ROTPRUNER: LARGE LANGUAGE MODEL PRUNING IN ROTATED SPACE

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

025

026

034

Paper under double-blind review

ABSTRACT

Network pruning is a crucial technique for compressing large language models with billions of parameters, aiming to reduce memory and computational costs with minimal performance degradation. However, existing pruning methods for LLMs often focus on heuristic metrics or layer-wise reconstruction losses, neglecting the impact on the overall model output, which can lead to suboptimal result. Additionally, these methods operate directly on the original weight and activation spaces, which may not be ideal for pruning. In this paper, we propose that the original parameter space is not optimal for pruning and present a novel training-based pruning framework called RotPruner. RotPruner rotates the spaces of weight matrices and activations in linear layers, and applies existing pruning methods in a rotated space that is more suitable for pruning. We introduce an efficient algorithm to identify an appropriate rotation that preserves the performance of pruned LLMs. RotPruner is capable of integrating with other pruning methods and supporting unstructured, semi-structured, and structured pruning. We evaluate RotPruner on several large language models, including OPT, LLaMA-2, and LLaMA-3, and demonstrate state-of-the-art performance on both language modeling and zero-shot tasks.

028 1 INTRODUCTION

Recently, Large Language models (LLMs) have became a milestone in natural language processing, achieving great results in various tasks (Zhao et al., 2023). However, the success of these models results from an increase in scale and computational complexity, making the storage and time consuming of LLMs challenging. Model compression, as a post-training technique, has arouse great

Model compression techniques usually include three types: distillation, pruning and quantization (Zhu et al., 2023; Gholami et al., 2022; Hoefler et al., 2021). In this work, we focus on pruning, which sets the several elements in the weight matrices of model to zero. Traditional pruning techniques often requires a post-pruning re-training to recover the performance after pruning (Ma et al., 2023; Huang et al., 2020; Han et al., 2015). However, this is challenging in LLMs due to its model size. To address this limitation, post-training method without re-retaining, such as Wanda (Sun et al., 2023) and SparseGPT (Frantar & Alistarh, 2023) are proposed.

interest since it can reduce the memory and computational requirements of these models.

Current pruning methods face two major challenges. First, traditional pruning methods focus on heuristic metric or individual layer's reconstruction loss and ignore inter-layer interaction, leading a high accumulative error. In contrast, block-wise, or model-wise pruning, considering a block's or the whole model's reconstruction loss, can reduce the error accumulation. The larger the pruning group, the more difficult the optimization. Second, current pruning methods run the algorithm in the original weight space, which is not the optimal space to prune. Changing the pruning space will give a better result of the same pruning method.

To address the above challenges, we introduce RotPruner, a novel training-based pruning framework
 of LLMs. Figure 1 provides the overview of our method, which applies rotation matrix to activations
 and weights, and runs pruning method on the rotated activations and weights. This approach does
 not update the weights and thus can preserve the knowledge of the pretrained model. On the rotated
 weights and activation, our method can leverage any other one-shot pruning methods and can realize
 both structured and unstructured pruning. Our further exploration finds that if we run the pruning



Figure 1: Overview of RotPruner. RotPruner first applies rotation to activation X and weight W, and conduct pruning method on the rotated space. The rotation matrix is learned via cayley SGD.

method on random rotated space, the performance is catastrophic. Therefore, we propose to learn the
rotation matrices utilizing cayley SGD (Li et al., 2020), which is an efficient optimizer for training
orthonormal matrices, to optimize the rotation matrices to minimize the loss of the pruned network.
The weights of the model are fixed and thus this do not change the result of dense model.

We empirically evaluate our method on widely adopted open-source LLM: OPT (Zhang et al., 2022),
LLaMA-2 (Touvron et al., 2023) and LLaMA-3 (Dubey et al., 2024) families in the setting of unstructured, 2:4 semi-structured and structured pruning. Our approach exceeds the performance of
state-of-the-art methods such as SparseGPT (Frantar & Alistarh, 2023), Wanda (Sun et al., 2023)
and SliceGPT (Ashkboos et al., 2024a) across various language benchmarks. We show that the space
of the original weight matrices is not the optimal one to be pruned on and introduce an approach to
find a better pruning space.

088

090

092

074

075 076

2 RELATED WORK

091 2.1 NETWORK PRUNING

Network pruning is a widely used technique to reduce the model size and speed up the computation 093 of neural network by generating sparse weight matrices. Pruning can be categorized into unstruc-094 tured, semi-structured, and structured pruning based on different sparsity patterns. Unstructured 095 pruning eliminates the entries in the weight matrix without any structured pattern. While it can de-096 crease the entries of network, it can not get any inference speedup. Semi-structured pruning with 097 N:M sparsity (Zhou et al., 2021) requires N non-zero entries in every continuous M weights. It can 098 leverage NVIDIA's Sparse Tensor Cores to accelerate matrix computation. Structured pruning elim-099 inates the entries in entire rows or columns in the weight matrix and can reduce the dimension of 100 hidden state. It can reach significant computational and memory reducing with greater performance 101 loss.

