

# MOONSHOT: A Framework for Multi-Objective Pruning of Vision and Large Language Models

Anonymous authors  
Paper under double-blind review

## Abstract

Weight pruning is a common technique for compressing large neural networks. We focus on the challenging post-training one-shot setting, where a pre-trained model is compressed without any retraining. Existing one-shot pruning methods typically optimize a single objective, such as a layer-wise reconstruction loss or a second-order Taylor approximation of the training loss. We highlight that neither objective alone is consistently the most effective across architectures and sparsity levels. Motivated by this insight, we propose *MOONSHOT*, a general and flexible framework that extends any single-objective pruning method into a multi-objective formulation by jointly optimizing both the layer-wise reconstruction error and second-order Taylor approximation of the training loss. *MOONSHOT* acts as a wrapper around existing pruning algorithms. To enable this integration while maintaining scalability to billion-parameter models, we propose modeling decisions and introduce an efficient procedure for computing the inverse Hessian, preserving the efficiency of state-of-the-art one-shot pruners. When combined with state-of-the-art pruning methods on Llama-3.2 and Llama-2 models, *MOONSHOT* reduces C4 perplexity by up to 32.6% at 2:4 sparsity and improves zero-shot mean accuracy across seven classification benchmarks by up to 4.9 points. On Vision Transformers, it improves accuracy on ImageNet-1k by over 5 points at 70% sparsity, and on ResNet-50, it yields a 4-point gain at 90% sparsity. Our code will be made publicly available in the camera-ready version if the paper is accepted.

## 1 Introduction

Contemporary vision and language models have huge parameter counts (He et al., 2016; Dosovitskiy et al., 2021; Zhang et al., 2022), incurring significant computational costs during the inference phase. Pruning is a common strategy for compressing large neural networks. The aim is to remove a subset of weights by setting them to zero while maintaining relatively high predictive performance. Pruning can be either a) unstructured, where any individual weight can be set to zero (Han et al., 2015; Benbaki et al., 2023; Frantar et al., 2022; Kuznedelev et al., 2023; Frantar & Alistarh, 2023; Sun et al., 2024), b) structured, where entire rows and columns are set to zero (Ma et al., 2023; Meng et al., 2024b) or c) semi-structured where specific patterns are enforced, such as n:m sparsity, where n weights are set to zero within each block of m weights (Frantar et al., 2022; Kuznedelev et al., 2023; Frantar & Alistarh, 2023; Sun et al., 2024). In this work, we focus on these three different compression modes.

Various techniques have been proposed for pruning vision and large language models (Han et al., 2015; Frankle & Carbin, 2019; Yu et al., 2022; Frantar et al., 2022; Benbaki et al., 2023; Kuznedelev et al., 2023; Frantar & Alistarh, 2023; Meng et al., 2024a; Sun et al., 2024). Many existing methods rely on gradual pruning, where the model is fine-tuned on the original loss after every pruning stage to recover accuracy. However, for billion-scale models, such fine-tuning can be extremely expensive. In this context, recent works (Frantar & Alistarh, 2022; Frantar et al., 2022; Benbaki et al., 2023; Kuznedelev et al., 2023) have focused on the challenging task of *post-training* pruning in one-shot i.e., compressing a model without retraining based on a small amount of calibration data. In this paper, we focus on post-training one-shot pruning approaches, which are computationally attractive and particularly relevant for real-world applications.

When pruning a pre-determined fraction of the weights, various criteria are employed to preserve model accuracy or perplexity as much as possible, each leading to a different performance-sparsity trade-off. For example, weight magnitudes can be used as a criterion to decide which weights to prune and which to keep (Hanson & Pratt, 1988; Mozer & Smolensky, 1989; Gordon et al., 2020). However, magnitude-based pruning approaches rely extensively on expensive retraining to minimize the loss in performance. Another popular type of approach uses a local quadratic approximation of the original training loss to estimate the reduction in model performance. These approaches then approximately minimize this objective while imposing a sparsity constraint. This idea was introduced by LeCun et al. (1989b); Hassibi & Stork (1992b) through the Optimal Brain Surgeon (OBS) framework and built upon by various methods (Singh & Alistarh, 2020b; Frantar et al., 2021; Yu et al., 2022; Benbaki et al., 2023; Kuznedelev et al., 2023). A third prevalent criterion is based on the layer-wise OBS strategy (Dong et al., 2017; Frantar et al., 2022; Frantar & Alistarh, 2023; Sun et al., 2024; Meng et al., 2024b). In this approach, the pruning task is divided into layer-wise subproblems. For each layer, the goal is to minimize the squared reconstruction error between the original and pruned layer outputs subject to a sparsity constraint. While the OBS objective uses global information from the training loss of the pre-trained neural network to guide pruning, the layer-wise reconstruction loss uses more localized information in the embedding spaces.

To better understand the impact of pruning criteria on performance, we conducted a series of experiments across both vision and language models. On Vision Transformers, we evaluated CAP (Kuznedelev et al., 2023), which minimizes a second-order Taylor approximation of the training loss. On a convolutional neural network (ResNet-50 (He et al., 2015)), we considered OBC (Frantar et al., 2022), which uses the layer-wise reconstruction loss. For large language models, we evaluated SparseGPT (Frantar & Alistarh, 2023) and Wanda (Sun et al., 2024), both of which are designed around the layer-wise reconstruction objective. [These methods have support for both unstructured and semi-structured pruning.](#) To isolate the effect of the pruning criterion, we adapted each method to operate under the opposite objective: we evaluated CAP using the layer-wise reconstruction error, and OBC, CAP and Wanda using the second-order Taylor approximation of the training-loss. These comparisons, illustrated in Table 1, revealed that neither criterion is uniformly superior. Depending on the architecture, pruning method, and sparsity level, either the layer-wise reconstruction error or second-order Taylor approximation of the training-loss objective may yield better results. In several cases, the pruning methods performed better when paired with the objective they were not originally designed to optimize.

This indicates that *the two objectives capture complementary signals of parameter importance*, and relying on one alone can lead to suboptimal pruning decisions. Motivated by this insight, we propose a novel multi-objective optimization framework that jointly minimizes both the layer-wise reconstruction objective and second-order approximation of the training loss. This multi-objective optimization consistently improves the performance-sparsity trade-off of various state-of-the-art methods.

Extending state-of-the-art pruning methods to a multi-objective formulation introduces new challenges. These pruning algorithms typically approximate the objective using a quadratic form involving the Hessian of the model weights and require computing (or approximating) its inverse (Singh & Alistarh, 2020a; Frantar et al., 2022; Frantar & Alistarh, 2023;

Table 1: Comparison between the second-order Taylor approximation of the training loss and the layer-wise reconstruction error objectives across different pruning methods, models and sparsity regimes. We either keep the original objective, indicated with a star \*, as pruning criterion (approximation of the training loss for CAP and layer-wise reconstruction loss for OBC, Wanda and SparseGPT), or we replace it with the alternative single-objective criterion.

