mathbf{W}$  (thus,  $\mathbf{X} \Delta_{\mathbf{w}} = \Delta_{\mathbf{x}} \mathbf{W}$ ) such that

$$(\mathbf{X} + \mathbf{\Delta}_{\mathbf{x}})\mathbf{W} = \mathbf{X}\mathbf{W} + \mathbf{\Delta}_{\mathbf{x}}\mathbf{W} = \mathbf{X}(\mathbf{W} + \mathbf{\Delta}_{\mathbf{w}}).$$
 (4)

However, for nonlinear models  $\mathcal{F}(\mathbf{X}; \mathbf{W})$ , the dual property does not hold, i.e., for an input compensation  $\Delta_x$ , there does not exist  $\Delta_{\mathbf{w}}$  such that  $\mathcal{F}(\mathbf{X} + \Delta_{\mathbf{x}}; \mathbf{W}) = \mathcal{F}(\mathbf{X}; \mathbf{W} + \mathbf{X})$  $\Delta_{\mathbf{w}}$ ).

189 Let  $\mathbf{a}_i$  and  $\hat{\mathbf{b}}_i$  be the *i*-th column of  $\mathbf{A}$  and  $\hat{\mathbf{B}}$ , respec-190 tively. Then, according to (3), the *i*-th row of  $\Delta_{\mathbf{x}}$  can 191 be re-expressed as  $\sum_{j=1}^{r} (\mathbf{x}_i^{\mathsf{T}} \mathbf{a}_j) \hat{\mathbf{b}}_j$ , which is similar to the 192 attention mechanism (Vaswani et al., 2017):  $\{x_i\}$  are the 193 query,  $\{\mathbf{a}_i\}$  are the keys, and  $\{\hat{\mathbf{b}}_i\}$  are the values. This ob-194 servation inspires us to design a framework of IC for general 195 models. 196

#### 197 4.2. General Models 198

199 Building on the above insight from the linear layer, we propose a general framework based on the attention mecha-200 nism for general models. The proposed procedure is shown 202 in Algorithm 1. Figure 1 provides an overview of the IC framework, consisting of a frozen encoder  $\mathcal{E}(\cdot)$  and a learn-203 able compensation pool  $(\mathbf{K}, \mathbf{V})$  (where  $\mathbf{K} \in \mathbb{R}^{d_e \times r}$  and 204  $\mathbf{V} \in \mathbb{R}^{r \times d_i}$ ). The encoder, which can either be a small submodule of the pruned model or an identity function, maps 206 **X** into an embedding  $\mathbf{Q}_{\mathbf{x}} = \mathcal{E}(\mathbf{X}) \in \mathbb{R}^{N \times d_e}$ , while the compensation pool consists of r learnable compensations. 208 209 Based on attention mechanism, the input compensation is 210 then constructed as:

$$\mathbf{\Delta}_{\mathbf{x}} = \operatorname{softmax}\left(\frac{\mathbf{Q}_{\mathbf{x}}\mathbf{K}}{\sqrt{d_e}}\right)\mathbf{V}.$$
 (5)

We add  $\Delta_{\mathbf{x}}$  to the input and learn the compensation pool by minimizing the following supervised loss:

$$\min_{\mathbf{K},\mathbf{V}} \sum_{(\mathbf{X},\mathbf{Y})\in\mathcal{D}} \ell(\mathcal{F}(\mathbf{X} + \boldsymbol{\Delta}_{\mathbf{x}}; \hat{\mathbf{W}}), \mathbf{Y}),$$
(6)

Algorithm 1 Input Compensation.

- **Require:** pruned model  $\mathcal{F}(\cdot; \hat{\mathbf{W}})$ , training set  $\mathcal{D}$ , rank r, step size  $\eta$ , encoder  $\mathcal{E}(\cdot)$ ;
- 1: initialize  $\mathbf{K} \in \mathbb{R}^{d_e \times r}$  and  $\mathbf{V} \in \mathbb{R}^{r \times d_i}$ ;
- 2: for t = 1, ..., T do
- sample an input **X** from  $\mathcal{D}$ ; 3:
- compute query embedding  $\mathbf{Q}_{\mathbf{x}} = \mathcal{E}(\mathbf{X})$ ; 4:
- 5:
- compute compensation  $\Delta_{\mathbf{x}} = \operatorname{softmax}\left(\frac{\mathbf{Q}_{\mathbf{x}}\mathbf{K}}{\sqrt{d_{e}}}\right)\mathbf{V};$ compute gradient  $\mathbf{g} = \nabla_{(\mathbf{K},\mathbf{V})}\ell(\mathcal{F}(\mathbf{X}+\Delta_{\mathbf{x}};\hat{\mathbf{W}}),\mathbf{Y})$ 6: (or  $\nabla_{(\mathbf{K},\mathbf{V})} \| \mathcal{F}(\mathbf{X} + \boldsymbol{\Delta}_{\mathbf{x}}; \hat{\mathbf{W}}) - \mathcal{F}(\mathbf{X}; \mathbf{W}) \|^2$  if **Y** is unavailable); 7:  $(\mathbf{K} \mathbf{V})$  $(\mathbf{K} \mathbf{V})$

$$f: (\mathbf{K}, \mathbf{v}) \leftarrow (\mathbf{K}, \mathbf{v}) - \eta \mathbf{g}$$

8: end for

9: return (**K**, **V**).

where  $\ell(\cdot, \cdot)$  is a loss function. In cases where labels of X are unavailable, we can learn the compensation pool by minimizing the reconstruction loss:

$$\min_{\mathbf{K},\mathbf{V}} \sum_{(\mathbf{X},\cdot)\in\mathcal{D}} \|\mathcal{F}(\mathbf{X} + \boldsymbol{\Delta}_{\mathbf{x}}; \hat{\mathbf{W}}) - \mathcal{F}(\mathbf{X}; \mathbf{W})\|^{2}.$$
 (7)

For NLP tasks, inputs (i.e., sentences) are sequences of discrete tokens, making direct modification of inputs infeasible. To deal with this issue, input compensation is operated in the input embedding space. Let  $\mathbf{H}_{\mathbf{x}} \in \mathbb{R}^{N imes d_e}$  be the embeddings extracted by the input embedding layer of the pruned LLM. Similar to (5), we construct the input compensation as  $\Delta_{\mathbf{x}} = \operatorname{softmax}\left(\frac{\mathbf{H}_{\mathbf{x}}\mathbf{K}}{\sqrt{d_e}}\right)\mathbf{V}.$  The input embeddings are then adjusted as  $\mathbf{H} + \Delta_{\mathbf{x}}$  and we learn the compensation pool by minimizing the reconstruction loss of the last hidden states.

#### 5. Experiments

#### 5.1. Experiments on Image Classification

Datasets. We conduct image classification experiments on ten datasets: CIFAR100 (Krizhevsky & Hinton, 2009), Flowers (Nilsback & Zisserman, 2008), Food (Bossard et al., 2014), EuroSAT (Helber et al., 2019), SUN (Xiao et al., 2016), DTD (Cimpoi et al., 2014), UCF (Soomro et al., 2012), SVHN (Netzer et al., 2011), OxfordPets (Jawahar et al., 2012) (denoted by Pets), and RESISC45 (Cheng et al., 2017) (denoted by RESISC). A summary of the datasets is in Table 10 of Appendix B.6.

Implementation Details. We adopt CLIP ViT-B/32 and ViT-B/16 (Radford et al., 2021) as the base models, whose pruned image encoder is used as the encoder of IC. We initialize the K and V by the standard normal distribution and train the compensation pool for 30 epochs using the SGD optimizer with a learning rate of 40 and momentum

Limancing I I uncu mouchs by imput compensation	Enhancing	Pruned	Models	by Inp	ut Com	pensatior
---	-----------	--------	--------	--------	--------	-----------