102 103 104

2.2 ONE SHOT PRUNING AND TRAINING-BASED PRUNING

Traditional pruning requires re-training after pruning to recover the performance, which is challeng ing to Large Language Models due to its scale. One-shot pruning can prune LLMs in a single step
 without need for re-training, reducing the time and computational cost. For example, SparseGPT
 formalizes the problem of pruning LLMs by solving a local layer-wise reconstruction problem and

prunes the weight matrices based on the weight and inverse Hessian of the loss. Wanda prunes
 the weight matrices based on the product of weight magnitudes and norm of activation. SliceGPT
 (Ashkboos et al., 2024a) utilizes singular value decomposition to prune small singular vectors of
 activation and thus decrease the dimension of hidden state.

112 Training-based pruning includes mask update and weight update. ADMM-pruner (Boža) and 113 FISTAPruner (Zhao et al., 2024) convert the reconstruction error of a pruned model to a convex 114 problem and use classic convex optimization algorithm to solve the problem. They both update the 115 weight of the dense model, which may lose knowledge of the dense model. These methods use a 116 small size of calibration data, which is similar to one shot pruning. On the other hand, AST (Huang 117 et al., 2024) proposes a pruning framework to retrain pruned models efficiently, while using more 118 data than other methods. Our work is different from these methods. We do not change the weights of LLM to preserve the knowledge of the dense model and we use the same scale of data as one shot 119 pruning. 120

121 122

123 124

2.3 OUTLIERS IN LLM

Recent studies have found that LLM has a significant phenomenon of outliers (Puccetti et al., 2022; 125 Kovaleva et al., 2021; Timkey & Van Schijndel, 2021), whose magnitude is much more larger than 126 others'. The outliers occur in both weight and activation. Several works on LLM quantization 127 (Ashkboos et al., 2024b; Liu et al., 2024) have developed to efficiently quantize LLMs with little 128 performance loss. In the field of LLM pruning, OWL (Yin et al., 2023) addresses the emergent 129 outliers in LLMs and provides a new technique that leverages the distribution of outliers to guide 130 layer-wise sparsity assignment of LLM pruning. Our work further explore the application of outlier 131 distribution of LLMs. We develop a method to produce more outliers and larger variance in the 132 weight and activation of LLM and therefore improve the performance.

133 134 135

3 Method

136 137 138

In this section, we will introduce our method. First, we present our method motivated by outliers in LLMs. Next, we describe our method to combine orthonormal matrix with network pruning. Finally, we describe how to train the orthonormal matrix.

140 141 142

143

139

3.1 MOTIVATION

144 In network quantization, researchers have focused on the outliers of activation. Removal of outliers 145 improves the performance of quantized networks (Ashkboos et al., 2024b; Liu et al., 2024). In 146 the case of pruning, where weights are eliminated, we focus on weight outliers. It is intuitive that 147 preserving the outliers means small weights in matrix are eliminated, leading to small changes to the 148 weights matrix, and therefore can keep the performance of the network. Based on this motivation, pruning methods can obtain better results by defining better pruning metric to distinguish outliers in 149 LLMs. However, for a specific weight matrix, there must be a best sparsity mask that can achieve 150 the best result of pruning. We want to show that changing the distribution of weight matrix can get 151 a better result. 152

153 Consider a linear layer W and input X. The result of the linear layer is XW. In previous pruning 154 methods, sparsity mask M is obtained based on W and X, and the result of pruned linear layer is 155 $X(M \odot W)$. If we introduce a matrix A with full rank and apply it to X and W, $X' = XA^{-1}$, 156 W' = AW, the result of the dense layer is not changed. However, we can get a different pruning 157 result $X'(M' \odot W')$. Consider a simple example:

X

- 158
- 159

160

$$= \begin{bmatrix} 2 & 1 \end{bmatrix}, \boldsymbol{W} = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}$$



Figure 2: Distribution of weights in original space, learned rotated space and random rotated space.

Given the sparsity ratio 50%, it's obvious that pruning on W will always have a wrong result. However, if we apply a transformation,

$$\boldsymbol{A} = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{bmatrix} = \boldsymbol{A}^{-T}, \ \boldsymbol{X}' = \boldsymbol{X}\boldsymbol{A}^{-1} = \begin{bmatrix} \frac{3}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}, \ \boldsymbol{W}' = \boldsymbol{A}\boldsymbol{W} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

pruning on W' does not change the matrix, and thus does not change the result of matrix multiplication.