Domain	Model	Method	Sparsity	Second-Order Taylor Approx. of Training Loss	Layer-Wise Reconst. Error	
Language models	Llama-3.2-1B	SparseGPT	0.50	29.14±0.16	<b>27.15*</b> ±0.23	
		Wanda	0.50	<b>30.48</b> ±0.14	35.71*±0.21	
		Wanda	0.60	<b>88.87</b> ±1.62	117.71*±0.87	
	C4 perplexity (↓)	Llama-3.2-3B	SparseGPT	0.50	18.12±0.06	<b>17.61*</b> ±0.08
			Wanda	0.50	<b>18.2</b> ±0.04	18.88*±0.03
			Wanda	0.60	<b>40.28</b> ±0.67	41.98*±0.4
Vision models	DeiT-Tiny	CAP	0.60	<b>62.28*</b> ±0.05	54.18±0.15	
			2:4	<b>52.28*</b> ±0.04	47.65±0.11	
	ImageNet-1k accuracy (↑)	DeiT-Small	CAP	0.50	<b>77.27*</b> ±0.03	76.56±0.04
				2:4	69.65*±0.02	<b>70.25</b> ±0.04
				0.50	50.88±25.39	<b>76.63*</b> ±0.05
ResNet-50	OBC	0.70	48.94±24.42	<b>74.73*</b> ±0.03		

Kuznedev et al., 2023; Sun et al., 2024; Meng et al., 2024b). To make this calculation efficient, these methods rely on approximations or exploit the Hessian structure, typically using a block-diagonal approximation. However, the Hessians associated with the layer-wise reconstruction loss and the second-order Taylor approximation of the training loss exhibit different structures, and existing algorithms adopt distinct block-diagonal formulations depending on the specific objective. As a result, combining the Hessians from different objectives impacts directly the block-diagonal approximation, introducing new challenges in adapting the single-objective pruning methods. In *MOONSHOT*, we propose some modeling decisions (as described in Figure 2) to adapt the existing algorithms to the new multi-objective formulation.

While a relatively straightforward adaptation is possible for smaller architectures, additional computational complexity appears in large-scale models. SparseGPT (Frantar & Alistarh, 2023) and OSSCAR (Meng et al., 2024b), for instance, are state-of-the-art pruning methods for one-shot pruning of LLMs, which minimize the layer-wise reconstruction error. SparseGPT has support for both unstructured and semi-structured pruning while OSSCAR is designed for structured pruning. Both methods achieve their efficiency by exploiting the structure of the layer-wise reconstruction objective: the Hessian for each layer is naturally block-diagonal, with each block corresponding to the Hessian of a row in the weight matrix. Notably, all the blocks are identical and depend only on the input data. This structure allows for a very efficient computation of both the Hessian and its inverse. In contrast, *MOONSHOT* combines the reconstruction loss with a second-order approximation of the training loss, resulting in a Hessian that is no longer block-diagonal and particularly large, especially for models with billions of parameters. Even imposing a block-diagonal approximation on the Hessian of the multi-objective formulation is insufficient: since the blocks along the diagonal differ, inverting the Hessian requires computing a lot more matrix inversions, and a naive computation quickly becomes prohibitively expensive in practice. To address these challenges, we develop an efficient method to scale the multi-objective formulation to modern large architectures. In particular, we propose a fast approximate method for computing the Hessian and its inverse, enabling compatibility with high-performance state-of-the-art pruning methods such as SparseGPT and OSSCAR.

Our framework is very flexible and can handle different pruning patterns. *MOONSHOT* can be used for any of the sparsity patterns supported by the underlying single-objective baseline. In this work, we consider unstructured, semi-structured 2:4, and structured sparsity. Unstructured pruning offers strong memory savings and can yield speedups on CPUs (NeuralMagic, 2021) and specialized hardware accelerators (Han et al., 2015; Dave et al., 2021), but typically requires high sparsity ratios to achieve speedups on GPUs (Gale et al., 2020), often at the cost of model performance. In contrast, n:m sparsity also enables efficient execution on modern GPUs even at moderate sparsity levels (Mishra et al., 2021), making it particularly well-suited for the sparsity regimes commonly used in large language models (Frantar & Alistarh, 2023). Finally, structured pruning yields direct speedups on GPUs and CPUs (Kurtic et al., 2023; Meng et al., 2024b), but is typically applied at much lower sparsity ratios, since maintaining accuracy becomes increasingly difficult at higher structured sparsity levels.

*MOONSHOT* is orthogonal to the other existing methods designed to improve single-objective pruning. In particular, prior works (Frantar et al., 2022; Kuznedev et al., 2023; Lu et al., 2024; Yin et al., 2024) have shown that, in the case of unstructured pruning, a more principled distribution of the sparsity budget across layers can improve the performance of single-objective pruning approaches. In vision models, CAP (Kuznedev et al., 2023) improves top-1 accuracy of DeiT-Tiny on ImageNet-1k by nearly 10 points (a relative gain of 22%) at 70% sparsity with non-uniform sparsity allocation. These improvements are even more critical for large language models, which typically suffer severe performance degradation when pruned beyond 50% sparsity unless the sparsity is distributed non-uniformly across layers. For example, Yin et al. (2024) report that, with OWL, a carefully optimized layer-wise sparsity allocation, it is possible to achieve 71.38% smaller perplexity on WikiText using Wanda (Sun et al., 2024) at 70% sparsity. We show that when our method, *MOONSHOT*, is combined with non-uniform sparsity allocation strategies, we achieve additional improvements in the performance of the pruned model. On Llama-3.2-1B and Llama-3.2-3B across both SparseGPT and Wanda pruning baselines, *MOONSHOT* reduces C4 perplexity by up to an additional 25% and improves zero-shot mean accuracy by up to 1 additional point in comparison to the baselines with non-uniform sparsity allocation alone.

**Contributions.** We propose a novel optimization-based framework, which extends existing single-objective pruning approaches to a multi-objective formulation, enabling improved accuracy-sparsity trade-offs in the post-training one-shot pruning setting. Our contributions can be summarized as given below.

- We show that the layer-wise reconstruction loss and second-order Taylor approximation of the training loss result in different sparsity-accuracy trade-offs across architectures and sparsity levels. To the best of our knowledge, this is the first work to highlight those differences and systematically compare these two pruning objectives side by side.
- Motivated by this insight, we introduce a novel multi-objective optimization formulation to simultaneously minimize two objectives: a local quadratic approximation of the training loss and the layer-wise reconstruction error, subject to sparsity constraints. While these objectives have been considered in isolation, considering them simultaneously is new. Our framework, *MOONSHOT* (Multi-Objective **ONE-SHOT** pruning), provides a principled extension to existing single-objective pruning approaches, enabling them to operate under a multi-objective formulation.
- We introduce a set of modeling choices and algorithmic adaptations that extend single-objective pruning methods to a multi-objective setting. For applications to large language models, we propose an efficient procedure for computing the inverse Hessian in our multi-objective formulation. This fast computation is essential for preserving the scalability of existing pruning methods. Our implementation of *MOONSHOT*-SparseGPT prunes Llama-3.2-3B in under 40 minutes and Llama-3.2-1B in 8 minutes on a single GPU, showing that the multi-objective formulation remains efficient at the relevant billion-parameter scale.
- We validate our proposed method across diverse domains and applications:
  - (i) We evaluate *MOONSHOT* on large language models, including Llama-3.2-1B, Llama-3.2-3B (Grattafiori et al., 2024) and Llama-2-13b-chat-hf (Touvron et al., 2023). *MOONSHOT* improves the performance of state-of-the-art pruning methods such as SparseGPT (Frantar & Alistarh, 2023) and Wanda (Sun et al., 2024) in the post-training one-shot setting, under (a) unstructured sparsity (including non-uniform allocations via OWL (Yin et al., 2024) and AlphaPruning (Lu et al., 2024)), (b) semi-structured  $n:m$  sparsity. It also improves OSCCAR (Meng et al., 2024b) in the structured pruning case. On Llama-3.2-1B and 3B, *MOONSHOT* reduces C4 test perplexity (Raffel et al., 2020) of Wanda by up to 32.6% at 2:4 sparsity, achieves over 20% perplexity reduction for both SparseGPT and Wanda at 60% and 2:4 sparsity, and improves the mean accuracy across seven classification benchmarks by up to 1.5 points (see Figure 1 and Table 3). At 10% structured sparsity, *MOONSHOT* reduces C4 perplexity by up to 11% and improves the mean accuracy by up to 4.9 points. For Llama-2-13b-chat-hf, *MOONSHOT* reduces C4 perplexity by up to 14% and similarly improves the mean accuracy by up to 1.5 points at 70% unstructured sparsity (see Table 3).
  - (ii) We also evaluate *MOONSHOT* on computer vision benchmarks, including Vision Transformers (DeiT-Tiny, DeiT-Small, and DeiT-Base) and a convolutional model (ResNet-50). Across these models, our approach improves state-of-the-art methods such as CAP (Kuznedelev et al., 2023) and OBC (Frantar et al., 2022) in the unstructured and  $n:m$  sparsity regimes. In particular, on ImageNet-1k (Deng et al., 2009), it improves CAP by over 5 points in accuracy at 70% sparsity and 2 points at 2:4 sparsity, and improves OBC by 4 points at 90% sparsity (see Figure 1 and Table 2).