	Sparsity	#Nonzero Params	CIFAR100	Flowers	Food	EuroSAT	SUN	UCF	SVHN	Pets	DTD	RESISC	Ave
Magnitude	50%	110M	33.9	26.1	34.2	45.6	30.8	35.4	45.3	38.7	27.9	55.4	37 3
Magnitude	48%	112M	43.4	30.9	44.5	61.3	36.3	43.9	56.7	48.9	31.4	64.8	46.2
Magnitude + IC	50%	112M	73.0	62.9	72.4	96.5	48.9	63.1	94.4	69.2	44.1	87.1	71.2
Wanda	50%	110M	75.0	56.4	74.1	95.2	50.8	59.7	91.8	57.6	43.4	84.4	68.9
Wanda	48%	112M	77.6	63.9	77.5	95.9	55.2	65.0	89.3	65.1	45.7	87.9	72.3
Wanda + IC	50%	112M	80.1	76.4	80.4	97.9	54.7	69.1	96.1	77.5	49.8	91.6	77.4
SparseGPT	50%	110M	83.3	69.1	81.6	97.9	58.0	68.5	93.7	59.4	48.2	89.8	74.9
SparseGPT	48%	112M	84.5	72.4	82.7	97.9	60.4	71.7	91.2	66.3	50.2	91.3	76.9
SparseGPT + IC	50%	112M	82.9	76.2	83.1	98.2	57.2	71.0	96.7	79.7	53.8	92.9	79.2
Magnitude	4:8	110M	49.0	25.9	36.5	45.1	32.8	37.8	60.8	45.3	27.1	60.2	42.1
Magnitude + IC	4:8	112M	72.9	62.4	72.1	96.5	48.2	62.9	94.3	68.3	44.8	87.4	71.0
Wanda	4:8	110M	60.9	30.5	59.2	83.1	37.2	43.2	74.4	47.0	30.9	68.7	53.5
Wanda + IC	4:8	112M	76.9	71.4	77.3	97.2	50.2	64.2	95.2	75.3	51.1	89.5	74.8
SparseGPT	4:8	110M	80.2	55.7	79.6	96.6	52.3	61.4	85.8	58.5	42.3	86.6	69.9
SparseGPT + IC	4:8	112M	81.8	72.6	81.6	98.1	51.7	65.8	96.6	78.1	49.1	92.2	76.8
Magnitude	2:4	110M	30.7	11.2	19.8	42.7	19.4	23.3	27.5	25.8	15.6	35.3	25.1
Magnitude + IC	2:4	112M	78.6	66.5	77.9	97.4	53.9	67.5	96.1	70.9	44.3	90.2	74.3
Wanda	2:4	110M	39.6	14.5	35.1	46.1	21.7	25.5	42.4	25.2	20.9	41.4	31.2
Wanda + IC	2:4	112M	80.6	74.5	80.7	97.7	54.9	68.1	96.5	74.1	51.4	91.6	77.0
SparseGPT	2:4	110M	75.5	40.2	73.0	94.3	44.5	52.6	61.3	45.7	33.5	81.6	60.2
puise of 1	2	110101	10.0		1010	2110		02.0	0110	1017	0010	0110	
SparseGPT + IC	2:4	112M	82.2	79.4	82.6	98.4	57.7	69.9	96.8	76.4	54.3	92.5	79.0
SparseGPT + IC	2:4	112M	82.2	79.4	82.6	98.4	57.7	69.9	96.8	76.4	54.3	92.5	79.(
SparseGPT + IC	2:4 Table	112M 2: Testing acc	<sup>82.2</sup> uracy on i	<sup>79.4</sup> mage cl	82.6 assifi	98.4 cation ta	57.7 sks us	<sup>69.9</sup>	96.8 CLIP Vi	76.4 iT-B/	<sup>54.3</sup>	92.5	79.0
SparseGPT + IC	2:4 Table Sparsity	112M 2: Testing accu #Nonzero Params	82.2 aracy on i CIFAR100	79.4 mage cl Flowers	82.6 assifi Food	98.4 cation ta EuroSAT	57.7 sks us sun	69.9 sing C UCF	96.8 CLIP Vi SVHN	76.4 iT-B/ Pets	54.3 16. DTD	92.5 RESISC	79.0 Avg
SparseGPT + IC	2:4 Table Sparsity 50%	112M 2: Testing accu #Nonzero Params 109M	82.2 1racy on i CIFAR100 76.9	79.4 mage cl Flowers 56.5	82.6 assifi Food 78.3	98.4 cation ta EuroSAT 90.7	57.7 sks us SUN 51.2	69.9 5 OCF 65.6	96.8 CLIP Vi SVHN 95.3	76.4 iT-B/ Pets 62.9	54.3 16. DTD 42.8	92.5 RESISC 82.1	79.0 Avg 70.2
SparseGPT + IC Magnitude Magnitude	2:4 Table Sparsity 50% 47%	112M 2: Testing accu #Nonzero Params 109M 111M	82.2 Jracy on i CIFAR100 76.9 80.9	79.4 mage cl Flowers 56.5 64.6	82.6 assifi Food 78.3 82.1	98.4 cation ta EuroSAT 90.7 94.4	57.7 sks us sun 51.2 56.5	69.9 5000 CF 65.6 70.1 75.1	96.8 CLIP Vi SVHN 95.3 95.9	76.4 iT-B/ Pets 62.9 68.5	54.3 16. DTD 42.8 47.3 (1.1)	92.5 RESISC 82.1 86.9	79.0 Avg 70.2 74.2
SparseGPT + IC Magnitude Magnitude Magnitude + IC	2:4 Table Sparsity 50% 47% 50%	112M 2: Testing acct #Nonzero Params 109M 111M 111M	82.2 1racy on i CIFAR100 76.9 80.9 82.9	79.4 mage cl Flowers 56.5 64.6 86.7	82.6 Assific Food 78.3 82.1 84.7	98.4 cation ta EuroSAT 90.7 94.4 97.6	57.7 sks us sun 51.2 56.5 60.1	69.9 5000 CF 65.6 70.1 75.1	96.8 CLIP Vi SVHN 95.3 95.9 97.1	76.4 iT-B/ Pets 62.9 68.5 82.5	54.3 16. DTD 42.8 47.3 61.1	92.5 RESISC 82.1 86.9 92.8	79.0 Avg 70.2 74.7 82.1
SparseGPT + IC Magnitude Magnitude Magnitude + IC Wanda	2:4 Table Sparsity 50% 47% 50% 50%	112M 2: Testing acct #Nonzero Params 109M 111M 111M 109M	82.2 1racy on i CIFAR100 76.9 80.9 82.9 84.1 85.6	79.4 mage cl Flowers 56.5 64.6 86.7 78.1	82.6 Food 78.3 82.1 84.7 85.5	98.4 cation ta EuroSAT 90.7 94.4 97.6 97.6 97.6	57.7 sks us SUN 51.2 56.5 60.1 59.5 62.2	69.9 5000 CF 65.6 70.1 75.1 68.9 71.7	96.8 CLIP Vi SVHN 95.3 95.9 97.1 96.9 96.9 96.9	76.4 iT-B/ Pets 62.9 68.5 82.5 72.7 72.7	54.3 16. DTD 42.8 47.3 61.1 51.8 57.2	92.5 RESISC 82.1 86.9 92.8 91.2	79.0 Avg 70.2 74.7 82.1 78.0
SparseGPT + IC Magnitude Magnitude Magnitude + IC Wanda Wanda	2:4 Table Sparsity 50% 47% 50% 50% 47% 50%	112M 2: Testing acct #Nonzero Params 109M 111M 111M 109M 111M	82.2 Jracy on i CIFAR100 76.9 80.9 82.9 84.1 85.6 86.2	79.4 mage cl Flowers 56.5 64.6 86.7 78.1 81.7 82.8	82.6 assifi Food 78.3 82.1 84.7 85.5 86.8 87.8	98.4 cation ta EuroSAT 90.7 94.4 97.6 97.6 98.0 98.4	57.7 sks us SUN 51.2 56.5 60.1 59.5 62.8 63.8	69.9 UCF 65.6 70.1 75.1 68.9 71.7 75.5	96.8 <b>CLIP V</b> SVHN 95.3 95.9 97.1 96.9 96.8 97.6	76.4 iT-B/ Pets 62.9 68.5 82.5 72.7 76.8 83.6	54.3 16. DTD 42.8 47.3 61.1 51.8 55.3 63.5	92.5 RESISC 82.1 86.9 92.8 91.2 92.9 94.7	79.0 Avg 70.2 74.7 82.1 78.0 80.8 83.4
Magnitude Magnitude Magnitude + IC Wanda Wanda + IC	2:4 Table Sparsity 50% 47% 50% 47% 50% 47% 50%	112M 2: Testing acct #Nonzero Params 109M 111M 109M 111M 109M 111M	82.2 Iracy on i CIFAR100 76.9 80.9 82.9 84.1 85.6 86.2	79.4 mage cl Flowers 56.5 64.6 86.7 78.1 81.7 82.8	82.6 Assift Food 78.3 82.1 84.7 85.5 86.8 87.8	98.4 cation ta EuroSAT 90.7 94.4 97.6 98.0 98.4 98.4	57.7 sks us SUN 51.2 56.5 60.1 59.5 62.8 63.8 (2.6)	69.9 <b>UCF</b> 65.6 70.1 75.1 68.9 71.7 75.5 <b>72.6</b>	96.8 <b>CLIP Vi</b> SVHN 95.3 95.9 97.1 96.9 96.8 97.6 02.0	76.4 iT-B/ Pets 62.9 68.5 82.5 72.7 76.8 83.6	54.3 16. DTD 42.8 47.3 61.1 51.8 55.3 63.5 55.4	92.5 RESISC 82.1 86.9 92.8 91.2 92.9 94.7	79.0 Avg 70.2 74.3 82.1 78.0 80.8 83.4
SparseGPT + IC Magnitude Magnitude + IC Wanda Wanda + IC SparseGPT SparseGPT	2:4 Table Sparsity 50% 47% 50% 50% 47% 50% 47%	112M 2: Testing acct #Nonzero Params 109M 111M 109M 111M 109M 111M	82.2 1racy on i CIFAR100 76.9 80.9 82.9 84.1 85.6 86.2 87.2 87.7	79.4 mage cl Flowers 56.5 64.6 86.7 78.1 81.7 82.8 80.2 82.2	82.6 assifu Food 78.3 82.1 84.7 85.5 86.8 87.8 88.1 88.1 89.9	98.4 cation ta EuroSAT 90.7 94.4 97.6 98.0 98.4 98.0 98.4	57.7 sks us SUN 51.2 56.5 60.1 59.5 62.8 63.8 63.8 63.8	69.9 <b>Sing C</b> UCF 65.6 70.1 75.1 68.9 71.7 75.5 73.8 76.2	96.8 CLIP VI SVHN 95.3 95.9 97.1 96.9 96.8 97.6 97.0 97.0	76.4 <b>iT-B/</b> Pets 62.9 68.5 82.5 72.7 76.8 83.6 75.6 80.4	54.3 16. DTD 42.8 47.3 61.1 51.8 55.3 63.5 56.4 50.8	92.5 RESISC 82.1 86.9 92.8 91.2 92.9 94.7 93.7 94.7	Avg           70.1           74.1           82.1           78.0           83.0           81.4           82.2
SparseGPT + IC Magnitude Magnitude + IC Wanda Wanda + IC SparseGPT SparseGPT - SparseGPT - SparseGPT - SparseGPT - SparseGPT - IC	2:4 Table Sparsity 50% 47% 50% 50% 47% 50% 50% 47% 50%	112M 2: Testing acct #Nonzero Params 109M 111M 109M 111M 109M 111M 109M 111M	82.2 1racy on i CIFAR100 76.9 80.9 82.9 84.1 85.6 86.2 87.2 87.7 86.1	79.4 mage cl Flowers 56.5 64.6 86.7 78.1 81.7 82.8 80.2 82.3 86.0	82.6 Food 78.3 82.1 84.7 85.5 86.8 87.8 88.1 88.8 88.1 88.8 87.9	98.4 cation ta EuroSAT 90.7 94.4 97.6 97.6 98.0 98.4 98.0 98.3 98.4	57.7 sks us SUN 51.2 56.5 60.1 59.5 62.8 63.8 63.8 63.4 64.4	69.9 Sing C UCF 65.6 70.1 75.1 68.9 71.7 75.5 73.8 76.2 76.2	96.8 <b>CLIP V</b> <b>SVHN</b> 95.3 95.9 97.1 96.9 96.8 97.6 97.0 97.0 97.0 97.0	76.4 <b>iT-B/</b> Pets 62.9 68.5 82.5 72.7 76.8 83.6 75.6 80.4 85.2	54.3 16. DTD 42.8 47.3 61.1 51.8 55.3 63.5 56.4 59.8 66.4	92.5 RESISC 82.1 86.9 92.8 91.2 92.9 94.7 93.7 94.6 95.0	Avg           70.2           74.7           82           78.0           83.4           81.4           82.9           84
Magnitude Magnitude Magnitude + IC Wanda Wanda + IC SparseGPT SparseGPT + IC	2:4 Table Sparsity 50% 47% 50% 50% 47% 50% 47% 50%	112M 2: Testing acct #Nonzero Params 109M 111M 109M 111M 109M 111M 109M 111M	82.2 Iracy on i CIFAR100 76.9 80.9 82.9 84.1 85.6 86.2 87.2 87.7 86.1 77.7 86.1	79.4 mage cl Flowers 56.5 64.6 86.7 78.1 81.7 82.8 80.2 82.3 86.0 86.0	82.6 assifi Food 78.3 82.1 84.7 85.5 86.8 87.8 88.1 88.8 88.1 88.8 87.9	98.4 cation ta EuroSAT 90.7 94.4 97.6 97.6 98.0 98.4 98.0 98.3 98.4	57.7 sks us SUN 51.2 56.5 60.1 59.5 62.8 63.8 63.8 63.4 64.4	69.9 UCF 65.6 70.1 75.1 68.9 71.7 75.5 73.8 76.2 76.2 76.2	96.8 CLIP Vi SVHN 95.3 95.9 97.1 96.9 96.8 97.6 97.0 97.0 97.0 97.6	76.4 <b>iT-B/</b> Pets 62.9 68.5 82.5 72.7 76.8 83.6 75.6 80.4 85.2	54.3 16. DTD 42.8 47.3 61.1 51.8 55.3 63.5 56.4 59.8 66.4	92.5 RESISC 82.1 86.9 92.8 91.2 92.9 94.7 93.7 94.6 95.0	79.0 Avg 70.1 74.1 82.1 83.4 81.4 82.9 84.1 84.1
Magnitude Magnitude Magnitude + IC Wanda Wanda + IC SparseGPT SparseGPT + IC Magnitude + IC	2:4 Table Sparsity 50% 47% 50% 50% 47% 50% 50% 47% 50% 50% 47% 50% 50% 50% 50% 50% 50% 50% 50	112M 2: Testing acci #Nonzero Params 109M 111M 109M 111M 109M 111M 109M 111M 111	82.2 Iracy on i CIFAR100 76.9 80.9 82.9 84.1 85.6 86.2 87.2 87.7 86.1 75.8 81.5	79.4 mage cl Flowers 56.5 64.6 86.7 78.1 81.7 82.8 80.2 82.3 86.0 52.0	82.6 assifi Food 78.3 82.1 84.7 85.5 86.8 87.8 88.1 88.8 87.9 75.4 82.2	98.4 cation ta EuroSAT 90.7 94.4 97.6 97.6 98.0 98.4 98.0 98.3 98.4 98.4 98.5 98.4	57.7 sks us SUN 51.2 56.5 60.1 59.5 62.8 63.8 63.8 63.4 64.4 50.0 57.6	69.9 UCF 65.6 70.1 75.1 68.9 71.7 75.5 73.8 76.2 76.2 61.4 72.5	96.8 SVHN 95.3 95.9 97.1 96.9 96.8 97.6 97.0 97.0 97.0 97.6 77.4 96.9 97.0 97.0 97.6 97.0 97.6	76.4 Pets 62.9 68.5 82.5 72.7 76.8 83.6 75.6 80.4 85.2 68.8 81.6	54.3 16. DTD 42.8 47.3 61.1 51.8 55.3 63.5 56.4 59.8 66.4 41.4 51.4	92.5 RESISC 82.1 86.9 92.8 91.2 92.9 94.7 93.7 94.6 95.0 79.1	Avg           70.           74.           82.           78.           83.           81           82.           84.           67
SparseGPT + IC Magnitude Magnitude + IC Wanda Wanda + IC SparseGPT SparseGPT + IC Magnitude + IC	2:4 Table Sparsity 50% 47% 50% 50% 47% 50% 47% 50% 47% 50% 47% 50% 47% 50% 47% 50% 47% 50% 47% 50% 47% 50% 47% 50% 50% 50% 50% 50% 50% 50% 50	112M           2: Testing accl           #Nonzero Params           109M           111M	82.2 Iracy on i CIFAR100 76.9 80.9 82.9 84.1 85.6 86.2 87.2 87.7 86.1 75.8 81.5	79.4 mage cl Flowers 56.5 64.6 86.7 78.1 81.7 82.8 80.2 82.3 86.0 52.0 84.2	82.6 Food 78.3 82.1 84.7 85.5 86.8 87.8 88.1 88.8 87.9 75.4 83.2	98.4 cation ta EuroSAT 90.7 94.4 97.6 97.6 98.0 98.4 98.0 98.3 98.4 89.6 97.5 97.5	57.7 <b>Sks us</b> <b>SUN</b> 51.2 56.5 60.1 59.5 62.8 63.8 63.8 63.8 63.4 64.4 50.0 57.6 57.6	69.9 UCF 65.6 70.1 75.1 68.9 71.7 75.5 73.8 76.2 76.2 61.4 72.5	96.8 SULIP Vi SVHN 95.3 95.9 97.1 96.9 96.8 97.6 97.0 97.0 97.0 97.6 77.4 96.9 97.4 96.9 97.4 96.9 97.1 96.9 97.1 96.9 97.1 96.9 97.1 96.9 97.1 96.9 97.1 97.0 97	76.4 Pets 62.9 68.5 82.5 72.7 76.8 83.6 75.6 80.4 85.2 68.8 81.6	54.3 16. DTD 42.8 47.3 61.1 51.8 55.3 63.5 56.4 59.8 66.4 41.4 54.4	92.5 RESISC 82.1 86.9 92.8 91.2 92.9 94.7 93.7 94.6 95.0 79.1 91.9 91.9 95.0	Av           70.           74.           82.           78.           80.           83.           81.           82.           84.           67.           80.