186 This simple motivated example implies that pruning on a weight matrix that have more zeros is 187 better. If we can find an A with full rank that minimize $\|AW\|_0$, we can get the best result of pruning. However this problem is hard to solved. Moreover, the optimal L^0 norm may not satisfy 188 the required sparsity ratio. This problem can be approximately converted to minimize $||AW||_1$ 189 (Candès et al., 2006), which can produce small elements in the weights. But eliminating small 190 elements of W' can harm the performance of the pruned model. The reason is that the distribution 191 of activation is also important to the result of pruned model. In LLMs, the activations are not 192 uniformly distributed, but also have emergent outliers (Ashkboos et al., 2024b). 193

Figure 2 shows the distribution of the weight matrices in a OPT-125M. We find that the distribution shows the occurrence of outliers. Random rotation will reduce the outliers and harm the performance of pruning (see Table 1). But with specially learned transform, the distribution shows more outliers and produce larger variance.

Model	OPT-125M	OPT-1.3B	OPT-2.7B	OPT-6.7B
original space	37.41	19.01	14.60	12.38
random rotated space	12783	17786	48.40	24.61

199 200 201

174

175 176

177

178 179

181 182 183

202 203

204

205

Table 1: WikiText2 perplexity (\downarrow) of runing Wanda in original space and rotated space.

206

3.2 ROTATION

207 208

Based on the analysis above, we introduce our method to combine orthonormal matrix with networkpruning.

We first restrict the transform matrix A to orthonormal matrix (denote as Q below). This can introduce two advantages. Firstly, it's easy to calculate the reversal of $Q(Q^{-1} = Q^T)$. Secondly, purely applying Q introduces additional parameters. Based on the computational invariant in Transformer (Vaswani, 2017) from SliceGPT (Ashkboos et al., 2024a), the orthonormal matrix can be fused into linear layers with a RMSNorm connected. The first step of our method is to convert LayerNorm (Ba, 2016) to RMSNorm (Zhang & Sennrich, 2019) (if the norm is LayerNorm) and fuse the scaling coefficient to weight matrices.



Figure 3: LLaMA block in rotated space. Input matrices W_{in} are pre-multiplied by Q^T and out matrices W_{out} are post-multiplied by Q. Layers filled with dotted line are pruned.

234 In each block of LLM decoder, we introduce the orthonormal matrix Q_s on the linear layers, as 235 illustrated in Figure 3. Take LLaMA family as an example. There are 7 linear layers in every 236 decoder layer. In each block (self attention or MLP), the weight matrices can be divided into two 237 groups: input matrix W_{in} and output matrix W_{out} . Linear layers in the group share the same Q 238 and layers with RMSNorm connected have orthonormal Q_s . The input matrices are pre-multiplied 239 by Q^T and the out matrices are post-multiplied by Q. Specifically, the embedding layer is postmultiplied by Q and the head projection is pre-multiplied by Q^T . The total number of Qs is two 240 times the number of layers plus one. Once we apply the orthonormal transform to the weight matrix, 241 the pruning method is conducted on them. 242

243 Compared with original network, we add an extra matrix multiplication on the residuals, which can 244 increase the size and computation. If the rotation matrices share the same parameters, for example 245 in Figure 3, $Q_1 = Q_2$, the residual is multiplied by an identity matrix, which is the same as the 246 original network. The more rotation matrices share parameters, the size and computation of the pruned network is smaller but the performance is worst. Formally, if there are l layers in the decoder 247 layer and every k Qs share parameters, there will be n = |2l/k| + 1 individual Qs in total. The 248 number of extra parameter is $(n-1) \cdot d^2$ and extra computation is $2(n-1) \cdot bsd^2$, where b, s, d249 represent the batch size, sequence length and dimension of hidden state. The ablation experiment on 250 shared rotation shows that although this will get a worst result, but can still outperform other method 251 and also save memory and computational cost. 252