## 2 Multi-Objective Pruning

Our framework aims to improve existing single-objective network pruning methods by considering both the layer-wise reconstruction error and second-order approximation of the training loss. We first introduce these single-objective loss functions. We then formulate a new multi-objective optimization problem and propose *MOONSHOT*, which adapts state-of-the-art algorithms to minimize this new objective. To maintain the efficiency of the original LLM pruning methods, we finally propose a highly efficient approach for computing the inverse Hessian, a key component for algorithms like SparseGPT (Frantar & Alistarh, 2023) and OSSCAR (Meng et al., 2024b).

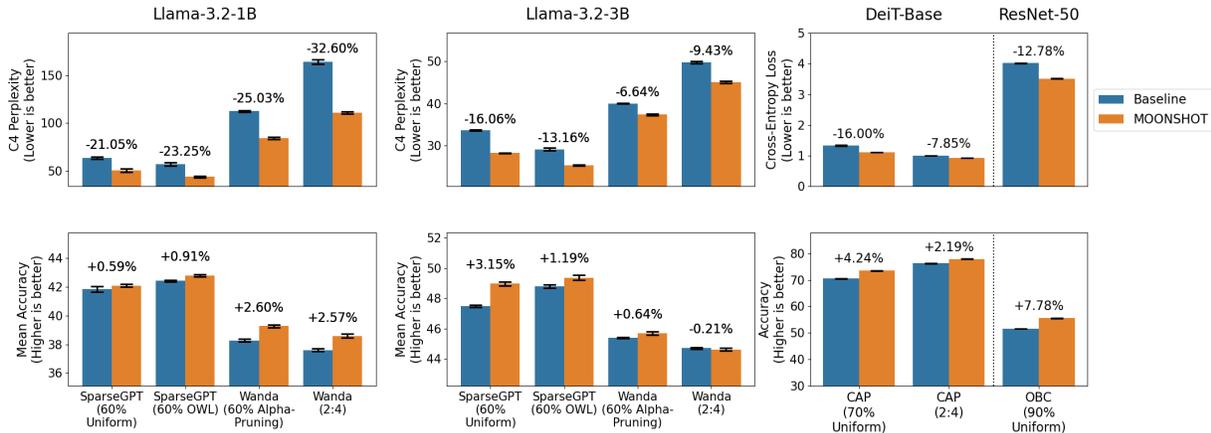


Figure 1: Impact of *MOONSHOT* on SparseGPT/Wanda (Llama-3.2) and CAP/OBC (DeiT-Base, ResNet-50) across sparsity regimes. For vision models, mean cross-entropy and ImageNet-1k accuracy are reported; for LLMs, perplexity on C4 along with mean zero-shot accuracy over seven classification tasks. Results are averaged over three seeds with standard errors.

## 2.1 Layer-wise pruning objectives

Consider the task of pruning a neural network with  $L$  layers. For any given layer  $l \in [L]$ , layer-wise pruning approaches (Nagel et al., 2020; Frantar et al., 2022; Kuznedelev et al., 2023; Frantar & Alistarh, 2023; Sun et al., 2024) aim to zero out some parameters and potentially adjust the remaining weights in the  $l$ -th layer to minimize the performance drop as much as possible. More formally, given the original pre-trained weights in the  $l$ -th layer, layer-wise pruning targets the following discrete optimization problem:

$$\min_{W^{(l)}} \mathcal{L}(W^{(l)}, \widehat{W}^{(l)}) \quad \text{s.t.} \quad \mathcal{S}(W^{(l)}) \leq S^{(l)}, \quad (1)$$

where  $\mathcal{S}(W^{(l)}) \leq S^{(l)}$  denotes the sparsity constraint, which depends on the sparsity type (unstructured, structured or n:m) and budget, and  $\mathcal{L}(W^{(l)}, \widehat{W}^{(l)})$  the loss function, which measures the performance drop when  $\widehat{W}^{(l)}$  in the  $l$ -th layer is replaced by  $W^{(l)}$ . Usually, the loss function uses a set of  $N$  training samples  $\{X_i\}_{i=1}^N$ . As we describe below, there are two commonly used loss functions  $\mathcal{L}$  in the one-shot pruning literature.

**Layer-wise reconstruction error.** Various existing layer-wise compression frameworks (Dong et al., 2017; He et al., 2017; Hubara et al., 2021; Frantar et al., 2022; Frantar & Alistarh, 2023; Sun et al., 2024) evaluate the pruned network’s performance by examining changes in the output of the pruned layer. Their goal is to minimize the squared error loss between the layer’s outputs generated by  $W^{(l)}$  and  $\widehat{W}^{(l)}$  on a training set. For a linear layer<sup>1</sup>  $l$  with input dimension  $d_{\text{in}}^{(l)}$  and output dimension  $d_{\text{out}}^{(l)}$ , we represent its input over  $N$  training samples as a  $d_{\text{in}}^{(l)} \times N$  matrix  $X^{(l)}$ . This reconstruction loss ( $\mathcal{L}_R^{(l)}$ ) can be written as follows:

$$\mathcal{L}_R^{(l)}(W^{(l)}) := \left\| W^{(l)} X^{(l)} - \widehat{W}^{(l)} X^{(l)} \right\|_F^2. \quad (2)$$

By ensuring that the outputs of the pruned layers remain close to those of the corresponding dense layers, this objective preserves the functional behavior of each layer, thereby maintaining the overall integrity of the model. As the entire network is constructed through the composition of its individual layers, the layer-wise reconstruction error can be viewed as a local approximation of the training loss.

**Second-order Taylor approximation and Fisher loss.** Another line of work (Hassibi & Stork, 1992a; Singh & Alistarh, 2020b; Benbaki et al., 2023) considers the impact of pruning weights on the (global)

<sup>1</sup>Convolutional layers can be processed in similar ways.

training loss. They consider a second-order Taylor approximation of the training loss  $\mathcal{L}_{\text{Tr}}$  around the pre-trained weights  $\widehat{W}$  using a second-order Taylor expansion. Typically, we set  $\nabla \mathcal{L}_{\text{Tr}}(\widehat{W}) = 0$  as  $\widehat{W}$  is assumed to be a stationary point of the training loss.