Magnitude Magnitude Magnitude + IC Wanda Wanda + IC SparseGPT SparseGPT + IC Magnitude Magnitude + IC Wanda	2:4 Table Sparsity 50% 47% 50% 50% 47% 50% 50% 47% 50% 50% 50% 50% 50% 50% 50% 50	112M 2: Testing acci #Nonzero Params 109M 111M 109M 111M 109M 111M 109M 111M 109M 111M	82.2 Iracy on i CIFAR100 76.9 80.9 82.9 84.1 85.6 86.2 87.2 87.7 86.1 75.8 81.5 78.6 78.6	79.4 mage cl Flowers 56.5 64.6 86.7 78.1 81.7 82.8 80.2 82.3 86.0 52.0 84.2 63.2	82.6 Food 78.3 82.1 84.7 85.5 86.8 87.8 88.1 88.8 87.9 75.4 83.2 81.4	98.4 cation ta EuroSAT 90.7 94.4 97.6 97.6 98.0 98.4 98.0 98.3 98.4 89.6 97.5 96.0 96.0 96.1	57.7 <b>Sks us</b> <b>SUN</b> 51.2 56.5 60.1 59.5 62.8 63.8 63.8 63.8 63.4 64.4 50.0 57.6 50.6 50.6	69.9 UCF 65.6 70.1 75.1 68.9 71.7 75.5 73.8 76.2 76.2 61.4 72.5 61.3 71.3	96.8 SVHN 95.3 95.9 97.1 96.9 96.8 97.6 97.0 97.0 97.6 77.4 96.9 77.4 96.9 77.4 96.9 78.2 77.4	76.4 Pets 62.9 68.5 82.5 72.7 76.8 83.6 75.6 80.4 85.2 68.8 81.6 69.5	54.3 16. DTD 42.8 47.3 61.1 51.8 55.3 63.5 56.4 59.8 66.4 41.4 54.4 42.8 41.4 54.4 42.8 41.4 54.8 55.3 55.4 55.8 55.8 55.4 55.8 55	92.5 RESISC 82.1 86.9 92.8 91.2 92.9 94.7 93.7 94.6 95.0 79.1 91.9 87.7	Av           70.           74.           82.           78.           80.           83.           81.           82.           84.           67.           80.           70.           70.           70.           70.           70.           70.           70.
SparseGPT + IC Magnitude Magnitude Magnitude + IC Wanda Wanda + IC SparseGPT SparseGPT + IC Magnitude Magnitude + IC Wanda Wanda + IC	2:4 Table Sparsity 50% 47% 50% 50% 50% 50% 47% 50% 47% 50% 47% 50% 478 4:8 4:8 4:8 4:8	112M           2: Testing acci           #Nonzero Params           109M           111M	82.2 Iracy on i CIFAR100 76.9 80.9 82.9 84.1 85.6 86.2 87.2 87.7 86.1 75.8 81.5 78.6 84.7	79.4           mage cl           Flowers           56.5           64.6           86.7           78.1           81.7           82.8           80.2           82.3           86.0           52.0           84.2           63.2           82.1	82.6 Food 78.3 82.1 84.7 85.5 86.8 87.8 88.1 88.8 87.9 75.4 83.2 81.4 86.8	98.4 cation ta EuroSAT 90.7 94.4 97.6 97.6 98.0 98.4 98.0 98.3 98.4 98.6 97.5 96.0 98.4	57.7           sks us           sUN           51.2           56.5           60.1           59.5           62.8           63.8           63.4           64.4           50.0           57.6           50.6           60.6	69.9 UCF 65.6 70.1 75.1 68.9 71.7 75.5 73.8 76.2 76.2 61.4 72.5 61.3 74.5	96.8 <b>CLIP V</b> <b>SVHN</b> 95.3 95.9 97.1 96.9 96.8 97.6 97.0 97.0 97.0 97.0 97.0 97.0 97.0 97.0 97.0 97.0 97.0 97.0 97.0 97.0 97.0 97.0 97.0 97.0 97.2 97.4 96.9	76.4 Pets 62.9 68.5 82.5 72.7 76.8 83.6 75.6 80.4 85.2 68.8 81.6 69.5 82.0	54.3           16.           DTD           42.8           47.3           61.1           51.8           55.3           63.5           56.4           59.8           66.4           41.4           54.3           61.3	92.5 RESISC 82.1 86.9 92.8 91.2 92.9 94.7 93.7 94.6 95.0 79.1 91.9 87.7 94.3	Avv           70.           74.           82.           78.           80.           81.           82.           67.           80.           70.           82.
Magnitude Magnitude Magnitude Magnitude + IC Wanda Wanda + IC SparseGPT SparseGPT + IC Magnitude Magnitude + IC Wanda Wanda + IC SparseGPT	2:4 Table Sparsity 50% 47% 50% 50% 50% 50% 47% 50% 47% 50% 47% 50% 478 418 418 418 418 418 418	112M           2: Testing accl           #Nonzero Params           109M           111M	82.2 Iracy on i CIFAR100 76.9 80.9 82.9 84.1 85.6 86.2 87.2 87.7 86.1 75.8 81.5 78.6 84.7 85.1 85.1	79.4 mage cl Flowers 56.5 64.6 86.7 78.1 81.7 82.8 80.2 82.3 86.0 52.0 84.2 63.2 82.1 74.0	82.6 Food 78.3 82.1 84.7 85.5 86.8 87.8 88.1 88.8 87.9 75.4 83.2 81.4 86.8 87.0	98.4 cation ta EuroSAT 90.7 94.4 97.6 97.6 98.0 98.4 98.0 98.3 98.4 98.6 97.5 96.0 98.4 99.6 95.6 95.6	57.7 <b>sks us</b> <b>sun</b> 51.2 56.5 60.1 59.5 62.8 63.8 63.8 63.8 63.4 64.4 50.0 57.6 50.6 60.6 60.6 60.4 60.4	69.9 UCF 65.6 70.1 75.1 68.9 71.7 75.5 73.8 76.2 76.2 76.2 61.4 72.5 61.3 74.5 69.8	96.8 SULIP Vi SVHN 95.3 95.9 97.1 96.9 96.8 97.6 97.0 97.0 97.0 97.0 97.0 97.6 77.4 96.9 78.2 97.4 78.2 97.4 72.6	76.4 TT-B/ Pets 62.9 68.5 82.5 72.7 76.8 83.6 75.6 80.4 85.2 68.8 81.6 69.5 82.0 78.2.0 78.2 78.2 78.2 78.2 78.2 78.2 78.2 79.7 75.6 75.7 75.6 75.6 75.7 75.6 75.7 75.	54.3           16.           DTD           42.8           47.3           61.1           51.8           55.3           63.5           56.4           59.8           66.4           41.4           54.3           50.5           50.5	92.5 RESISC 82.1 86.9 92.8 91.2 92.9 94.7 93.7 94.6 95.0 79.1 91.9 87.7 94.3 93.7 94.3	Av,           70.           74.           82.           78.           83.           81.           82.           84.           677.           80.           70.           70.           81.           82.           84.           677.           80.           70.           82.           76.           76.
SparseGPT + IC Magnitude Magnitude + IC Wanda Wanda + IC SparseGPT SparseGPT + IC Magnitude + IC Wanda Magnitude + IC Wanda Wanda + IC SparseGPT SparseGPT SparseGPT SparseGPT + IC	2:4 Table Sparsity 50% 47% 50% 50% 47% 50% 47% 50% 47% 50% 47% 50% 478 418 418 418 418 418 418 418 41	112M           2: Testing accl           #Nonzero Params           109M           111M	82.2 Iracy on i CIFAR100 76.9 80.9 82.9 84.1 85.6 86.2 87.2 87.7 86.1 75.8 81.5 78.6 84.7 85.1 84.7	79.4           mage cl           Flowers           56.5           64.6           86.7           78.1           81.7           82.8           80.2           82.3           86.0           52.0           84.2           63.2           82.1           74.0           84.9	82.6 assifi Food 78.3 82.1 84.7 85.5 86.8 87.8 88.1 88.8 87.9 75.4 83.2 81.4 86.8 87.0 87.1	98.4 EuroSAT 90.7 94.4 97.6 97.6 98.0 98.4 98.0 98.3 98.4 98.6 97.5 96.0 98.4 99.5 96.0 98.4 99.5 96.0 98.4	57.7 <b>sks us</b> <b>sun</b> 51.2 56.5 60.1 59.5 62.8 63.8 63.8 63.4 64.4 50.0 57.6 50.6 60.6 60.6 60.4 60.9	69.9           UCF           65.6           70.1           75.7           68.9           71.7           75.5           73.8           76.2           61.4           72.5           61.3           74.5           69.8           74.5	96.8 <b>CLIP V</b> <b>SVHN</b> 95.3 95.9 97.1 96.9 96.8 97.6 97.0 97.0 97.0 97.0 97.0 97.0 97.0 97.0 97.0 97.0 97.0 97.0 97.0 97.0 97.0 97.0 97.0 97.0 97.1 97.0 97.2 97.4 96.9 97.4 96.9 97.2 97.4 97.5 97.4 97.5 97.4 97.5 97.4 97.5 97.4 97.5	76.4 T-B/ Pets 62.9 68.5 82.5 72.7 76.8 83.6 75.6 80.4 85.2 68.8 81.6 69.5 82.0 78.2 83.6	54.3           16.           DTD           42.8           47.3           61.1           51.8           55.3           63.5           56.4           59.8           66.4           41.4           54.3           50.5           62.2	92.5 RESISC 82.1 86.9 92.8 91.2 92.9 94.7 93.7 94.6 95.0 79.1 91.9 87.7 94.3 93.7 94.4	Av           70.           74.           82.           78.           80.           83.           81.           82.           84.           67.           80.           70.           70.           82.           84.           67.           80.           70.           82.           76.           82.
SparseGPT + IC Magnitude Magnitude Magnitude + IC Wanda Wanda + IC SparseGPT SparseGPT + IC Magnitude + IC Wanda Wanda + IC SparseGPT SparseGPT SparseGPT SparseGPT + IC Magnitude Magnitude	2:4 Table Sparsity 50% 47% 50% 50% 47% 50% 50% 47% 50% 50% 50% 50% 50% 50% 50% 50	112M           2: Testing acci           #Nonzero Params           109M           111M	82.2 Iracy on i CIFAR100 76.9 80.9 82.9 84.1 85.6 86.2 87.2 87.7 86.1 75.8 81.5 78.6 84.7 85.1 84.7 85.1 84.7 68.3	79.4           mage cl           Flowers           56.5           64.6           86.7           78.1           81.7           82.8           80.2           82.3           86.0           52.0           84.2           63.2           82.1           74.0           84.9           39.9	82.6 assift Food 78.3 82.1 84.7 85.5 86.8 87.8 88.1 88.8 87.9 75.4 83.2 81.4 86.8 87.0 87.1 64.5	98.4 EuroSAT 90.7 94.4 97.6 97.6 98.0 98.4 98.0 98.3 98.4 98.6 97.5 96.0 98.4 95.6 98.3 78.6	57.7 <b>sks us</b> <b>sun</b> 51.2 56.5 60.1 59.5 62.8 63.8 63.8 63.4 64.4 50.0 57.6 50.6 60.6 60.6 60.4 60.9 39.4	69.9           UCF           65.6           70.1           75.7           68.9           71.7           75.5           73.8           76.2           61.4           72.5           61.3           74.5           69.8           74.9           53.2	96.8 <b>CLIP V</b> <b>SVHN</b> 95.3 95.9 97.1 96.9 96.8 97.6 97.0 97.6 77.4 96.9 78.2 97.4 97.2 97.4 95.3 95.9 97.3 97.0 97.0 97.0 97.0 97.0 97.0 97.0 97.0 97.0 97.0 97.0 97.0 97.0 97.0 97.0 97.5 97.2 97.4 95.2 97.4 97.5 97.4 97.5 97.5 97.5 95.3	76.4 T-B/ Pets 62.9 68.5 82.5 72.7 76.8 83.6 75.6 80.4 85.2 68.8 81.6 69.5 82.0 78.2 83.6 78.2 83.6 56.9 57.9 57.6 57.6 57.6 57.6 57.6 57.6 57.6 57.6 57.6 57.6 57.6 57.6 57.6 57.6 57.6 57.6 57.6 57.6 57.7 57.6 57.6 57.6 57.6 57.9 57.9 57.6 57.9 57.6 57.9 5	54.3         116.         DTD         42.8         47.3         61.1         51.8         55.3         63.5         56.4         59.8         66.4         41.4         54.3         50.5         62.2         32.9	92.5 RESISC 82.1 86.9 92.8 91.2 92.9 94.7 93.7 94.6 95.0 79.1 91.9 87.7 94.3 93.7 94.3 93.7 94.4 67.6	Av           70.           74.           82.           78.           80.           83.           81.           82.           84.           67.           80.           70.           82.           76.           82.           76.           82.           59.