3.3 **OPTIMIZING ROTATION MATRIX** 254

To train the orthonormal matrix, we define the optimization objective as the performance of the 256 pruned network. The optimization objective comprise the auto-regression training loss of the pruned network. We also add a distillation loss that measures the difference between dense model 258 and sparse model. We consider three types of distillation loss: L^2 -distance, cosine distance and 259 Jensen–Shannon divergence. The loss is finally defined as: 260

$$\arg\min \mathcal{L}(\boldsymbol{Q}_i; \boldsymbol{W}_i, \boldsymbol{M}_i, \boldsymbol{X}) = \mathcal{L}_{AR} + \alpha \mathcal{L}_{distill}$$

$$\int \mathcal{L}_{L^2}(Y, \hat{Y}) = \|Y - \hat{Y}\|_2$$

$$\mathcal{L}_{\text{distill}}(Y, \hat{Y}) = \begin{cases} \mathcal{L}(Y, \hat{Y}) = 1 - \frac{\hat{Y} \cdot \hat{Y}}{\|Y\|_2 \|\hat{Y}\|_2} \\ \mathcal{L}_{\text{IS}}(Y, \hat{Y}) = D_{\text{KL}}(\theta_Y, \theta_{\hat{Y}}) + D_{\text{KL}}(\theta_{\hat{Y}}, \theta_Y) \end{cases}$$

266 267 268

253

255

257

261 262

264 265

230

231

232 233

The pretrained weight W is fixed and we only train the Qs. To optimize the orthonormal matrix, 269 we use Cayley SGD method (Li et al., 2020), which can optimize on the Stiefel manifold efficiently.

5

270	Algorithm 1 RotPruner
2/1	Inputs: original model weights $\{W_i\}$ and input X, sparsity ratio s, one-shot pruning method \mathcal{M}
272	or fixed masks $\{M_i\}$, epochs
273	initialize $\{Q_i\}$
274	for epoch in epochs do
275	apply $\{Q_i\}$ on $\{W_i\}$
276	if given one-shot pruning method \mathcal{M} then
277	$M_i = \mathcal{M}(W'_i)$
278	end if
279	optimize $\{Q_i\}$ to minimize $\mathcal{L}(Q_i; W_i, M_i, X)$ using cayley SGD
280	end for

Specifically, in each iteration, the rotation matrices are updated by

$$\boldsymbol{Q}_{k+1} = (\boldsymbol{I} - \frac{\alpha}{2}\boldsymbol{Y})^{-1}(\boldsymbol{I} + \frac{\alpha}{2}\boldsymbol{Y})\boldsymbol{Q}_k$$

Y is a skew-symmetric matrix and is chosen to $Y = \hat{Y} - \hat{Y}^T$, where $\hat{Y} = GQ^T - \frac{1}{2}QQ^TGQ^T$, $G = \nabla f(Q)$. Moreover, the update matrix can be computed via fixed-point iteration to prevent matrix reversing. This optimizer keeps the Qs' orthonormality with approximately 2 times of the standard SGD.

When training the model with sparsity matrix M, there exists a problem that the gradient cannot be passed through the mask. Previous works use straight-through estimator (STE (Bengio et al., 2013)) to allow gradient to pass through mask by ignoring the mask in the backward pass. The backward of STE can be written as

$$\boldsymbol{W}_{t+1} = \boldsymbol{W}_t - \gamma_t(\boldsymbol{g}(\tilde{\boldsymbol{W}}_t))$$

Sparse-refined straight-through (SR-STE (Zhou et al., 2021)) estimator introduced a sparse-refined
 regularization term to the gradient, which can prevent mask oscillation. The backward of SR-STE can be written as

 $\boldsymbol{W}_{t+1} = \boldsymbol{W}_t - \gamma_t (g(\tilde{\boldsymbol{W}}_t) + \lambda_{\boldsymbol{W}}(\bar{\boldsymbol{M}} \odot \boldsymbol{W}_t))$

We adopt SR-STE to train the orthonormal matrix.

3.4 PROCEDURE

Algorithm 1 present the procedure of our method. Our method includes a basic pruning method without weight update and training of orthonormal matrix. For unstructured and semi-structured pruning, we choose two activation-free methods (magnitude based (Han et al., 2015) and Wanda) as the basic pruning method, since they need at most one forward pass and are efficient in time and memory. Activation-based methods need a backward pass and are time consuming. For structured pruning, we fixed the mask to the bottom rows or right columns.

In every training epoch, we apply the orthonormal matrix to the weight matrices and get sparsity masks on every layers based on the basic pruning method. Then we apply the masks and train the orthonormal matrices using STE with sparse-refined regularization term. After several iterations, update the masks based on new orthonormal matrix. For unstructured and semi-structured pruning, the orthonormal matrices are initialized to identity matrix; for structured pruning, they are initialized to the matrices that computed by SliceGPT via Principal Components Analysis.