Since computing the full Hessian is expensive, earlier work (Hassibi & Stork, 1992a) uses an approximation based on the empirical Fisher information matrix:  $\nabla^2 \mathcal{L}_{\text{Tr}}(\widehat{W}) \approx H = \frac{1}{N} \sum_{i=1}^N \nabla \ell_i(\widehat{W}) \nabla \ell_i(\widehat{W})^\top$ , where  $\nabla \ell_i(\widehat{W})$  denotes the gradient of the network for weights  $\widehat{W}$  on the  $i$ -th training sample. In the case of layer-wise pruning, we prune a layer  $l$  while keeping the weights for all the other layers fixed, and the second-order Taylor approximation leads to the following Fisher loss:

$$\mathcal{L}_F^{(l)}(W^{(l)}) := \left( \text{vec}(W^{(l)}) - \text{vec}(\widehat{W}^{(l)}) \right)^\top H^{(l)} \left( \text{vec}(W^{(l)}) - \text{vec}(\widehat{W}^{(l)}) \right), \quad (3)$$

where  $\text{vec}(W^{(l)})$  and  $\text{vec}(\widehat{W}^{(l)})$  indicate the vector form of the pruned and dense weights in layer  $l$  respectively, and  $H^{(l)}$  denotes the submatrix of the approximated Hessian corresponding to the weights in layer  $l$ .

The Fisher loss  $\mathcal{L}_F^{(l)}(W^{(l)})$  provides a more global view of the network’s behavior (since this is based on computing the gradient of the entire pre-trained model) as opposed to the layer-wise reconstruction error which focuses on the outputs of individual layers.

**Pruning objective proposed in MOONSHOT.** Leveraging the merits of pruning based on the reconstruction of the layer outputs and the second-order Taylor approximation of the training loss, our framework introduces the task of pruning layer  $l$  as a multi-objective optimization problem, defined as follows:

$$\min_{W^{(l)}} \left( \mathcal{L}_R^{(l)}(W^{(l)}), \mathcal{L}_F^{(l)}(W^{(l)}) \right) \quad \text{s.t.} \quad \mathcal{S}(W^{(l)}) \leq S^{(l)} \quad (4)$$

This approach offers two benefits: (i) Targeting multiple objectives enhances the accuracy of the pruned networks beyond what is achievable with a single-objective (e.g. layer-wise reconstruction error or Fisher loss). (ii) By considering multiple objectives simultaneously and therefore leveraging more information, pruning becomes more robust, maintaining high performance even when one of the single objectives,  $\mathcal{L}_R^{(l)}(W^{(l)})$  or  $\mathcal{L}_F^{(l)}(W^{(l)})$ , fails to accurately capture the network’s overall performance.

## 2.2 Reformulation as cardinality constrained convex quadratic problem

We consider a weighted combination of the two individual objectives to address the multi-objective pruning formulation in equation 4. To ensure a balanced consideration of  $\mathcal{L}_R^{(l)}(W^{(l)})$  and  $\mathcal{L}_F^{(l)}(W^{(l)})$ , which might differ widely in magnitude, we normalize these objectives relative to their values at  $\mathbf{0}$ , the weight matrix filled with zeros. For  $\lambda \in [0, 1]$ , we set the objective as:

$$\mathcal{L}_\lambda^{(l)} := (\lambda / \mathcal{L}_R^{(l)}(\mathbf{0})) \mathcal{L}_R^{(l)} + ((1 - \lambda) / \mathcal{L}_F^{(l)}(\mathbf{0})) \mathcal{L}_F^{(l)} \quad (5)$$

In the following, we show how existing baselines can be adapted to the multi-objective formulation: we first explain how the block-diagonal structure of the Hessian arises in these methods, then how it can be preserved in the multi-objective case, and finally show how to reduce the multi-objective formulation to a quadratic problem under sparsity constraints, which can be addressed by existing pruning algorithms.

**Block-diagonal representation.** The layer-wise reconstruction loss can be rewritten (Frantar et al., 2022):

$$\mathcal{L}_R^{(l)}(W^{(l)}) = \sum_{i=1}^{d_{out}} \left\| W_{i,:}^{(l)\top} X^{(l)} - \widehat{W}_{i,:}^{(l)\top} X^{(l)} \right\|_2^2 \quad (6)$$

with  $W_{i,:}^{(l)}$  and  $\widehat{W}_{i,:}^{(l)}$  the  $i$ -th row of  $W^{(l)}$  and  $\widehat{W}^{(l)}$  respectively.

This allows us to express the layer-wise reconstruction loss in the following quadratic form:

$$\mathcal{L}_R^{(l)}(W^{(l)}) = \left( \text{vec}(W^{(l)}) - \text{vec}(\widehat{W}^{(l)}) \right)^\top H_R^{(l)} \left( \text{vec}(W^{(l)}) - \text{vec}(\widehat{W}^{(l)}) \right) \quad (7)$$

with  $H_R^{(l)} = \text{Diag}(X^{(l)}(X^{(l)})^T, \dots, X^{(l)}(X^{(l)})^T)$ , a block-diagonal matrix containing  $X^{(l)}(X^{(l)})^T$   $d_{out}$  times.

This exact block-diagonal structure enables Hessian computations to scale to large architectures, where a dense Hessian would be intractable. It also makes the objective separable, which can be exploited to improve the efficiency of pruning algorithms (Frantar et al., 2022; Frantar & Alistarh, 2023).

Similarly, for computational reasons, the existing single-objective pruning algorithms using the Fisher loss also assume  $H^{(l)}$  in equation 3 to be block-diagonal (Singh & Alistarh, 2020b; Benbaki et al., 2023; Kuznedelev et al., 2023). Specifically, we can write  $H^{(l)}$  in Fisher loss as  $\text{Diag}(H_1^{(l)}, H_2^{(l)}, \dots, H_K^{(l)})$  with  $K$  the number of blocks. Using the quadratic expressions from equation 3 and equation 7, the weighted loss from equation 5 becomes:

$$\begin{aligned} \mathcal{L}_\lambda^{(l)}(W^{(l)}) &= \left( \text{vec}(W^{(l)}) - \text{vec}(\widehat{W}^{(l)}) \right)^\top \left( \frac{\lambda}{\mathcal{L}_R^{(l)}(\mathbf{0})} H_R^{(l)} + \frac{1-\lambda}{\mathcal{L}_F^{(l)}(\mathbf{0})} H^{(l)} \right) \left( \text{vec}(W^{(l)}) - \text{vec}(\widehat{W}^{(l)}) \right) \quad (8) \\ &= \frac{\lambda}{\mathcal{L}_R^{(l)}(\mathbf{0})} \sum_{i=1}^{d_{out}} \left( W_{i,:}^{(l)} - \widehat{W}_{i,:}^{(l)} \right)^\top X^{(l)}(X^{(l)})^T \left( W_{i,:}^{(l)} - \widehat{W}_{i,:}^{(l)} \right) + \frac{1-\lambda}{\mathcal{L}_F^{(l)}(\mathbf{0})} \sum_{k=1}^K (w_k^{(l)} - \widehat{w}_k^{(l)})^\top H_k^{(l)} (w_k^{(l)} - \widehat{w}_k^{(l)}) \quad (9) \end{aligned}$$

where  $w_k^{(l)}$  and  $\widehat{w}_k^{(l)}$  denote the weights of  $W^{(l)}$  and  $\widehat{W}^{(l)}$  corresponding to block  $k$  in  $H^{(l)}$  respectively.