SparseGPT + IC Magnitude Magnitude Magnitude + IC Wanda Wanda + IC SparseGPT SparseGPT + IC Magnitude + IC Wanda + IC SparseGPT + IC SparseGPT + IC SparseGPT + IC SparseGPT + IC Magnitude Magnitude + IC	2:4 Table Sparsity 50% 47% 50% 50% 50% 50% 47% 47% 50% 47% 47% 47% 47% 47% 47% 47% 47	112M           2: Testing acci           #Nonzero Params           109M           111M	82.2 Iracy on i CIFAR100 76.9 80.9 82.9 84.1 85.6 86.2 87.2 87.7 86.1 75.8 81.5 78.6 84.7 85.1 84.7 68.3 83.9	79.4           mage cl           Flowers           56.5           64.6           86.7           78.1           81.7           82.8           80.2           82.3           86.0           52.0           84.2           63.2           82.1           74.0           84.9           39.9           83.8	82.6 assift Food 78.3 82.1 84.7 85.5 86.8 87.8 88.1 88.8 87.9 75.4 83.2 81.4 86.8 87.0 87.1 64.5 85.9	98.4 EuroSAT 90.7 94.4 97.6 97.6 98.0 98.4 98.0 98.3 98.4 98.6 97.5 96.0 98.4 95.6 98.3 98.4 95.6 98.3 78.6 97.9	57.7 <b>sks us</b> <b>sun</b> 51.2 56.5 60.1 59.5 62.8 63.8 63.8 63.4 64.4 50.0 57.6 50.6 60.6 60.6 60.4 60.9 39.4 61.4	69.9           UCF           65.6           70.1           75.7           68.9           71.7           75.5           73.8           76.2           61.4           72.5           61.3           74.5           69.8           74.9           53.2           75.7	96.8 <b>CLIP V</b> <b>SVHN</b> 95.3 95.9 97.1 96.9 96.8 97.6 97.0 97.5 97.1 96.9 97.0 97.0 97.0 97.0 97.0 97.0 97.0 97.5 97.4 96.9 97.4 96.9 97.4 97.5 97.4 97.5 97.4 97.5 97.4 97.5 97.4 97.5 97.4 97.5 97.3 97.4 97.5 97.4 97.5 97.4 97.5 97.4 97.5 97.4 97.5 97.4 97.5 97.4 97.5 97.3 97.4 97.5 97.4 97.5 97.3 97.4 97.5 97.5 97.3 97.4	76.4 Pets 62.9 68.5 82.5 72.7 76.8 83.6 75.6 80.4 85.2 68.8 81.6 69.5 82.0 78.2 83.6 78.2 83.6 56.9 84.0 56.9 84.0	54.3           16.           DTD           42.8           47.3           61.1           51.8           55.3           63.5           56.4           59.8           66.4           41.4           54.3           50.5           62.2           32.9           57.3	92.5 RESISC 82.1 86.9 92.8 91.2 92.9 94.7 93.7 94.6 95.0 79.1 91.9 87.7 94.3 93.7 94.3 93.7 94.4 67.6 93.3	Av;           70.           74.           82.           78.           80.           83.           81.           82.           67.           80.           76.           82.           59.           82.
SparseGPT + IC Magnitude Magnitude Magnitude + IC Wanda Wanda + IC SparseGPT SparseGPT + IC Magnitude + IC Wanda Wanda + IC SparseGPT + IC SparseGPT + IC SparseGPT + IC Magnitude Magnitude + IC Magnitude Magnitude + IC Magnitude + IC Magnitude + IC Magnitude + IC Magnitude + IC Wanda	2:4 Table Sparsity 50% 47% 50% 50% 50% 50% 47% 47% 50% 47% 47% 47% 47% 47% 47% 47% 47	112M         2: Testing accl         #Nonzero Params         109M         111M         109M	82.2 Iracy on i CIFAR100 76.9 80.9 82.9 84.1 85.6 86.2 87.2 87.7 86.1 75.8 81.5 78.6 84.7 85.1 84.7 68.3 83.9 71.5	79.4           mage cl           Flowers           56.5           64.6           86.7           78.1           81.7           82.8           80.2           82.3           86.0           52.0           84.2           63.2           82.1           74.0           84.9           39.9           83.8           46.4	82.6 Food 78.3 82.1 84.7 85.5 86.8 87.8 88.1 88.8 87.9 75.4 83.2 81.4 86.8 87.0 87.1 64.5 85.9 71.0	98.4 EuroSAT 90.7 94.4 97.6 97.6 98.0 98.4 98.0 98.4 98.0 98.3 98.4 98.6 97.5 96.0 98.4 95.6 98.3 78.6 97.9 89.7	57.7 <b>sks us</b> <b>sun</b> 51.2 56.5 60.1 59.5 62.8 63.8 63.8 63.4 64.4 50.0 57.6 50.6 60.6 60.6 60.4 60.9 39.4 61.4 40.6	69.9           UCF           65.6           70.1           75.5           73.8           76.2           61.4           72.5           61.3           74.5           69.8           74.5           53.2           75.7           50.4	96.8 <b>CLIP V</b> <b>SVHN</b> 95.3 95.9 97.1 96.9 96.8 97.6 97.0 97.5 97.1 96.9 97.0 97.0 97.0 97.0 97.0 97.0 97.2 97.4 96.5 97.4 96.5 97.4 97.5 97.4 97.5 97.4 97.4 97.5 97.4 97.5 97.4 97.2 97.4 97.5 97.4 97.5 97.4 97.2 97.4 97.2 97.4 97.5 97.4 97.5 97.4 97.5 97.4 97.2 97.4 97.4 97.5 97.4 97.5 97.4 97.5 97.4 97.4 97.5 97.4 97.4 97.5 97.4 97.4 97.5 97.4 97.4 97.5 97.4 97.4 97.5 97.4 97.4 97.5 97.4 97.4 97.5 97.4 97.4 97.5 97.4 97.4 97.5 97.4 97.4 97.5 97.4 97.5 97.4 97.5 97.4 97.4 97.5 97.4 97.4 97.5 97.4 97.4 97.4 97.4 97.5 97.4 97.4 97.4 97.4 97.4 97.4 97.5 97.4	76.4 Pets 62.9 68.5 82.5 72.7 76.8 83.6 75.6 80.4 85.2 68.8 81.6 69.5 82.0 78.2 83.6 94.0 56.9 84.0 64.4 64.9 84.0 64.9 84.0 64.5 84.0 64.5 85.6 85.8 85.5 85.0 85.5 85.0 85.5 85.0 85.5 85.0 85.5 85.0 85.5 85.0 85.5 85.0 85.5 85.0 85.5 85.0 85.5 85.0 85.5 85.0 85.5 85.0 85.6	54.3         DTD         42.8         47.3         61.1         51.8         55.3         63.5         56.4         59.8         66.4         41.4         54.3         50.5         62.2         32.9         57.3         33.5	92.5 RESISC 82.1 86.9 92.8 91.2 92.9 94.7 93.7 94.6 95.0 79.1 91.9 87.7 94.3 93.7 94.3 93.7 94.4 67.6 93.3 75.2	Avg           70           82           7883           81           82           84           67           80           70           82           84           67           80           70           82           82           63
SparseGPT + IC Magnitude Magnitude Magnitude + IC Wanda Wanda + IC SparseGPT SparseGPT + IC Magnitude + IC Wanda Wanda + IC SparseGPT + IC SparseGPT SparseGPT + IC Magnitude Magnitude + IC Magnitude Magnitude + IC Magnitude + IC Magnitude + IC Magnitude + IC Wanda Wanda + IC	2:4 Table Sparsity 50% 47% 50% 50% 50% 50% 47% 47% 50% 47% 47% 27% 27% 27% 27% 27% 27% 27% 2	112M         2: Testing accl         #Nonzero Params         109M         111M	82.2 ITACY ON i CIFAR100 76.9 80.9 82.9 84.1 85.6 86.2 87.2 87.7 86.1 75.8 81.5 78.6 84.7 85.1 84.7 68.3 83.9 71.5 84.9	79.4           mage cl           Flowers           56.5           64.6           86.7           78.1           81.7           82.8           80.2           82.3           86.0           52.0           84.2           63.2           82.1           74.0           84.9           39.9           83.8           46.4           87.2	82.6 Food 78.3 82.1 84.7 85.5 86.8 87.8 88.1 88.8 87.9 75.4 83.2 81.4 86.8 87.0 87.1 64.5 85.9 71.0 86.9	98.4 EuroSAT 90.7 94.4 97.6 97.6 98.0 98.4 98.0 98.3 98.4 98.6 97.5 96.0 98.4 95.6 98.3 78.6 97.9 89.7 98.2	57.7 <b>sks us</b> <b>sun</b> 51.2 56.5 60.1 59.5 62.8 63.8 63.8 63.4 64.4 50.0 57.6 50.6 60.6 60.6 60.4 60.9 39.4 61.4 40.6 62.7	69.9           UCF           65.6           70.1           75.5           73.8           76.2           61.4           72.5           61.3           74.5           69.8           74.5           69.8           74.9           53.2           75.4           75.4	96.8 <b>CLIP V</b> <b>SVHN</b> 95.3 95.9 97.1 96.9 96.8 97.6 97.0 97.2 97.4 96.9 97.4 96.9 97.4 96.9 97.4 97.2 97.4 97.2 97.4 97.4 97.2 97.4 97.2 97.4 97.2 97.4 97.4 97.2 97.4 97.2 97.4 97.2 97.4 97.2 97.4 97.2 97.4 97.2 97.4 97.2 97.4 97.2 97.4 97.2 97.4 97.2 97.4 97.2 97.4 97.2 97.4 97.2 97.4 97.2 97.4 97.2 97.4 97.2 97.4 97.2 97.4 97.2 97.4 97.2 97.4 97.4 97.4 97.4 97.4 97.5 97.4	76.4 Pets 62.9 68.5 82.5 72.7 76.8 83.6 75.6 80.4 85.2 68.8 81.6 69.5 82.0 78.2 83.6 56.9 84.0 61.4 84.3	54.3           16.           DTD           42.8           47.3           61.1           51.8           55.3           63.5           56.4           59.8           66.4           41.4           54.3           50.5           62.2           32.9           57.3           33.5           59.2	92.5 RESISC 82.1 86.9 92.8 91.2 92.9 94.7 93.7 94.6 95.0 79.1 91.9 87.7 94.3 93.7 94.3 93.7 94.4 67.6 93.3 75.2 94.1	Avg           79.1           70.2           74.7           74.7           82.1           83.4           81.4           80.1     <
SparseGPT + IC Magnitude Magnitude Magnitude + IC Wanda Wanda + IC SparseGPT SparseGPT + IC Magnitude + IC Wanda Wanda + IC SparseGPT + IC Magnitude Magnitude + IC Magnitude Magnitude + IC Magnitude Magnitude + IC SparseGPT + IC Magnitude Magnitude + IC SparseGPT + IC Magnitude + IC SparseGPT + IC	2:4 Table Sparsity 50% 47% 50% 50% 50% 50% 47% 50% 47% 50% 4:8 4:8 4:8 4:8 4:8 4:8 4:8 4:8	112M         2: Testing accl         #Nonzero Params         109M         111M         109M         111M	82.2 ITACY ON i CIFAR100 76.9 80.9 82.9 84.1 85.6 86.2 87.2 87.7 86.1 75.8 81.5 78.6 84.7 85.1 84.7 68.3 83.9 71.5 84.9 82.8	79.4           mage cl           Flowers           56.5           64.6           86.7           78.1           81.7           82.8           80.2           82.3           86.0           52.0           84.2           63.2           82.1           74.0           84.9           39.9           83.8           46.4           87.2           66.2	82.6 assifi Food 78.3 82.1 84.7 85.5 86.8 87.8 88.1 88.8 87.9 75.4 83.2 81.4 86.8 87.0 87.1 64.5 85.9 71.0 86.9 82.8	98.4 EuroSAT 90.7 94.4 97.6 97.6 98.0 98.4 98.0 98.3 98.4 98.0 98.3 98.4 98.0 98.3 98.4 98.0 98.3 98.4 95.6 97.5 96.0 98.4 95.6 98.3 78.6 97.9 89.7 98.2 94.4	57.7 <b>sks us</b> <b>sun</b> 51.2 56.5 60.1 59.5 62.8 63.8 63.8 63.4 64.4 50.0 57.6 50.6 60.6 60.6 60.4 60.9 39.4 61.4 40.6 62.7 53.6	69.9           UCF           65.6           70.1           75.5           73.8           76.2           76.2           61.4           72.5           61.3           74.5           69.8           74.9           53.2           75.4           63.0	96.8 <b>CLIP V</b> <b>SVHN</b> 95.3 95.9 97.1 96.9 96.8 97.6 97.0 97.2 97.4 96.5 97.4 96.5 97.4 97.2 97.4 97.2 97.4 96.5 97.4 97.2 97.4 97.4 97.2	76.4 Pets 62.9 68.5 82.5 72.7 76.8 83.6 75.6 80.4 85.2 68.8 81.6 69.5 82.0 78.2 83.6 56.9 84.0 64.4 85.2 78.2 83.6 78.2 83.6 78.2 83.6 78.2 83.6 79.7 76.8 81.6 75.6 81.6 75.6 81.6 75.6 81.6 75.6 81.6 75.6 81.6 75.6 81.6 75.6 81.6 75.6 81.6 75.6 81.6 75.6 81.6 75.6 81.6 75.6 81.6 75.6 81.6 75.6 81.6 75.6 81.6 75.6 81.6 75.6 75.6 81.6 75.6 81.6 75.6 75.6 75.6 75.6 75.6 75.6 75.6 75.6 75.6 81.6 75.6 75.6 81.6 75.6 81.6 75.6 75.6 81.6 75.6 75.6 81.6 75.6 81.6 75.6 81.6 75.6 81.6 75.6 81.6 75.6 81.6 75.6 81.6 75.6 81.6 75.6 81.6 75.6 81.6 75.6 81.6 75.6 81.6 75.6 81.6 75.6 75.6 81.6 75.6 75.6 75.6 81.6 75.7 75.7	54.3         DTD         42.8         47.3         61.1         51.8         55.3         63.5         56.4         59.8         66.4         41.4         54.3         50.5         62.2         32.9         57.3         33.5         59.2         43.1	92.5 RESISC 82.1 86.9 92.8 91.2 92.9 94.7 93.7 94.6 95.0 79.1 91.9 87.7 94.3 93.7 94.3 93.7 94.4 67.6 93.3 75.2 94.1 90.2	Avg           79.0           70.1           74.1           78.6           80.8           83.4           83.4           84.2           84.2           87.1           80.1           80.1           80.2           84.2           87.1           80.1           70.5           82.5           82.5           82.5           82.5           82.5           82.5           83.6           83.6           83.6           83.6           83.6           83.6           83.6           83.6           83.6           83.6           83.6           83.6