316

281 282

287

288

289

290

295

299 300

301 302

303

317 4 EXPERIMENTS 318

Models and evaluation We evaluate our method on widely adopted LLMs: OPT-125M/1.3B/2.7B/6.7B, LLaMA-2-7B, LLaMA-3-8B. We evaluate our method on perplexity and zero-shot task. Perplexity is measured on test set of WikiText-2 (Merity et al., 2016). We use LM Evaluation Harness (Gao et al., 2021) to evaluate zero-shot accuracy on seven benchmarks: Wino-Grande (Sakaguchi et al., 2021), Piqa (Bisk et al., 2020), RTE (Wang, 2018), ARC Easy (Clark et al., 2018), ARC Challenge (Clark et al., 2018), WNLI (Wang, 2018) and QNLI (Wang, 2018).

WinoGrande, Piqa, ARC-easy and ARC-challenge benchmark the ability for knowledge question answering and RTE, QNLI and WNLI benchmark the ability for natural language inference.

Setup We use WikiText2 as the calibration set. We sample 128 data with sequence length 2048.
We evaluate three types of pruning: unstructured, 2:4 semi-structured and structured. We train for 5 epochs for each model and set the initial learning rate to 1e-2. The coefficient of SR-STE is 2e-5.
We are able to prune an 8B model in an L40S GPU with 48GB memory in 1.5 hour.

Baseline For unstructured and 2:4 semi-structured pruning, we choose SparseGPT and Wanda as baseline. For structured pruning, we compare against SliceGPT.

334 335 4.1 M

331 332

333

336

337

- 4.1 MAIN RESULTS
- 4.1.1 PERPLEXITY RESULTS

In Table 2, we show the WikiText2 perplexity result for different pruning of various models. We
achieved 50% unstructured, 2:4 semi-structured or 30% sturctured sparsity by pruning linear layers
in the LLMs, except for embeddings and head. These results show that RotPruner surpasses existing
methods on perplexity. We also find that the pruned OPT-6.7B can outperform the dense model.

We also run experiments OPT-1.3B with different sparsity ratio to further analyse the performance
 of RotPruner. The results are shown in figure Figure 4. We find that our method can surpass the
 dense model when sparsity ratio is smaller than 50%.

			OPT			LLaMA-2	LLaMA-3
Method	Sparsity	125M	1.3B	2.7B	6.7B	7B	8B
Dense	0%	27.65	14.62	12.47	10.85	5.47	6.13
SparseGPT	50% unstructured	34.12	17.48	13.43	11.61	6.46	8.29
Wanda	50% unstructured	37.41	19.01	14.60	12.38	6.72	9.40
RotPruner	50% unstructured	30.45	14.95	13.05	10.41	6.42	8.50
SparseGPT	2:4 semi-structured	60.02	23.83	17.20	14.13	10.37	14.65
Wanda	2:4 semi-structured	80.32	28.25	21.25	15.90	11.34	21.21
RotPruner	2:4 semi-structured	43.09	17.34	16.31	13.01	9.20	11.65
SliceGPT	30% structured 30% structured	42.32	20.26	16.30	12.80	8.62	17.08
RotPruner		32.56	16.19	14.06	12.09	8.34	15.26

Table 2: WikiText2 perplexity (\downarrow) performance comparison for different pruning methods on LLMs

4.1.2 ZERO-SHOT TASKS

In Table 3, we present the results of zero-shot tasks of different pruning method on OPT-6.7B, LLaMA-2-7B and LLaMA-3-8B. RotPruner surpasses other methods on OPT-6.7B and LLaMA-3-8B on the average accuracy of zero-shot tasks.

366 367

360 361

362

364

365

368

4.2 INFERENCE SPEED

We evaluate the inference speed of our pruned models on RTX4090. As discussed in subsection 3.2, since we add extra parameters and matrix multiplications on the residuals, we do this experiment to test how these parameters affect the inference speed. We only test on semi-structured pruning, because unstructured pruning has no benefit to the inference speed and we hold the same setting of structured pruning as SliceGPT.

We use Torch's to_sparse_semi_structured() to accelerate the 2:4 structured sparse models and Torch's Timer to benchmark the inference time. We test a LLaMA layer on eight situations: dense layer and sparse layer with residual rotation on attention or MLP. The result in show in Table 4. Tick means adding the residual rotation and cross means not. Specially, dense model without any residual rotations is the original dense model, sparse model without any residual rotations is the sparse