**Adapting existing baselines.** If  $H_R^{(l)}$  and  $H^{(l)}$  use different block sizes, the original block structure would be altered. To isolate the impact of the multi-objective formulation while preserving the effectiveness of the baseline, *MOONSHOT* enforces the block size of the original baseline. As described in Figure 2, there are two main cases:

- For algorithms focusing only on the Fisher loss, such as CAP (Kuznedelev et al., 2023), we can apply a block-diagonal approximation to  $H_R$ , specifically to  $X^{(l)}(X^{(l)})^T$ , ensuring that each block of  $H_R^{(l)}$  aligns in size with the corresponding block of  $H^{(l)}$ . In this case, we maintain the original block-diagonal approximation of  $H^{(l)}$  assumed by the single-objective baseline.
- For algorithms that focus solely on the layer-wise reconstruction loss, such as OBC (Frantar et al., 2022), SparseGPT (Frantar & Alistarh, 2023) and OSSCAR (Meng et al., 2024b), we can set  $K$  to  $d_{out}$ , ensuring that each block of  $H^{(l)}$  matches the dimensions of  $X^{(l)}(X^{(l)})^T$ . In this case, we maintain the exact block-diagonal structure of  $H_R^{(l)}$  without further approximation, as in the original baseline.

Finally, for Wanda (Sun et al., 2024) which uses a diagonal approximation of  $X^{(l)}(X^{(l)})^T$ , we use a diagonal approximation of  $H^{(l)}$  as well.

In all cases, we write in the following  $H_R^{(l)} = \text{Diag}(L_1^{(l)}, \dots, L_K^{(l)})$  ( $L_k^{(l)}$  denotes either a block of the block-diagonal approximation of  $X^{(l)}(X^{(l)})^T$  or  $X^{(l)}(X^{(l)})^T$  itself, depending on the baseline).

**Quadratic formulation.** Let  $F_k^{(l)} = \frac{\lambda}{\mathcal{L}_R^{(l)}(\mathbf{0})} L_k^{(l)} + \frac{1-\lambda}{\mathcal{L}_F^{(l)}(\mathbf{0})} H_k^{(l)}$ . The multi-objective formulation in equation 4 can be reformulated as the following quadratic optimization problem under sparsity constraints:

$$\min_{W^{(l)}} \mathcal{L}_\lambda^{(l)}(W^{(l)}) = \sum_{k=1}^K (w_k^{(l)} - \widehat{w}_k^{(l)})^\top F_k^{(l)} (w_k^{(l)} - \widehat{w}_k^{(l)}) \quad \text{s.t.} \quad \mathcal{S}(W^{(l)}) \leq S^{(l)}, \quad (10)$$

Most single-objective pruning methods reduce to solving a separable quadratic problem with sparsity constraints (Frantar & Alistarh, 2023; Kuznedelev et al., 2023; Frantar et al., 2022; Meng et al., 2024b). At this point, the formulation resembles the single-objective case (equation 3 and equation 7), which means that existing single-objective baselines can, in principle, be applied to the new multi-objective setting. However, as we show in the following section, a direct application in the case of LLMs is intractable and requires additional adaptations.

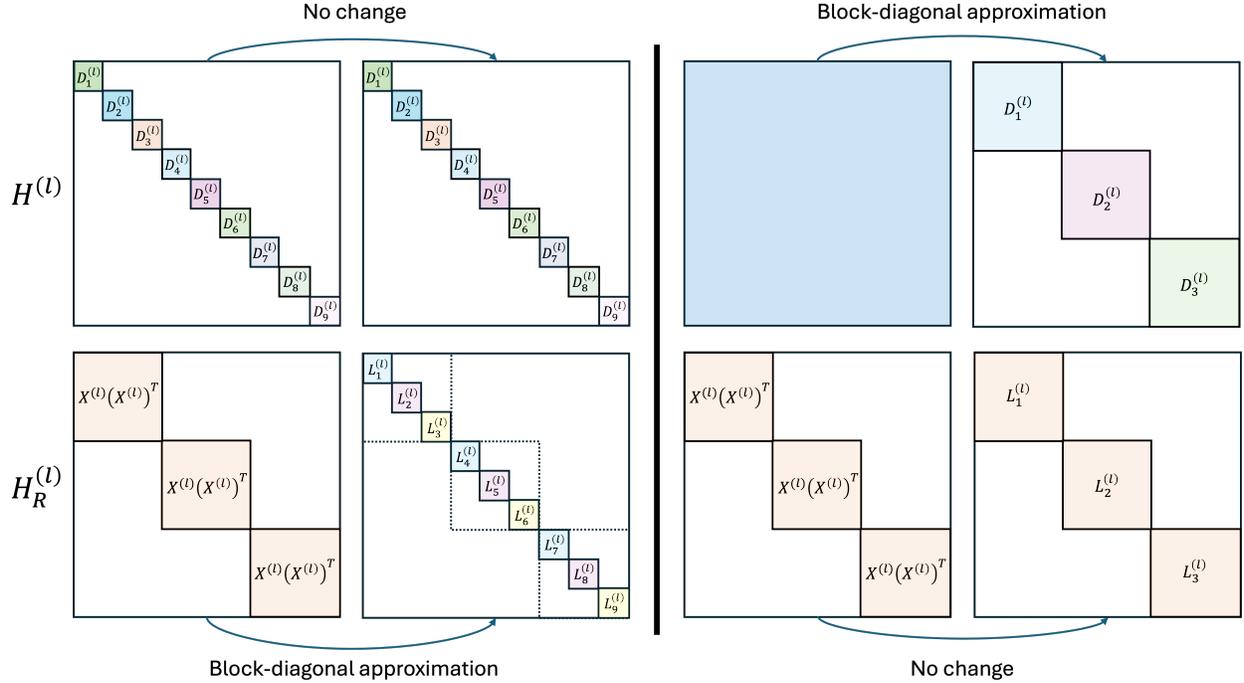


Figure 2: Depending on the block-diagonal approximation assumed by the single-objective algorithm, *MOONSHOT* matches the size of the original block-diagonal approximation to the Hessian of the other objective. [Left] Case 1: When adapting a Fisher-objective algorithm, we keep the block-diagonal approximation of  $H^{(l)}$  from the baseline unchanged, and perform a block-diagonal approximation on  $H_R^{(l)}$  (more precisely on  $X^{(l)}(X^{(l)})^T$ ). [Right] Case 2: When adapting a layer-wise reconstruction objective algorithm, we keep the exact block-diagonal form of  $H_R^{(l)}$ , as in the original baseline, and perform a block-diagonal approximation on  $H^{(l)}$ .

### 2.3 Efficient inverse Hessian computation for LLMs

Most state-of-the-art single-objective pruning methods (Frantar et al., 2022; Kuznedeleev et al., 2023; Frantar & Alistarh, 2023; Sun et al., 2024; Meng et al., 2024b) require access to the inverse Hessian  $(F_k^{(l)})^{-1}$  for each layer  $l$  and block  $k$  in order to compute the impact of pruning a weight  $w_p$  on the objective, as well as the potential update on the support  $\delta_p$ , as described in the OBS algorithm (Hassibi & Stork, 1992b):

$$p = \operatorname{argmin}_p \frac{[w_k^{(l)}]_p^2}{[(F_k^{(l)})^{-1}]_{p,p}}, \quad \delta_p = -\frac{[w_k^{(l)}]_p}{[(F_k^{(l)})^{-1}]_{p,p}} [(F_k^{(l)})^{-1}]_{:,p} \quad (11)$$