261 o: 262 20 263 p 264 sa 265 tv 266 n 267 n 268

274

of 0.9. The mini-batch size is 128. Following (Bahng et al., 2022),  $v_i$  is learnable padding pixels on all sides, where the padding size is set to 30. The rank r is chosen as 32 and a sensitivity analysis is provided in Section 6. We evaluate two types of sparsity: unstructured sparsity and structured m:n sparisty (Mishra et al., 2021) (4:8 and 2:4), i.e., at most m out of every n contiguous weights to be non-zero.

**Baselines.** The proposed IC can be integrated into any existing pruning methods. To verify its effectiveness, we consider three pruning methods: (i) Magnitude Pruning (Han et al., 2015) which discards weights based on their magnitudes; (ii) Wanda (Sun et al., 2024) designs a scoring metric as the weight magnitudes multiplied by the corresponding input activations on a per-output basis; (iii) SparseGPT (Frantar & Alistarh, 2023) which adjusts the unpruned weights by solving a layer-wise reconstruction problem using a second-order optimizer. SparseGPT is a weight compensation method, while Magnitude and Wanda design a scoring metric for pruning without updating weights. For all methods, the base models are fully finetuned on the training set of all tasks before pruning.

**Results.** Tables 1 and 2 show the testing accuracy on ten image classification tasks using CLIP ViT-B/32 and ViT-B/16, respectively. As can be seen, IC consistently

275 brings large improvements to existing pruning methods 276 in both unstructured (sparsity=50%) and structured (4:8 277 and 2:4) cases. Specifically, compared with Magnitude, 278 Magnitude + IC achieves improvements of 25% and 10% 279 on ViT-B/32 and ViT-B/16, respectively; Compared with 280 Wanda, Wanda + IC has improvements of about 2.5%; Com-281 pared with SparseGPT, SparseGPT + IC performs better by 282 an improvement of 2% on ViT-B/32. The large improve-283 ments contributed by IC verify that the learned compensa-284 tion pool is effective in constructing input compensation 285 for the pruned models. Moreover, SparseGPT + IC con-286 sistently performs the best, demonstrating that combining 287 both weight compensation and input compensation is more desirable. We can also observe that unstructured pruning 289 (sparsity=50%) achieves higher accuracy than structured 290 pruning (sparsity=4:8), which is aligned with findings in 291 previous works (Sun et al., 2024; Frantar & Alistarh, 2023; 292 Zhang et al., 2024).

293 In our IC framework, as the encoder is a submodule of 294 the pruned model, the additional parameters are only from 295 the very small compensation pool (only 2.3M). To verify 296 that IC's improvements over baselines are not due to these 297 additional parameters, we reduce the sparsity of baseline 298 methods to increase the number of nonzero parameters. As 299 shown in Tables 1 and 2, when using the same number of 300 non-zero parameters, combining Magnitude + IC, Wanda 301 + IC, and SparseGPT + IC consistently achieve higher av-302 erage accuracy than Magnitude, Wanda, and SparseGPT, 303 respectively. 304

#### 5.2. Experiments on Natural Language Processing

305

306

329

307 Models and Datasets. We evaluate IC on the LLaMA 308 model family, i.e., LLaMA-1 (7B) (Touvron et al., 2023a), 309 LLaMA-2 (7B) (Touvron et al., 2023b), and LLaMA-3.1 (8B) (Meta, 2024). Following (Sun et al., 2024; Fran-311 tar & Alistarh, 2023), 128 sequences sampled from the 312 first shard of the C4 dataset (Raffel et al., 2020) are used 313 as training data. We evaluate the pruned models on two 314 types of tasks: (i) language modeling task which evaluates 315 the perplexity on the held-out validation data of WikiText-316 2 (Merity et al., 2016); and (ii) seven zero-shot tasks in-317 clude BoolQ (Clark et al., 2019), RTE (Wang, 2018), Hel-318 laSwag (Zellers et al., 2019), WinoGrande (Sakaguchi et al., 319 2021), ARC-easy/challenging (Clark et al., 2018), and Open-320 bookQA (Mihaylov et al., 2018)) from the EleutherAI LM Harness package (Gao et al., 2024). As LLMs contain 322 billions of parameters, to make pruned models more com-323 pressed, we follow (Yin et al., 2024) and focus on the un-324 structured sparsity of 70% case. Implementation details are 325 in Appendix B.4.

Results on Language Modeling Task. Table 3 shows the WikiText validation perplexity. As can be seen, IC consis-

Table 3:	WikiText	validation	perplexity	of pruned	LLaMA
family	of models.				

	Sparsity	LLaMA-1 (7B)	LLaMA-2 (7B)	LLaMA-3.1 (8B)
Magnitude	70%	48432	52457	3483567
Magnitude + IC	70%	19678	8585	33194
Wanda	70%	85.0	74.4	99.7
Wanda + IC	70%	56.5	67.0	80.1
SparseGPT	70%	26.8	24.7	38.8
SparseGPT + IC	70%	17.7	18.3	27.5
Wanda + OWL	70%	24.6	30.9	70.7
Wanda + OWL + IC	70%	19.1	24.5	63.1

tently brings a significant improvement to existing pruning methods, verifying the effectiveness of compensating inputs for pruned LLMs. For example, SparseGPT + IC achieves a perplexity improvement of 6.0 over SparseGPT on all three LLaMA family of models, while Wanda + IC outperforms Wanda by a large margin of 7.0 on all three LLMs. Although Magnitude performs much worse, Magnitude + IC still effectively reduces the perplexity by over 60%.

**Results on Zero-shot Tasks.** Table 4 shows the testing accuracy of seven zero-shot tasks on the LLaMA family of models. As can be seen, IC consistently brings a noticeable improvement (averaged over all tasks) to all existing pruning methods. For example, Wanda + IC outperforms Wanda on LLaMA-3.1-8B, LLaMA-2-7B, and LLaMA-1-7B by margins of 1.09%, 1.73%, and 0.4%, respectively, indicating that the learned compensation pool can be effectively used to construct input compensation for pruned models without any weight update. Moreover, SparseGPT + IC consistently achieves the highest accuracy for all models, showing that learning  $\Delta_x$  and  $\Delta_w$  are complementary and thus can be combined together for boosting performance.

#### **5.3. Experiments on Image Generation**

Experimental Setting. We evaluate IC on Denoising Diffusion Probability Models (DDPM) (Ho et al., 2020). Following (Fang et al., 2023), the CIFAR-10 dataset (with the image size of  $32 \times 32$ ) (Krizhevsky & Hinton, 2009) and the off-the-shelf DDPM from (Ho et al., 2020) are used. K is initialized with zero and V is initialized randomly by a normal distribution with a standard deviation of 0.01, where the rank r is set to 128. We train K and V using the Adam optimizer (Kingma & Ba, 2015) with a learning rate of 0.002 over 100K steps. The mini-batch size is set to 128. The identity function is used as the encoder of IC to keep more original image information, which is crucial for image generation. Following (Fang et al., 2023), we focus on the sparsity of 30% case and compare IC with three pruning methods: Magnitude Pruning (Han et al., 2015), Taylor Pruning (Molchanov et al., 2022), and Diff-Pruning (Fang et al., 2023).