Model & sparsity	Method	WinoGrand	Piqa	RTE	ARC-e	ARC-c	WNLI	QNLI Mea
	Dense	65.51	76.27	50.50	65.66	30.46	46.48	50.87 55.8
OPT-6.7B 50%	SparseGPT	62.95	73.18	46.43	62.96	29.09	43.66	49.46 53.0
	Wanda	61.40	72.47	44.49	62.26	27.56	43.66	49.46 51.9
	RotPruner	62.72	74.01	45.42	64.20	28.84	45.07	49.46 53.3
	Dense	69.22	78.07	62.82	76.30	43.43	45.07	49.97 60.7
LLaMA2-7B 50%	SparseGPT	69.22	76.22	53.07	72.94	39.51	42.25	49.46 57.5
	Wanda	67.88	76.77	53.43	72.85	39.84	43.66	52.21 58.0
	RotPruner	65.27	73.99	54.15	68.85	37.03	43.66	55.45 56.9
LLaMA3-8B 50%	Dense	72.61	79.65	69.68	80.09	50.51	49.30	49.95 64.5
	SparseGPT	70.96	74.76	55.66	72.43	39.25	43.66	49.46 58.0
	Wanda	69.92	74.37	58.84	71.97	39.85	42.25	50.19 58.2
	RotPruner	67.40	75.35	63.53	71.76	38.99	40.85	50.41 58.

Table 3: Zero-shot tasks accuracy (\uparrow)

attention residual rotation	MLP residual rotation	Dense(ms)	Sparse(ms)
×	×	6.210	4.593 (1.352x)
\checkmark	×	6.638	4.619 (1.344x)
×	\checkmark	6.637	4.747 (1.308x)
\checkmark	\checkmark	7.045	4.837 (1.284x)

Table 4: Inference speed of sparse models

model obtained by traditional pruning method. Number in the round brackets means the speedup ratio compared with the dense model.

We find that the extra multiplications slow down the inference time but the gap is not significant especially with only attention residual rotation (1.006 times the time of no residual rotation). This is because other operations in attetion and MLP cost the majority of time and lower the affect of the residual rotation. Therefore, the residual rotation has little affect to the inference speed especially when we decrease the number of Q_s and add the residual rotation to the attention. If we use half a quarter of Q_s , which means every four Q_s share the same parameter, the inference time will be (1.003 times the time of no residual rotation). It's close to the sparse model using traditional pruning methods. Using less orthogonormal matrices will get closer to the traditional pruning methods. Actually, this operation not only has little affect to the performance of the sparse model, but also accelerate the convergence speed. We will provide experiment in the ablation experiments.

4.3 ABLATION EXPERIMENTS

We conduct ablation experiments to evaluate the performance of our method with different training losses, calibration set size, number of Q_s and basic pruning methods. The experiments are all conducted on a small model OPT-1.3B for short training time.

AR	L^2	cosine distance	JS-divergence	result
\checkmark				15.00
	\checkmark			20.01
		\checkmark		18.46
			\checkmark	18.80
\checkmark	\checkmark			14.85
\checkmark		\checkmark		14.33
\checkmark			\checkmark	15.03

Table 5: Ablation on training loss. Tick means using this loss.



Table 7: Ablation on Optimization method.

	number of Qs 49		25	13	7	4	
	perplexity 14.9	95	15.05	15.21	1 15.43	15.58	
	Table 8: A	Ablat	tion on 1	numbe	r of $oldsymbol{Q}$ s		
	method		OPT-1.	.3B	OPT-2.7B	LLaM/	A-2-7B
ma	magnitude agnitude+RotPruner		1712 17.1	2 5	265.14 21.08	19 18	.94 .84
Sparse	parseGPT w/o WR GPT w/o WR+RotPrun	er	23.0 15.2	1 0	23.27 14.26	7. 7.	89 32
	Wanda Wanda+RotPruner		19.0 14.8	1 6	14.60 13.05	6.' 6. '	72 42

Table 9: Ablation on different basic pruning method. WR stands for weight reconstruction.

Basic pruning method For pruning method, we choose magnitude, Wanda and SparseGPT. Magnitude and Wanda are fast, while SparseGPT takes a long time for pruning and it consists of pruning and weight construction. To shorten the training time, we use SparseGPT without weight reconstruction so that the sparse mask can be fixed. We do experiment on OPT-1.3B, OPT-2.7B and LLaMA-2-7B on the setting of 50% unstructured pruning. The experimental results Table 9 show that our method can improve the performance of one-shot pruning methods.

CONCLUSION

In this work, we introduce RotPruner, a novel training-based pruning framework for large language models. Unlike traditional pruning methods that operate directly in the original weight and input spaces, RotPruner employs a rotation that transforms the weight matrices and activations into an op-timal space for pruning. By doing this, RotPruner enables more effective pruning while minimizing the performance degradation typically associated with model compression. Our approach is com-patible with existing pruning methods, allowing for unstructured, semi-structured, and structured pruning.