In the case of structured pruning, this formula can be extended to determine the impact of pruning an entire column (Kurtic et al., 2023; Meng et al., 2024b) but still requires computing  $(F_k^{(l)})^{-1}$ . While this inverse Hessian can be computed efficiently for vision models (as in OBC (Frantar et al., 2022) and CAP (Kuznedeleev et al., 2023)), the computation becomes significantly more challenging in the context of LLMs. This is due to the larger number of blocks and the higher dimensions of  $F_k^{(l)}$ . Indeed, state-of-the-art layer-wise pruning algorithms like SparseGPT (Frantar & Alistarh, 2023) and OSSCAR (Meng et al., 2024b) use the layer-wise reconstruction loss with no further block-diagonal approximation:  $L_1^{(l)} = \dots = L_K^{(l)} = X^{(l)}(X^{(l)})^T$ . This design requires just one matrix inversion of  $X^{(l)}(X^{(l)})^T$ , which makes the algorithm practical even for large models. In our multi-objective formulation, by contrast, each block  $F_k^{(l)}$  is different due to the Fisher loss component. Therefore, a naive adaptation of such algorithms to the multi-objective formulation would require computing  $K$  matrix inversions (one for each block) instead of one, which would significantly slow

**Algorithm 1** Efficient Computation of the Block-Diagonal Hessian Inverse

**Input:** Layer input matrix  $X^{(l)} \in \mathbb{R}^{d_{\text{in}}^{(l)} \times N}$ , per-sample gradients  $A_k^{(l)} \in \mathbb{R}^{d_{\text{in}}^{(l)} \times N}$  for each block  $k = 1, \dots, K$ , multi-objective weight  $\lambda \in [0, 1]$ .

1: Compute base inverse:

$$J_0 \leftarrow \left( \frac{\lambda}{\mathcal{L}_R^{(l)}(\mathbf{0})} X^{(l)} (X^{(l)})^T \right)^{-1}$$

2: **for** each block  $k = 1, \dots, K$  **do**

3:   Compute  $G_k^{(l)}$  using the Woodbury identity (see Appendix A.3):

$$G_k^{(l)} \leftarrow J_0 - \left( \frac{1-\lambda}{N \mathcal{L}_F^{(l)}(\mathbf{0})} \right) J_0 A_k^{(l)} \left( I_N + \frac{1-\lambda}{N \mathcal{L}_F^{(l)}(\mathbf{0})} A_k^{(l)\top} J_0 A_k^{(l)} \right)^{-1} A_k^{(l)\top} J_0$$

**Output:** Block-diagonal Hessians inverse  $\{(F_k^{(l)})^{-1}\}_{k=1}^K = \{G_k^{(l)}\}_{k=1}^K$

down the algorithm. However, by leveraging the structure of our multi-objective formulation, we can compute the inverse Hessian efficiently. In particular, for a layer  $l$ , the Hessian component from the layer-wise reconstruction error  $X^{(l)}(X^{(l)})^T$  does not depend on the block  $k$ . In addition, the size of calibration sets  $N$  for LLMs is often small (128 for SparseGPT, Wanda and OSSCAR). Consequently, the Hessian component coming from the Fisher loss  $H_k^{(l)}$  is a matrix of low rank  $r \leq N \ll d_{\text{in}}^{(l)}$  that can be written as  $H_k^{(l)} = \frac{1}{N} A_k^{(l)} A_k^{(l)\top}$  with  $A_k^{(l)} = [\nabla \ell_{1,k}^{(l)} \quad \nabla \ell_{2,k}^{(l)} \quad \dots \quad \nabla \ell_{N,k}^{(l)}] \in \mathbb{R}^{d_{\text{in}}^{(l)} \times N}$  ( $d_{\text{in}}^{(l)}$  is the block-size in this case). We therefore propose to compute  $G_k^{(l)} = (F_k^{(l)})^{-1}$ , necessary for the state-of-the-art OBS strategy used in SparseGPT, OBC, Wanda, CAP, OSSCAR, following the procedure described in Algorithm 1. Our exact adaptation of the SparseGPT algorithm, denoted *MOONSHOT*-SparseGPT, is provided in Appendix A.2.

Here,  $I_N \in \mathbb{R}^{N \times N}$  is the identity matrix, and  $I_N + \left( \frac{1-\lambda}{N \mathcal{L}_F^{(l)}(\mathbf{0})} \right) A_k^{(l)} J_0 A_k^{(l)\top}$  is of size  $N \times N$ . Therefore, the  $N \times N$  matrix inversion and the Woodbury identity (Woodbury & of Statistics, 1950) can be computed very efficiently (at most 5-6 seconds for the largest layers of Llama-3.2-3B) in the case we are interested in ( $N = 128$ ). In particular, we observe a speedup of up to 6 times in comparison to inverting all the blocks using the standard matrix inversion via Cholesky decomposition (used in SparseGPT and OBC for example). While previous work also use the Woodbury identity to compute the Hessian inverse in the literature (Singh & Alistarh, 2020a; Kurtic et al., 2022), we extend in this paper its application to the billion-parameter scale and adapt it to the multi-objective pruning setting. The exact derivation of the update in Algorithm 1 is provided in Appendix A.3. It is important to note that the steps described above enable exact computation of the Hessian inverse for the block-diagonal approximation shown in Figure 2.

### 3 Experiments

#### 3.1 Models and Datasets

We evaluate the performance of our method on a wide range of models and baselines. Thus we prune:

- **Several Llama models:** Llama-3.2-1B (1B parameters) (Grattafiori et al., 2024), Llama-3.2-3B (3B parameters) (Grattafiori et al., 2024) and Llama-2-13b-chat-hf (Touvron et al., 2023) (13B parameters) using SparseGPT (Frantar & Alistarh, 2023), Wanda (Sun et al., 2024) and OSSCAR (Meng et al., 2024b)
- The DeiT Vision Transformers (Touvron et al., 2021): DeiT-Tiny (5.7M parameters), DeiT-Small (22.1M parameters) and DeiT-Base (86.6M parameters) using CAP (Kuznedelev et al., 2023)
- A Convolutional Neural Network (He et al., 2015): ResNet-50 (25.6M parameters) using OBC (Frantar et al., 2022)

We additionally prune the Instruct variants of Llama-3.2-1B and Llama-3.2-3B using SparseGPT and Wanda and report the results in Appendix A.9.

OSSCAR (Meng et al., 2024b) greedily prunes columns by optimizing a layer-wise reconstruction objective, with an optional local-search refinement. In our experiments, we use OSSCAR default hyperparameters and focus on OSSCAR greedy pruning step.

For pruning, we use 128 samples for the LLMs from the C4 (Raffel et al., 2020) dataset, and 4096 samples from ImageNet-1k (Deng et al., 2009) for the vision models. While we focus on the test accuracy on ImageNet-1k for the vision models, we focus on both perplexity and zero-shot performance for Llama-3.2. Following previous work (Frantar & Alistarh, 2023; Sun et al., 2024), we compute the test perplexity on WikiText2 (Merity et al., 2016) and PTB (Marcus et al., 1994). Additionally, we assess the zero-shot accuracy of the pruned LLMs on a variety of common-sense reasoning datasets, including BoolQ (Clark et al., 2019), PIQA (Bisk et al., 2020), HellaSwag (Zellers et al., 2019), WinoGrande (Sakaguchi et al., 2021), ARC-easy (Clark et al., 2018), ARC-challenge (Clark et al., 2018), and OpenBookQA (Mihaylov et al., 2018). In addition to mean performance across the seven classification benchmarks, we also report the win rate, i.e., the percentage of benchmarks on which one method outperforms the other.