**Results.** Table 5 shows the Frechet Inception Distance (FID) (Heusel et al., 2017). As can be seen, IC consis-

**Enhancing Pruned Models by Input Compensation** 

	Table 4: Testing accuracy of zero-shot tasks using LLaMA family of models.									
		Sparsity	BoolQ	RTE	HellaSwag	WinoGrande	ARC-e	ARC-c	OBQA	Avg
-	Magnitude	70%	38.29	52.71	25.62	51.14	26.64	19.71	11.60	32.24
(7B	Magnitude + IC	70%	55.99	52.35	25.33	48.38	25.93	21.93	15.00	34.99
A-1	Wanda	70%	57.16	54.87	28.73	50.91	32.15	18.86	13.80	36.64
aM.	Wanda + IC	70%	59.60	53.07	28.80	52.01	34.55	18.86	12.40	37.04
LL	SparseGPT	70%	63.43	56.32	33.89	58.96	44.07	23.63	17.80	42.58
	SparseGPT + IC	70%	66.06	54.87	37.47	60.06	48.40	25.51	18.20	44.37
	Magnitude	70%	37.95	53.07	25.95	49.25	27.74	22.78	16.80	33.36
7B)	Magnitude + IC	70%	42.57	52.35	25.77	49.33	25.84	22.10	16.20	33.45
-7 (	Wanda	70%	46.09	52.71	27.86	51.14	30.05	18.09	11.80	33.96
MA	Wanda + IC	70%	58.01	52.71	27.91	50.28	29.80	19.54	11.60	35.69
La	SparseGPT	70%	65.75	53.07	33.47	57.06	43.73	22.35	17.40	41.83
	SparseGPT + IC	70%	65.11	52.71	36.50	57.85	49.33	24.49	17.60	43.37
	Magnitude	70%	37.83	52.71	26.16	49.33	26.09	20.14	14.60	32.41
(8B)	Magnitude + IC	70%	37.83	53.79	25.71	49.88	25.21	22.78	15.20	32.91
3.1	Wanda	70%	56.27	52.71	27.51	47.83	32.20	17.66	13.00	35.31
-A1	Wanda + IC	70%	61.74	52.71	27.75	49.25	33.25	17.92	12.20	36.40
LaV	SparseGPT	70%	67.71	52.71	33.60	56.20	43.14	21.08	16.40	41.55
Ц	SparseGPT + IC	70%	67.71	54.15	34.25	57.62	46.63	22.78	15.60	42.68

Table 5: FID of pruned DDPMs on CIFAR-10.

	Sparsity	FID
Magnitude	30%	5.48
Magnitude + IC	30%	5.31
Taylor Pruning	30%	5.56
Taylor Pruning + IC	30%	5.21
Diff-Pruning	30%	5.29
Diff-Pruning + IC	30%	5.15

tently improves the existing pruning methods, demonstrating the effectiveness of compensating inputs for pruned LLMs. For instance, Taylor Pruning+IC achieves an FID improvement of 0.35 compared to Taylor Pruning. Similarly, Diff-Pruning+IC outperforms Diff-Pruning by 0.14.

#### 6. Analysis

330 331

333

335

345

353

354 355

361

362

363

364

367

368

369

In this section, we conduct empirical analyses to investigate the key components of IC, including rank r, sparsity, computation cost, sparse retraining, and input-dependent compensation. We adopt the experimental setting used in Section 5.1 with CLIP ViT-B/32.

375 Sensitivity of Rank. We conduct experiments to study the 376 sensitivity of rank r to the performance of Magnitude + IC, 377 where r is chosen from  $\{2, 4, 8, 16, 32, 64, 128\}$ . Figure 2 378 shows the testing accuracy with different ranks (detailed 379 results are shown in Table 8 of Appendix B.3). As can be 380 seen, a very small rank (e.g., 2) is not desirable. When the 381 rank is small ( $\leq$  16), increasing the rank leads to better 382 performance. A very large rank (e.g., 128) contains more 383 parameters but does not contribute to better performance. In 384



Figure 2: Performance of Magnitude + IC with different ranks on image classification tasks using CLIP ViT-B/32.

practice, we can choose the rank  $\in [16, 32]$ .

Sensitivity of Sparsity. We study the performance of Magnitude + IC with different sparsities. Figure 3 shows the trend of testing accuracy (averaged over ten tasks) w.r.t. sparsity (detailed results are shown in Table 6 of Appendix B.1). As can be seen, when the sparsity is high ( $\geq 40\%$ ), Magnitude + IC significantly outperforms Magnitude; When the sparsity is low ( $\leq 20\%$ ), Magnitude + IC and Magnitude perform comparably. In practice, a high sparsity is more desirable for pruning in order to reduce the model size; Thus, IC is practically useful for enhancing pruned models.

**Computation Cost.** Compared with existing pruning methods, the proposed IC needs two extra components (encoder and compensation pool) for adjusting the inputs. As the encoder is either a submodule of the pruned model or an identity function and the compensation pool is very small (only 2.3M in experiments), the extra storage is negligible and the computation cost is small. For example, in the



Figure 3: Performance of Magnitude and Magnitude + IC
 with different sparsities on image classification tasks using
 CLIP ViT-B/32.

397

398

399

400

401

402

403

404

405



406Figure 4: Testing accuracy (averaged over ten image classifi-<br/>cation tasks) using CLIP ViT-B/32 with sparse retraining

image classification experiment with CLIP ViT-B/32, compared with Magnitude, computation cost of Magnitude + IC increases by only 1% (from 305G to 309G FLOPs) and inference speed drops slightly from 665 to 535 images/second. Note that testing accuracy increases largely (Table 1), verifying that such additional computation cost is worthwhile.

Sparse Retraining with IC. In Section 5.1, we combine 415 416 IC with three pruning methods without sparse retraining (i.e., retraining the sparse model following the pruning step), 417 which can approach the performance of the dense model. 418 We conduct experiments to investigate whether IC is benefi-419 cial to pruning methods with sparse retraining. We retrain 420 retained parameters on the training data for 3 epochs using 421 the AdamW optimizer with a learning rate of 0.000001 and 422 weight decay of 0.01. Figure 4 shows the testing accuracy 423 (averaged over ten tasks) of pruning methods w/ or w/o IC 424 when sparse retraining is applied (detailed results are in 425 Table 7 of Appendix B.2). As can be seen, IC consistently 426 boosts existing pruning methods when sparse retraining is 427 applied. Moreover, sparse retraining achieves higher accu-428 429 racy than those without retraining (Table 1).

430 Input-dependent vs. Input-independent Compensation. 431 The design of our IC ensures the compensation  $\Delta_x$  is input-432 aware. A straightforward variant of IC is learning a globally 433 shared (i.e., input-independent) compensation  $\Delta$  for all 434 inputs. We conduct experiments to investigate the effective-435 ness of our input-dependent mechanism. Figure 5 shows the 436 testing accuracy (averaged over ten tasks). As can be seen, 437 IC performs much better than the input-independent vari-438 ant, demonstrating that input-dependent compensation is 439



(a) Sparsity=50%. (b) Sparsity=4:8. Figure 5: Testing accuracy (averaged over ten image classification tasks) of IC and an input-independent variant using CLIP ViT-B/32



Figure 6: Distribution of attention weights on image classification tasks using CLIP ViT-B/32.

more effective in reducing the error caused by the removed weights.

**Visualization.** In Section 5.1, we learn a compensation pool with r = 32 for ten image classification tasks, i.e.,  $\Delta_x$  is a weighted combination of  $32 v_i$ 's. Here, we study whether different tasks lead to different preferences for  $v_i$ 's. Figure 6 shows the average attention weights between  $v_i$  and testing samples belonging to different classes of three tasks (Flowers, Food, CIFAR100) (other tasks are not shown due to limited space). As can be seen, samples from Flowers prefer  $\{v_1, \ldots, v_5\}$ ; samples from Food prefer  $\{v_6, \ldots, v_{10}\}$ ; samples from CIFAR100 prefer  $\{v_{11}, \ldots, v_{17}\}$ .

### 7. Conclusion

In this paper, we proposed input compensation (IC) for enhancing pruned models by adjusting the inputs to compensate for the error caused by the pruned weights. A pool of multiple compensations is learned to construct inputdependent compensations. IC is designed in the input space while existing pruning methods are designed in the parameter space. Hence, IC can be integrated into any existing pruning methods. Extensive experiments on NLP and CV verify that IC largely enhances pruned models.

## 440 Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none of which we feel must be specifically highlighted here.

## References

441

442

443

444

445

446 447

448

449

450

451 452

453

454

455

456

457

458

459

460

461

- An, Y., Zhao, X., Yu, T., Tang, M., and Wang, J. Fluctuationbased adaptive structured pruning for large language models. In AAAI Conference on Artificial Intelligence, 2024.
- Bahng, H., Jahanian, A., Sankaranarayanan, S., and Isola, P. Exploring visual prompts for adapting large-scale models. Preprint arXiv:2203.17274, 2022.
- Blalock, D., Gonzalez Ortiz, J. J., Frankle, J., and Guttag,J. What is the state of neural network pruning? In *Proceedings of machine learning and systems*, 2020.
- Bossard, L., Guillaumin, M., and Van Gool, L. Food-101– mining discriminative components with random forests. In *European Conference on Computer Vision*, 2014.
- 463 Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., 464 Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., 465 Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., 466 Henighan, T., Child, R., Ramesh, A., Ziegler, D., Wu, J., 467 Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., 468 Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, 469 S., Radford, A., Sutskever, I., and Amodei, D. Language 470 models are few-shot learners. In Neural Information 471 Processing Systems, 2020. 472
- 473 Campos, J. and Lewis, F. Adaptive critic neural network for
  474 feedforward compensation. In *American Control Confer-*475 *ence*, 1999.
- Cheng, G., Han, J., and Lu, X. Remote sensing image scene classification: Benchmark and state of the art. In *Proceedings of the Institute of Electrical and Electronics Engineers*, 2017.
- 481 Cimpoi, M., Maji, S., Kokkinos, I., Mohamed, S., and
  482 Vedaldi, A. Describing textures in the wild. In *IEEE*483 *Conference on Computer Vision and Pattern Recognition*,
  484 2014.
- Clark, C., Lee, K., Chang, M.-W., Kwiatkowski, T., Collins,
  M., and Toutanova, K. BoolQ: Exploring the surprising difficulty of natural yes/no questions. Preprint arXiv:1905.10044, 2019.
- Clark, P., Cowhey, I., Etzioni, O., Khot, T., Sabharwal, A.,
  Schoenick, C., and Tafjord, O. Think you have solved question answering? try ARC, the AI2 reasoning challenge. Preprint arXiv:1803.05457, 2018.