We evaluate RotPruner on several large language models, including OPT, LLaMA-2, and LLaMA-3, achieving state-of-the-art results in both language modeling and various zero-shot tasks. These eval-uations demonstrate RotPruner's ability to enhance pruning performance. The results also highlight the significance of choosing an appropriate transformation space when applying pruning techniques. We hope this work can enhance the understanding of pruning in LLMs and the idea of rotated weight space can help improving the efficiency of neural networks.

REFERENCES

- Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat Behl, et al. Phi-3 technical report: A highly capable language model locally on your phone. arXiv preprint arXiv:2404.14219, 2024.
- Saleh Ashkboos, Maximilian L Croci, Marcelo Gennari do Nascimento, Torsten Hoefler, and James Hensman. Slicegpt: Compress large language models by deleting rows and columns. arXiv preprint arXiv:2401.15024, 2024a.
- Saleh Ashkboos, Amirkeivan Mohtashami, Maximilian L Croci, Bo Li, Martin Jaggi, Dan Alistarh, Torsten Hoefler, and James Hensman. Quarot: Outlier-free 4-bit inference in rotated llms. arXiv preprint arXiv:2404.00456, 2024b.

Jimmy Lei Ba. Layer normalization. arXiv preprint arXiv:1607.06450, 2016.

- 540 Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients 541 through stochastic neurons for conditional computation. arXiv preprint arXiv:1308.3432, 2013. 542 Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. Piqa: Reasoning about physical com-543 monsense in natural language. In Proceedings of the AAAI conference on artificial intelligence, 544 volume 34, pp. 7432–7439, 2020. 546 Vladimír Boža. Fast and effective weight update for pruned large language models. Transactions 547 on Machine Learning Research. 548 Emmanuel J Candès, Justin Romberg, and Terence Tao. Robust uncertainty principles: Exact signal 549 reconstruction from highly incomplete frequency information. IEEE Transactions on information 550 theory, 52(2):489-509, 2006. 551 552 Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and 553 Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. 554 arXiv preprint arXiv:1803.05457, 2018. 555 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha 556 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. arXiv preprint arXiv:2407.21783, 2024. 558 559 Elias Frantar and Dan Alistarh. Sparsegpt: Massive language models can be accurately pruned in one-shot. In International Conference on Machine Learning, pp. 10323–10337. PMLR, 2023. 561 Leo Gao, Jonathan Tow, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence 562 Golding, Jeffrey Hsu, Kyle McDonell, Niklas Muennighoff, et al. A framework for few-shot 563 language model evaluation. Version v0. 0.1. Sept, 10:8-9, 2021. 564 565 Amir Gholami, Sehoon Kim, Zhen Dong, Zhewei Yao, Michael W Mahoney, and Kurt Keutzer. A 566 survey of quantization methods for efficient neural network inference. In Low-Power Computer 567 Vision, pp. 291–326. Chapman and Hall/CRC, 2022. 568 Song Han, Jeff Pool, John Tran, and William Dally. Learning both weights and connections for 569 efficient neural network. Advances in neural information processing systems, 28, 2015. 570 571 Torsten Hoefler, Dan Alistarh, Tal Ben-Nun, Nikoli Dryden, and Alexandra Peste. Sparsity in deep 572 learning: Pruning and growth for efficient inference and training in neural networks. Journal of 573 Machine Learning Research, 22(241):1–124, 2021. 574 Weiyu Huang, Guohao Jian, Yuezhou Hu, Jun Zhu, and Jianfei Chen. Pruning large language models 575 with semi-structural adaptive sparse training. arXiv preprint arXiv:2407.20584, 2024. 576 577 Zhongzhan Huang, Wenqi Shao, Xinjiang Wang, Liang Lin, and Ping Luo. Convolution-578 weight-distribution assumption: Rethinking the criteria of channel pruning. arXiv preprint 579 arXiv:2004.11627, 2020. 580 Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, 581 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 582 Mistral 7b. arXiv preprint arXiv:2310.06825, 2023. 583 584 Olga Kovaleva, Saurabh Kulshreshtha, Anna Rogers, and Anna Rumshisky. Bert busters: Outlier 585 dimensions that disrupt transformers. arXiv preprint arXiv:2105.06990, 2021. 586 Jun Li, Li Fuxin, and Sinisa Todorovic. Efficient riemannian optimization on the stiefel manifold 587 via the cayley transform. arXiv preprint arXiv:2002.01113, 2020. 588 589 Zechun Liu, Changsheng Zhao, Igor Fedorov, Bilge Soran, Dhruv Choudhary, Raghuraman Krish-590 namoorthi, Vikas Chandra, Yuandong Tian, and Tijmen Blankevoort. Spinquant-Ilm quantization with learned rotations. arXiv preprint arXiv:2405.16406, 2024. 592
- 593 Xinyin Ma, Gongfan Fang, and Xinchao Wang. Llm-pruner: On the structural pruning of large language models. *Advances in neural information processing systems*, 36:21702–21720, 2023.