### 3.2 Setup

We prune Llama-3.2-3B and Llama-2-13b-chat-hf on a single NVIDIA A100 GPU (80 GB) and Llama-3.2-1B on a single NVIDIA L40 GPU (40 GB). For the vision models, we use four NVIDIA L40 GPUs (40 GB each) and prune the layers across the four devices. *MOONSHOT* is implemented in PyTorch (Paszke et al., 2019).

### 3.3 Implementation Details

**Pruning blocks of rows with *MOONSHOT-SparseGPT* and *MOONSHOT-OSSCAR*.** Unlike the single-objective setting where the Hessian is block-diagonal with identical blocks repeated,  $H = \text{Diag}(XX^\top, \dots, XX^\top)$ , *MOONSHOT-SparseGPT* and *MOONSHOT-OSSCAR* with  $\lambda \neq 1$  uses a Hessian with row-dependent blocks,  $H = \text{Diag}(F_1, \dots, F_K)$ . In the case of the largest layers, storing all the blocks simultaneously can become infeasible under GPU memory constraints. Our adaption of SparseGPT and OSSCAR prunes the rows by blocks of size  $K_p$ . The exact adaptation of SparseGPT is described in Appendix A.2.

**Efficient Backsolve for OSSCAR.** After the greedy column-selection step, OSSCAR performs a backsolve (i.e., it computes the optimal weights on the support). To preserve efficiency in the multi-objective setting, we exploit problem structure to perform this backsolve efficiently. Additional details are provided in Appendix A.4.

**Pruning Efficiency.** In the case of the attention layers, a sufficiently large  $K_p$  can be used; however, for the much larger projection layers,  $K_p$  often needs to be reduced, which can increase runtime due to reduced parallelism. Moreover, the projection layers typically require a larger  $H$ , further increasing computational cost. To maintain the efficiency of the original method in the case of LLMs, only the self-attention layers are pruned using the multi-objective formulation. Concretely, this corresponds `q_proj`, `k_proj`, `v_proj`, and `o_proj` for SparseGPT, and `o_proj` for OSSCAR (OSSCAR prunes only the `down_proj` and `o_proj` matrices.) The projection layers are pruned using the layer-wise reconstruction loss only. For SparseGPT, we select  $K_p$  such that at least 50% of rows are pruned at a time for Llama-3.2-1B and Llama-3.2-3B, and fix  $K_p = 512$  for Llama-2-13b-chat-hf. For OSSCAR,  $K_p$  corresponds to 50% of the rows for Llama-3.2-1B and 25% of the rows for Llama-3.2-3B. An evaluation of *MOONSHOT*’s effectiveness with the pruning times of each method is included in Appendix A.5.

We additionally evaluate the impact of applying *MOONSHOT* across both the attention and projection layers of Llama, both in terms of performance and computational cost, in Appendix A.8.

**Hessian Recomputation.** Fisher-based methods typically compute  $H^{(l)}$  once, as recomputation requires per-sample gradients, and they report results with Hessian recomputation as a more costly alternative (Benbaki et al., 2023; Kuznedelev et al., 2023). In contrast, layer-wise reconstruction-based methods like SparseGPT and Wanda recompute  $H_R^{(l)}$  after each block of layers, as  $H_R^{(l)}$  depends only on the input data and can be recomputed with relatively low overhead (Frantar & Alistarh, 2023; Sun et al., 2024). In

this paper, due to the multi-objective formulation, we follow the standard approach in the Fisher-based literature, and compute the Hessian (inverse) once. We also report results with Hessian recomputation after each block for SparseGPT and Wanda in Appendix A.7.

**Selecting  $\lambda$  in equation 10.** The results provided in Tables 3, 2 and in Figure 4 are obtained by selecting the best value of  $\lambda$  based on the training loss for vision models and training perplexity for LLMs. For the vision models, we evaluate  $\lambda \in \{0.0, 0.25, 0.5, 0.75, 1.0\}$ , and for the [Llama-3.2 models](#), we test  $\lambda \in \{0.0, 0.1, 0.25, 0.5, 0.75, 0.9, 1.0\}$ . [For Llama-2-13b-chat-hf, we only test  \$\lambda \in \{0.9, 1.0\}\$ .](#)

While  $\lambda$  is selected via tuning, Section 4.1 shows that  $\lambda \in (0, 1)$  - that is, beyond the standard single-objective baselines ( $\lambda = 0$  or  $1$ ) - almost always lead to better performance. In addition,  $\lambda = 0.5$  [for vision models](#) and  $\lambda = 0.9$  [for LLMs](#) serve as simple and effective defaults in resource- or time-constrained scenarios. For typical pruning use cases, where pruning is performed once offline and the resulting model is used across multiple downstream tasks or applications, investing in hyperparameter tuning can further enhance the performance gains achieved by *MOONSHOT*.

**Hyperparameters.** Additional information on the hyperparameters used for the baselines and *MOONSHOT* is provided in Appendix A.6.

### 3.4 Main Results

Tables 2 and 3 report results at relevant sparsity levels: [10% structured and 60%/70% unstructured](#) for LLMs, [70% unstructured](#) for DeiT models, and [90% unstructured](#) for ResNet-50, in addition to 2:4 semi-structured sparsity. Across all settings, *MOONSHOT* consistently outperforms the baseline, yielding statistically significant improvements. Comprehensive results across architectures, sparsity regimes and  $\lambda$  values are available in Appendix A.12.

**Vision Models.** Table 2 shows that on DeiT-Small, *MOONSHOT* improves test accuracy on ImageNet-1k by up to 5.5 points compared to CAP at [70% unstructured](#) sparsity. This indicates that the Fisher-based Hessian used in CAP is insufficiently informative at this level of compression. In contrast, our multi-objective formulation yields a more stable and informative Hessian, resulting in a much higher quality pruned model. DeiT-Tiny and DeiT-Base also show consistent improvements of 1–3 points across both unstructured and semi-structured sparsity settings. For ResNet-50, *MOONSHOT* improves accuracy by 4 points at [90% unstructured](#) sparsity compared to OBC, and further improves performance at 2:4 sparsity.

**Language Models.** Table 3 shows that on Llama-3.2-1B, *MOONSHOT* lowers test perplexity on C4 by up to 54 points with Wanda at 2:4 sparsity and by 13 points with SparseGPT at [60% unstructured](#) sparsity. These improvements extend across other language modeling benchmarks (WikiText2, PTB) and generalize to downstream classification tasks, [where mean accuracy often improves by up to 1 point](#). Similar results are observed for [Llama-3.2-3B and Llama-2-13b-chat-hf](#), and *MOONSHOT* [improves the mean accuracy of these models by up to 1.5 points at 60% and 70% unstructured sparsity respectively](#).

Importantly, *MOONSHOT* complements existing sparsity allocation strategies. When combined with AlphaPruning or OWL, it yields additional performance gains. For example, on Llama-3.2-1B at [60% unstruc-](#)

Table 2: Impact of *MOONSHOT* on CAP for the DeiT models (left) and OBC for ResNet-50 (right) across unstructured and 2:4 sparsity levels. ImageNet-1k accuracies over 3 seeds are averaged with standard errors.