- Das, R. J., Ma, L., and Shen, Z. Beyond size: How gradients shape pruning decisions in large language models. Preprint arXiv:2311.04902, 2023.
- Dettmers, T., Lewis, M., Belkada, Y., and Zettlemoyer, L. GPT3.Int8 (): 8-bit matrix multiplication for transformers at scale. In *Neural Information Processing Systems*, 2022.
- Ding, N., Hu, S., Zhao, W., Chen, Y., Liu, Z., Zheng, H.-T., and Sun, M. OpenPrompt: An open-source framework for prompt-learning. In *Annual Meeting of the Association* for Computational Linguistics, 2022.
- Ding, N., Lv, X., Wang, Q., Chen, Y., Zhou, B., Liu, Z., and Sun, M. Sparse low-rank adaptation of pre-trained language models. In *Conference on Empirical Methods in Natural Language Processing*, 2023.
- Dong, P., Li, L., Tang, Z., Liu, X., Pan, X., Wang, Q., and Chu, X. Pruner-Zero: Evolving symbolic pruning metric from scratch for large language models. In *International Conference on Machine Learning*, 2024.
- Fan, Y. and Hunter, A. Understanding the cooking process with english recipe text. In *Findings of the Association for Computational Linguistics*, 2023.
- Fang, G., Ma, X., and Wang, X. Structural pruning for diffusion models. In *Neural Information Processing Systems*, 2023.
- Franklin, G. F., Powell, J. D., Emami-Naeini, A., and Powell, J. D. *Feedback control of dynamic systems*. Prentice Hall Upper Saddle River, 2002.
- Frantar, E. and Alistarh, D. Optimal brain compression: A framework for accurate post-training quantization and pruning. In *Neural Information Processing Systems*, 2022.
- Frantar, E. and Alistarh, D. SparseGPT: Massive language models can be accurately pruned in one-shot. In *International Conference on Machine Learning*, 2023.
- Frantar, E., Kurtic, E., and Alistarh, D. M-FAC: Efficient matrix-free approximations of second-order information. In *Neural Information Processing Systems*, 2021.
- Gao, L., Tow, J., Abbasi, B., Biderman, S., Black, S., DiPofi, A., Foster, C., Golding, L., Hsu, J., Le Noac'h, A., Li, H., McDonell, K., Muennighoff, N., Ociepa, C., Phang, J., Reynolds, L., Schoelkopf, H., Skowron, A., Sutawika, L., Tang, E., Thite, A., Wang, B., Wang, K., and Zou, A. A framework for few-shot language model evaluation. Technical report, 2024.
- Guo, R., Xu, W., and Ritter, A. Meta-tuning LLMs to leverage lexical knowledge for generalizable language

- 495 496 497 498 499 500 501 502 503 504 505 506 507 508 509 510 511 512 513 514 515 516 517 518 519 520 521 522 523 524 525 526 527 528 529 530 531 532 533 534 535 536 537 538 539 540 541 542 543 544 545 546 547
  - style understanding. In Annual Meeting of the Association
     for Computational Linguistics, 2024.
  - Han, S., Pool, J., Tran, J., and Dally, W. Learning both
     weights and connections for efficient neural network. In
     *Nneural Information Processing Systems*, 2015.
  - Hassibi, B., Stork, D. G., and Wolff, G. J. Optimal brain surgeon and general network pruning. In *IEEE International Conference on Neural Networks*, 1993.
  - Helber, P., Bischke, B., Dengel, A., and Borth, D. EuroSAT:
    A novel dataset and deep learning benchmark for land
    use and land cover classification. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 2019.
  - Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., and
    Hochreiter, S. GANs trained by a two time-scale update
    rule converge to a local nash equilibrium. In *Neural Information Processing Systems*, 2017.
  - Ho, J., Jain, A., and Abbeel, P. Denoising diffusion probabilistic models. In *Neural Information Processing Systems*, 2020.
  - Hou, Y., Dong, H., Wang, X., Li, B., and Che, W.
     MetaPrompting: Learning to learn better prompts. In International Conference on Computational Linguistics, 2022.
  - Jawahar, C., Zisserman, A., Vedaldi, A., and Parkhi, O. Cats and dogs. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2012.
  - Jiang, W., Zhang, Y., and Kwok, J. Effective structured prompting by meta-learning and representative verbalizer. In *International Conference on Machine Learning*, 2023.
  - Kim, S., Hooper, C. R. C., Gholami, A., Dong, Z., Li, X., Shen, S., Mahoney, M. W., and Keutzer, K. SqueezeLLM: Dense-and-sparse quantization. In *International Conference on Machine Learning*, 2024.
  - Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015.
  - Kojima, T., Gu, S. S., Reid, M., Matsuo, Y., and Iwasawa,Y. Large language models are zero-shot reasoners. In *Neural Information Processing Systems*, 2022.
  - Krizhevsky, A. and Hinton, G. Learning multiple layers of features from tiny images. Technical report, 2009.
- Krstic, M. Input delay compensation for forward complete and strict-feedforward nonlinear systems. *IEEE Transactions on Automatic Control*, 2009.

- Kuo, B. C. and Golnaraghi, M. F. *Automatic control systems*. Prentice Hall Englewood Cliffs, NJ, 1995.
- Lester, B., Al-Rfou, R., and Constant, N. The power of scale for parameter-efficient prompt tuning. In *Empirical Methods in Natural Language Processing*, 2021.
- Li, J., Li, D., Xiong, C., and Hoi, S. BLIP: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International Conference on Machine Learning*, 2022.
- Li, J., Li, D., Savarese, S., and Hoi, S. BLIP-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International Conference on Machine Learning*, 2023a.
- Li, J., Tang, T., Zhao, W. X., Nie, J.-Y., and Wen, J.-R. Pretrained language models for text generation: A survey. *ACM Computing Surveys*, 2024.
- Li, X. L. and Liang, P. Prefix-Tuning: Optimizing continuous prompts for generation. In *Annual Meeting of the Association for Computational Linguistics*, 2021.
- Li, Y., Yu, Y., Zhang, Q., Liang, C., He, P., Chen, W., and Zhao, T. LoSparse: Structured compression of large language models based on low-rank and sparse approximation. In *International Conference on Machine Learning*, 2023b.
- Liang, C., Zuo, S., Zhang, Q., He, P., Chen, W., and Zhao, T. Less is more: Task-aware layer-wise distillation for language model compression. In *International Conference* on Machine Learning, 2023.
- Lin, J., Tang, J., Tang, H., Yang, S., Chen, W.-M., Wang, W.-C., Xiao, G., Dang, X., Gan, C., and Han, S. AWQ: Activation-aware weight quantization for on-device LLM compression and acceleration. In *Proceedings of Machine Learning and Systems*, 2024.
- Liu, J., Shen, D., Zhang, Y., Dolan, W. B., Carin, L., and Chen, W. What makes good in-context examples for GPT-3? In *Proceedings of Deep Learning Inside Out*, 2022.
- Liu, P., Yuan, W., Fu, J., Jiang, Z., Hayashi, H., and Neubig, G. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. ACM Computing Surveys, 2023.
- Liu, X., Zheng, Y., Du, Z., Ding, M., Qian, Y., Yang, Z., and Tang, J. GPT understands, too. Preprint arXiv:2103.10385, 2021.
- Loshchilov, I. and Hutter, F. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019.

- Merity, S., Xiong, C., Bradbury, J., and Socher, R. Pointer
  sentinel mixture models. Preprint arXiv:1609.07843,
  2016.
- 554 Meta. The LLaMA 3 herd of models. Preprint 555 arXiv:2407.21783, 2024.
- Mihaylov, T., Clark, P., Khot, T., and Sabharwal, A. Can a suit of armor conduct electricity? a new dataset for open book question answering. Preprint arXiv:1809.02789, 2018.
- Mishra, A., Latorre, J. A., Pool, J., Stosic, D., Stosic, D.,
  Venkatesh, G., Yu, C., and Micikevicius, P. Accelerating
  sparse deep neural networks. Preprint arXiv:2104.08378,
  2021.
- Molchanov, P., Tyree, S., Karras, T., Aila, T., and Kautz, J.
  Pruning convolutional neural networks for resource efficient inference. In *International Conference on Learning Representations*, 2022.
- Netzer, Y., Wang, T., Coates, A., Bissacco, A., Wu, B.,
  and Ng, A. Y. Reading digits in natural images with
  unsupervised feature learning. In *Neural Information Processing Systems Workshop*, 2011.
- Nilsback, M.-E. and Zisserman, A. Automated flower classification over a large number of classes. In *Indian Conference on Computer Vision, Graphics & Image Processing*, 2008.
- Podell, D., English, Z., Lacey, K., Blattmann, A., Dockhorn, T., Müller, J., Penna, J., and Rombach, R. SDXL: Improving latent diffusion models for high-resolution image synthesis. In *International Conference on Learning Representations*, 2024.
- Polino, A., Pascanu, R., and Alistarh, D. Model compression via distillation and quantization. In *International Conference on Learning Representations*, 2018.
- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., and
  Sutskever, I. Language models are unsupervised multitask
  learners. OpenAI Blog, 2019.
- Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G.,
  Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark,
  J., Krueger, G., and Sutskever, I. Learning transferable
  visual models from natural language supervision. In *International Conference on Machine Learning*, 2021.
- Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., and Liu, P. J. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 2020.

- Rombach, R., Blattmann, A., Lorenz, D., Esser, P., and Ommer, B. High-resolution image synthesis with latent diffusion models. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022.
- Sakaguchi, K., Bras, R. L., Bhagavatula, C., and Choi, Y. WinoGrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 2021.
- Shao, W., Chen, M., Zhang, Z., Xu, P., Zhao, L., Li, Z., Zhang, K., Gao, P., Qiao, Y., and Luo, P. OmniQuant: Omnidirectionally calibrated quantization for large language models. In *International Conference on Learning Representations*, 2024.
- Singh, S. P. and Alistarh, D. WoodFisher: Efficient secondorder approximation for neural network compression. In *Neural Information Processing Systems*, 2020.
- Soomro, K., Zamir, A. R., and Shah, M. UCF101: A dataset of 101 human actions classes from videos in the wild. Preprint arXiv:1212.0402, 2012.
- Sun, M., Liu, Z., Bair, A., and Kolter, J. Z. A simple and effective pruning approach for large language models. In *International Conference on Learning Representations*, 2024.
- Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.-A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., Rodriguez, A., Joulin, A., Grave, E., and Lample, G. LLAMA: Open and efficient foundation language models. Preprint arXiv:2302.13971, 2023a.
- Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlvkov, N., Batra, S., Bhargava, P., Bhosale, S., Bikel, D., Blecher, L., Ferrer, C. C., Chen, M., Cucurull, G., Esiobu, D., Fernandes, J., Fu, J., Fu, W., Fuller, B., Gao, C., Goswami, V., Goyal, N., Hartshorn, A., Hosseini, S., Hou, R., Inan, H., Kardas, M., Kerkez, V., Khabsa, M., Kloumann, I., Korenev, A., Koura, P. S., Lachaux, M.-A., Lavril, T., Lee, J., Liskovich, D., Lu, Y., Mao, Y., Martinet, X., Mihaylov, T., Mishra, P., Molybog, I., Nie, Y., Poulton, A., Reizenstein, J., Rungta, R., Saladi, K., Schelten, A., Silva, R., Smith, E. M., Subramanian, R., Tan, X. E., Tang, B., Taylor, R., Williams, A., Kuan, J. X., Xu, P., Yan, Z., Zarov, I., Zhang, Y., Fan, A., Kambadur, M., Narang, S., Rodriguez, A., Stojnic, R., Edunov, S., and Scialom, T. LLaMA 2: Open foundation and finetuned chat models. Preprint arXiv:2307.09288, 2023b.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I. Attention is all you need. In *Neural Information Processing Systems*, 2017.

- 605 Wang, A. GLUE: A multi-task benchmark and analysis 606 platform for natural language understanding. Preprint 607 arXiv:1804.07461, 2018. 608
- Wang, J., Bao, W., Sun, L., Zhu, X., Cao, B., and Philip, S. Y. 609 Private model compression via knowledge distillation. In 610 AAAI Conference on Artificial Intelligence, 2019.

611

621

622

623

624

625

630

631

632

633

- 612 Wang, W., Chen, W., Luo, Y., Long, Y., Lin, Z., Zhang, 613 L., Lin, B., Cai, D., and He, X. Model compression and 614 efficient inference for large language models: A survey. 615 Preprint arXiv:2402.09748, 2024.
- 616 Wei, J., Wang, X., Schuurmans, D., Bosma, M., brian ichter, 617 Xia, F., Chi, E. H., Le, Q. V., and Zhou, D. Chain of 618 thought prompting elicits reasoning in large language 619 models. In Neural Information Processing Systems, 2022. 620
  - Xiao, G., Lin, J., Seznec, M., Wu, H., Demouth, J., and Han, S. SmoothQuant: Accurate and efficient post-training quantization for large language models. In International Conference on Machine Learning, 2023.
- Xiao, J., Ehinger, K. A., Hays, J., Torralba, A., and Oliva, 626 627 A. SUN database: Exploring a large collection of scene 628 categories. International Journal of Computer Vision, 2016. 629
  - Xu, M., Yin, W., Cai, D., Yi, R., Xu, D., Wang, Q., Wu, B., Zhao, Y., Yang, C., Wang, S., et al. A survey of resource-efficient LLM and multimodal foundation models. Preprint arXiv:2401.08092, 2024.
- 635 Yao, Z., Yazdani Aminabadi, R., Zhang, M., Wu, X., Li, 636 C., and He, Y. ZeroQuant: Efficient and affordable post-637 training quantization for large-scale transformers. In 638 Neural Information Processing Systems, 2022.
- 639 Yin, L., Wu, Y., Zhang, Z., Hsieh, C.-Y., Wang, Y., Jia, Y., 640 Pechenizkiy, M., Liang, Y., Wang, Z., and Liu, S. Out-641 lier weighed layerwise sparsity (OWL): A missing secret 642 sauce for pruning LLMs to high sparsity. In International 643 Conference on Machine Learning, 2024. 644
- 645 Yu, L., Jiang, W., Shi, H., Yu, J., Liu, Z., Zhang, Y., Kwok, 646 J., Li, Z., Weller, A., and Liu, W. MetaMath: Bootstrap 647 your own mathematical questions for large language mod-648 els. In International Conference on Learning Representa-649 tions, 2024. 650
- Yu, X., Liu, T., Wang, X., and Tao, D. On compressing deep 651 models by low rank and sparse decomposition. In IEEE 652 Conference on Computer Vision and Pattern Recognition, 653 2017. 654
- 655 Zellers, R., Holtzman, A., Bisk, Y., Farhadi, A., and Choi, Y. 656 HellaSwag: Can a machine really finish your sentence? 657 In Annual Meeting of the Association for Computational 658 Linguistics, 2019. 659

- Zhang, H., Zhang, X., Huang, H., and Yu, L. Promptbased meta-learning for few-shot text classification. In Conference on Empirical Methods in Natural Language Processing, 2022a.
- Zhang, N., Li, L., Chen, X., Deng, S., Bi, Z., Tan, C., Huang, F., and Chen, H. Differentiable prompt makes pre-trained language models better few-shot learners. In International Conference on Learning Representations, 2022b.
- Zhang, T., Yu, T., Hashimoto, T., Lewis, M., Yih, W.-t., Fried, D., and Wang, S. Coder reviewer reranking for code generation. In International Conference on Machine Learning, 2023.
- Zhang, Y., Bai, H., Lin, H., Zhao, J., Hou, L., and Cannistraci, C. V. Plug-and-play: An efficient post-training pruning method for large language models. In International Conference on Learning Representations, 2024.