- 594 Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture 595 models. arXiv preprint arXiv:1609.07843, 2016. 596 Giovanni Puccetti, Anna Rogers, Aleksandr Drozd, and Felice Dell'Orletta. Outliers dimensions 597 that disrupt transformers are driven by frequency. arXiv preprint arXiv:2205.11380, 2022. 598 Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adver-600 sarial winograd schema challenge at scale. Communications of the ACM, 64(9):99-106, 2021. 601 602 Mingjie Sun, Zhuang Liu, Anna Bair, and J Zico Kolter. A simple and effective pruning approach for large language models. arXiv preprint arXiv:2306.11695, 2023. 603 604 William Timkey and Marten Van Schijndel. All bark and no bite: Rogue dimensions in transformer 605 language models obscure representational quality. arXiv preprint arXiv:2109.04404, 2021. 606 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-607 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda-608 tion and fine-tuned chat models. arXiv preprint arXiv:2307.09288, 2023. 609 610 A Vaswani. Attention is all you need. Advances in Neural Information Processing Systems, 2017. 611 Alex Wang. Glue: A multi-task benchmark and analysis platform for natural language understand-612 ing. arXiv preprint arXiv:1804.07461, 2018. 613 614 An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, 615 Chengyuan Li, Dayiheng Liu, Fei Huang, et al. Qwen2 technical report. arXiv preprint 616 arXiv:2407.10671, 2024. 617 Lu Yin, You Wu, Zhenyu Zhang, Cheng-Yu Hsieh, Yaqing Wang, Yiling Jia, Mykola Pechenizkiy, 618 Yi Liang, Zhangyang Wang, and Shiwei Liu. Outlier weighed layerwise sparsity (owl): A missing 619 secret sauce for pruning llms to high sparsity. arXiv preprint arXiv:2310.05175, 2023. 620 621 Biao Zhang and Rico Sennrich. Root mean square layer normalization. Advances in Neural Infor-622 mation Processing Systems, 32, 2019. 623 Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christo-624 pher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer 625 language models. arXiv preprint arXiv:2205.01068, 2022. 626 627 Pengxiang Zhao, Hanyu Hu, Ping Li, Yi Zheng, Zhefeng Wang, and Xiaoming Yuan. A convexoptimization-based layer-wise post-training pruner for large language models. arXiv preprint 628 arXiv:2408.03728, 2024. 629 630 Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, 631 Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. arXiv 632 preprint arXiv:2303.18223, 2023. 633 Aojun Zhou, Yukun Ma, Junnan Zhu, Jianbo Liu, Zhijie Zhang, Kun Yuan, Wenxiu Sun, and Hong-634 sheng Li. Learning n: m fine-grained structured sparse neural networks from scratch. arXiv 635 preprint arXiv:2102.04010, 2021. 636 637 Xunyu Zhu, Jian Li, Yong Liu, Can Ma, and Weiping Wang. A survey on model compression for 638 large language models. arXiv preprint arXiv:2308.07633, 2023. 639 640 641 А APPENDIX 642 643 A.1 EXPERIMENT ON MORE LLMS 644 We also add other LLMs: Phi-3-3.8B (Abdin et al., 2024), Mistral-7B (Jiang et al., 2023) and 645 Qwen2-7B (Yang et al., 2024) for comparison. The results are shown in Table 10. Our method 646
- 647 perform better in most of the models and pruning settings.

653					
654					
655					
656					
657					
658					
659					
660					
661					
662					
663					
664					
665					
666					
667					
668					
669	Method	Sparsity	Phi-3-3.8B	Mistral-7B	Qwen2-7B
670	Dense	0%	6.01	5.25	7.15
671	Crasses CDT	5001	0.56	5.00	7.02
672	SparseGP1 Wondo	50% unstructured	8.30	5.99	7.92
673	PotPruper	50% unstructured	9.37	0.20 5 03	8.40 8.12
674	Kotr fuller	J0% unstructured	0.45	5.75	0.12
675	SparseGPT	2:4 semi-structured	11.99	8.16	9.30
676	Wanda	2:4 semi-structured	20.00	10.70	12.19
677	RotPruner	2:4 semi-structured	12.53	7.31	9.25
678	SliceGPT	30% structured	10.65	8.94	10.73
679	RotPruner	30% structured	10.16	7.42	9.64
680					

Table 10: WikiText2 perplexity (\$\$\$\$\$) performance comparison for different pruning methods on LLMs