# Appendix

# 664 A. Additional Related Works on Foundation Models

Foundation Models are large pre-trained models designed to serve as base models for various downstream tasks. These models are typically trained on a large amount of data and contain massive of parameters. Notable examples include Large Language Models (LLMs) like LLaMA series (Touvron et al., 2023a;b; Meta, 2024), which have promising performance in natural language processing tasks such as text generation (Li et al., 2024; Zhang et al., 2023), understanding (Guo et al., 2024; Fan & Hunter, 2023), and reasoning (Wei et al., 2022; Yu et al., 2024). In the realm of computer vision (CV), models like CLIP (Contrastive Language-Image Pretraining) (Radford et al., 2021) use multimodal learning to bridge textual and visual information, enhancing various CV tasks such as image classification (Radford et al., 2021), image captioning and visual question answering (Li et al., 2022; 2023a). Additionally, diffusion models like DDPM (Ho et al., 2020), Stable Diffusion (Rombach et al., 2022), and SDXL (Podell et al., 2024) have revolutionized image generation by employing a process of gradually transforming noise into images, showing the diverse applications of foundation models in creative applications. 

# **B. Additional Experimental Results**

## B.1. Sensitivity of Sparsity

 Table 6 shows the testing accuracy of Magnitude and Magnitude + IC with different sparsities. We can see that Magnitude + IC significantly outperforms Magnitude when the sparsity is high ( $\geq 40\%$ ). In practice, a high sparsity is more desirable for pruning to reduce the model size. Hence, IC is practically useful for boosting the performance of pruned models.

687			Table 6: Pe	rformance	of Magr	nitude and M	lagnitud	e + IC v	vith differ	ent spai	rsities.		
688	Sparsity	IC	CIFAR100	Flowers	Food	EuroSAT	SUN	UCF	SVHN	Pets	DTD	RESISC	Avg
689 690	10%	X	88.2	97.7	89.1	98.8	73.6	86.0	97.1	91.6	74.6	96.0	89.3
691	10%	1	87.7	97.6	88.8	98.5	73.3	85.1	97.0	91.7	73.8	96.0	88.9
692	20%	X	87.5	96.5	88.4	98.6	72.9	84.1	96.9	90.4	72.8	95.4	88.3
693	20%	1	87.2	96.3	88.1	98.5	72.6	83.2	96.9	90.4	72.2	95.3	88.1
694 695	30%	X	84.8	88.8	84.2	97.6	68.6	79.2	96.5	88.3	67.7	94.1	85.0
696	30%	1	84.9	91.0	85.3	98.3	68.5	78.9	96.7	88.3	67.0	94.2	85.3
697	40%	X	71.1	57.0	70.0	93.5	54.6	65.8	93.8	73.2	49.7	87.7	71.6
698	40%	1	79.8	81.9	80.7	97.4	60.4	72.9	96.2	81.4	58.9	91.6	80.1
700	50%	X	33.9	26.1	34.2	45.6	30.8	35.4	45.3	38.7	27.9	55.4	37.3
701	50%	1	73.0	62.9	72.4	96.5	48.9	63.1	94.4	69.2	44.1	87.1	71.2
702	60%	X	5.6	5.2	4.3	4.3	5.5	4.7	9.1	9.0	10.0	10.2	6.8
703	60%	1	57.5	25.0	52.0	89.4	26.1	37.0	85.6	29.0	16.5	74.1	49.2
704	70%	X	1.7	2.7	2.0	13.0	0.8	2.1	7.5	2.9	2.8	3.3	3.9
706	70%	1	19.3	6.1	13.8	75.4	4.0	7.0	51.6	4.1	4.0	26.6	21.2
707	80%	X	1.0	1.0	0.8	13.0	0.3	1.4	6.5	2.8	1.7	2.4	3.1
708	80%	1	6.9	5.1	5.2	68.8	1.0	3.5	38.8	3.8	2.8	22.1	15.8
709	90%	X	1.1	0.5	1.0	6.9	0.2	0.7	6.4	2.6	2.1	2.1	2.4
711	90%	1	3.5	2.7	2.2	49.5	0.5	2.1	6.7	3.5	3.4	7.5	8.2

### **B.2. Sparse Retraining with IC**

We conduct experiments to study the performance of IC when sparse retraining is applied. Table 7 shows the testing accuracy on image classification tasks using CLIP ViT-B/32. As shown, IC brings a significant improvement to existing pruning methods when sparse retraining is used.

Table 7: Testing accuracy on image classification tasks using CLIP ViT-B/32 with sparse retraining.

	0		0			0			1		0	
	Sparsity	CIFAR100	Flowers	Food	EuroSAT	SUN	UCF	SVHN	Pets	DTD	RESISC	Avg
Magnitude	50%	79.3	80.8	81.2	87.8	61.6	73.8	95.5	78.2	55.1	90.6	78.4
Magnitude + IC	50%	82.4	82.6	83.0	98.2	63.1	76.3	96.8	80.0	57.7	92.5	81.3
SparseGPT	50%	80.3	85.4	83.4	83.1	63.1	74.7	96.1	82.0	58.2	91.3	79.8
SparseGPT + IC	50%	84.5	88.3	86.0	98.4	65.7	77.5	96.8	83.3	61.5	93.9	83.6
Wanda	50%	81.2	84.8	83.5	88.1	62.8	73.8	96.0	81.8	59.3	92.0	80.3
Wanda + IC	50%	84.6	87.4	85.3	98.4	64.7	76.6	96.8	82.6	61.5	94.0	83.2
Magnitude	4:8	75.9	70.4	78.7	77.0	57.1	68.1	94.9	76.1	47.9	88.9	73.5
Magnitude + IC	4:8	81.1	76.7	81.2	97.8	59.5	72.0	96.5	78.0	51.8	91.6	78.6
SparseGPT	4:8	79.6	80.3	82.7	80.9	60.0	70.0	95.7	81.2	55.6	90.8	77.7
SparseGPT + IC	4:8	83.8	84.0	84.8	98.3	62.2	74.3	96.7	82.7	58.7	93.7	81.9
Wanda	4:8	78.9	76.9	81.4	80.3	58.0	68.7	95.6	78.6	54.6	90.2	76.3
Wanda + IC	4:8	83.6	81.4	83.3	98.0	60.5	71.9	96.7	80.0	57.3	92.8	80.6

## **B.3. Sensitivity of Rank**

Table 8 shows the testing accuracy and number of parameters of  $(\mathbf{K}, \mathbf{V})$  with different ranks. We can see that a very small rank (e.g., 2) is not desirable. In practice, we can choose the rank from 16 to 32.

Table 8: Testing accuracy of Magnitude + IC (sparsity=50%) with different ranks on image classification tasks using CLIP ViT-B/32.

	1 2/01												
]	Rank	#Params	CIFAR100	Flowers	Food	EuroSAT	SUN	UCF	SVHN	Pets	DTD	RESISC	Avg
	2	0.14M	70.9	51.7	69.5	94.9	47.1	58.2	93.8	57.2	38.1	83.5	66.5
	4	0.29M	73.5	56.6	71.6	95.8	47.9	60.3	94.2	61.9	40.6	85.7	68.8
	8	0.58M	73.9	62.0	72.2	96.4	49.4	62.7	94.3	65.8	42.7	86.4	70.6
	16	1.15M	73.5	62.0	72.9	96.9	49.6	63.8	94.4	69.8	44.0	87.2	71.4
	32	2.30M	73.0	62.9	72.4	96.5	48.9	63.1	94.4	69.2	44.1	87.1	71.2
	64	4.60M	73.1	64.3	73.1	96.8	50.8	64.4	94.4	65.5	42.4	87.9	71.3
	128	9.20M	73.5	56.7	72.6	96.6	49.8	62.0	94.2	62.7	40.4	87.0	69.5

#### **B.4. Implementation Details of NLP Experiments**

**Implementation Details.** We randomly initialize  $\mathbf{K}$  and  $\mathbf{V}$  by a normal distribution with zero mean and standard deviation 0.01, where the rank r is set to 32. We train  $\mathbf{K}$  and  $\mathbf{V}$  using the AdamW optimizer (Loshchilov & Hutter, 2019) with a learning rate of 0.001 and a linear warmup scheduler over 20 epochs. The mini-batch size is set to 1, with a gradient accumulation of 2. The input embedding layer is used as the encoder of IC.

#### B.5. Language Modeling Task with Sparsity=50%

For language modeling experiments in Section 5.2, we focus on sparsity of 70% for making models more compressed. Here, we also follow Wanda (Sun et al., 2024) and conduct experiments to evaluate IC under the 50% sparsity. The table below shows the WikiText validation perplexity of pruned LLaMA family of models (results with <sup>†</sup> are reported in the original publication). As can be seen, IC consistently brings a significant improvement to existing pruning methods, verifying the effectiveness of compensating inputs for pruned LLMS again.

**Enhancing Pruned Models by Input Compensation** 

Table 9: Wiki	Table 9: WikiText validation perplexity of pruned LLaMA family of models.								
	Sparsity	LLaMA-1 (7B)	LLaMA-2 (7B)	LLaMA-3.1 (8B)					
Magnitude <sup>†</sup>	50%	17.29	14.89	-					
Magnitude	50%	17.29	14.90	134.26					
Magnitude + IC	50%	14.60	12.44	80.37					
Wanda <sup>†</sup>	50%	7.26	6.41	-					
Wanda	50%	7.26	6.42	9.88					
Wanda + IC	50%	7.17	6.33	8.80					
SparseGPT <sup>†</sup>	50%	7.22	6.51	-					
SparseGPT	50%	7.22	6.52	9.46					
SparseGPT + IC	50%	7.07	6.38	8.51					

## **B.6. Statistics of Image Classification Datasets**

Table 10 shows the statistics of the image classification datasets.

Dataset	Training Size	Testing Size	#Classe
CIFAR100 (Krizhevsky & Hinton, 2009)	50,000	10,000	100
Flowers (Nilsback & Zisserman, 2008)	4,093	2,463	10
Food (Bossard et al., 2014)	50,500	30,300	10
EuroSAT (Helber et al., 2019)	13,500	8,100	1
SUN (Xiao et al., 2016)	15,888	19,850	39
DTD (Cimpoi et al., 2014)	2,820	1,692	4
UCF (Soomro et al., 2012)	7,639	3,783	10
SVHN (Netzer et al., 2011)	73,257	26,032	1
Pets (Jawahar et al., 2012)	2,944	3,669	3
RESISC (Cheng et al., 2017)	18,900	6,300	4

Table 10: Summary of ten image classification datasets.