Sparsity	Method	DeiT Tiny	DeiT Small	DeiT Base	Sparsity	Method	ResNet-50
Dense	-	72.14	79.83	81.80	Dense	-	77.11
0.7	CAP	44.22 ± 0.32	57.50 ± 0.83	70.44 ± 0.15	0.9	OBC	51.52 ± 0.07
	<i>MOONSHOT</i> -CAP	<b>45.05 ± 0.20</b>	<b>62.97 ± 0.15</b>	<b>73.42 ± 0.06</b>		<i>MOONSHOT</i> -OBC	<b>55.52 ± 0.09</b>
2:4	CAP	52.28 ± 0.04	69.65 ± 0.02	76.21 ± 0.07	2:4	OBC	75.46 ± 0.03
	<i>MOONSHOT</i> -CAP	<b>54.20 ± 0.15</b>	<b>71.54 ± 0.08</b>	<b>77.88 ± 0.05</b>		<i>MOONSHOT</i> -OBC	<b>75.50 ± 0.03</b>

Table 3: Impact for the Llama-3.2 models of *MOONSHOT* on SparseGPT/Wanda at 60% unstructured sparsity (including with OWL and AlphaPruning), 2:4 sparsity, and OSSCAR at 10% structured sparsity. We also include Llama-2-13b-chat-hf at 70% unstructured sparsity. The perplexities on C4, WikiText2 and PTB, as well as the zero-shot accuracies, are averaged over 3 seeds with standard errors. Mean performance and win rate are computed over the 7 zero-shot downstream classification tasks.

(a) Llama-3.2-1B

Sparsity	Method	C4 ↓	WikiText2 ↓	PTB ↓	BoolQ ↑	HellaSwag ↑	WinoGrande ↑	ARC-e ↑	ARC-c ↑	PIQA ↑	OBQA ↑	Mean ↑	Win Rate ↑
Dense	-	14.02	9.75	17.59	64.01	47.73	60.14	65.15	31.23	74.32	26.4	52.71	-
0.1 (structured)	OSSCAR	43.00 ± 1.28	43.91 ± 0.60	106.62 ± 8.57	52.61 ± 1.63	35.82 ± 0.29	54.17 ± 1.11	45.19 ± 7.83	24.23 ± 2.04	65.65 ± 3.68	15.93 ± 0.33	41.94 ± 1.98	19.05 ± 12.60
	MOONSHOT-OSSCAR	<b>38.04 ± 0.58</b>	<b>30.93 ± 0.92</b>	<b>86.73 ± 9.88</b>	<b>56.55 ± 0.90</b>	<b>37.89 ± 1.11</b>	<b>55.51 ± 0.37</b>	<b>58.84 ± 0.38</b>	<b>28.58 ± 0.51</b>	<b>71.49 ± 0.06</b>	<b>18.93 ± 0.48</b>	<b>46.83 ± 0.42</b>	<b>80.95 ± 12.60</b>
0.6	SparseGPT	63.63 ± 1.18	54.60 ± 1.00	81.11 ± 3.99	60.67 ± 0.59	32.16 ± 0.20	54.46 ± 0.53	44.94 ± 0.11	21.47 ± 0.48	62.21 ± 0.20	17.07 ± 0.41	41.85 ± 0.20	42.86 ± 8.25
	MOONSHOT-SparseGPT	<b>50.28 ± 1.99</b>	<b>39.13 ± 1.54</b>	<b>60.14 ± 2.90</b>	<b>62.36 ± 0.12</b>	<b>32.49 ± 0.13</b>	53.09 ± 0.18	<b>46.49 ± 0.38</b>	21.30 ± 0.24	<b>63.22 ± 0.17</b>	15.73 ± 0.55	<b>42.10 ± 0.11</b>	<b>57.14 ± 8.25</b>
0.6 (Alpha-Pruning)	SparseGPT	61.05 ± 0.77	52.80 ± 0.43	78.27 ± 3.64	62.08 ± 0.26	32.00 ± 0.08	<b>53.88 ± 0.25</b>	45.29 ± 0.49	<b>22.01 ± 0.59</b>	62.02 ± 0.42	<b>17.60 ± 0.42</b>	42.13 ± 0.06	33.33 ± 9.52
	MOONSHOT-SparseGPT	<b>49.31 ± 1.10</b>	<b>38.44 ± 0.52</b>	<b>60.32 ± 1.20</b>	<b>62.29 ± 0.05</b>	<b>32.53 ± 0.06</b>	53.70 ± 0.55	<b>46.30 ± 0.30</b>	21.99 ± 0.15	<b>63.13 ± 0.10</b>	16.60 ± 0.12	<b>42.36 ± 0.14</b>	<b>66.67 ± 9.52</b>
0.6 (OWL)	SparseGPT	56.82 ± 1.68	49.54 ± 1.22	68.73 ± 3.48	<b>62.20 ± 0.11</b>	32.86 ± 0.05	<b>53.67 ± 0.33</b>	44.14 ± 0.71	<b>23.46 ± 0.05</b>	62.59 ± 0.21	<b>18.00 ± 0.76</b>	42.42 ± 0.07	33.33 ± 17.17
	MOONSHOT-SparseGPT	<b>43.58 ± 0.91</b>	<b>35.72 ± 0.17</b>	<b>53.02 ± 0.91</b>	62.20 ± 0.04	<b>33.66 ± 0.16</b>	53.54 ± 0.55	<b>46.37 ± 0.23</b>	23.41 ± 0.08	<b>64.04 ± 0.03</b>	16.40 ± 0.20	<b>42.80 ± 0.07</b>	<b>66.67 ± 17.17</b>
2:4	SparseGPT	53.59 ± 0.35	42.56 ± 0.37	63.79 ± 0.19	61.42 ± 0.21	<b>31.68 ± 0.05</b>	<b>53.83 ± 0.37</b>	44.04 ± 0.23	<b>21.47 ± 0.44</b>	61.79 ± 0.40	<b>15.00 ± 0.20</b>	<b>41.32 ± 0.08</b>	<b>57.14 ± 14.29</b>
	MOONSHOT-SparseGPT	<b>50.99 ± 0.47</b>	<b>38.00 ± 0.58</b>	<b>59.32 ± 1.55</b>	<b>61.98 ± 0.29</b>	31.47 ± 0.21	53.09 ± 0.27	<b>45.37 ± 0.50</b>	20.28 ± 0.58	<b>62.13 ± 0.19</b>	14.33 ± 0.35	41.24 ± 0.20	42.86 ± 14.29
0.6	Wanda	117.71 ± 0.87	84.73 ± 0.73	119.64 ± 1.00	58.96 ± 1.39	28.86 ± 0.03	51.35 ± 0.49	38.82 ± 0.32	18.94 ± 0.26	59.05 ± 0.18	<b>13.93 ± 0.24</b>	38.56 ± 0.12	19.05 ± 4.76
	MOONSHOT-Wanda	<b>86.55 ± 1.67</b>	<b>63.57 ± 1.61</b>	<b>98.44 ± 3.98</b>	<b>61.56 ± 0.29</b>	<b>29.53 ± 0.06</b>	<b>51.64 ± 0.23</b>	<b>40.40 ± 0.17</b>	<b>19.60 ± 0.06</b>	<b>61.12 ± 0.08</b>	13.40 ± 0.20	<b>39.61 ± 0.07</b>	<b>80.95 ± 4.76</b>
0.6 (Alpha-Pruning)	Wanda	112.33 ± 1.03	80.75 ± 1.03	120.89 ± 0.73	58.01 ± 1.10	28.92 ± 0.09	51.22 ± 0.47	37.56 ± 0.16	19.51 ± 0.21	59.32 ± 0.04	<b>13.40 ± 0.40</b>	38.28 ± 0.10	14.29 ± 8.25
	MOONSHOT-Wanda	<b>84.21 ± 1.04</b>	<b>63.30 ± 0.92</b>	<b>101.96 ± 2.98</b>	<b>61.69 ± 0.24</b>	<b>29.56 ± 0.08</b>	<b>52.33 ± 0.37</b>	<b>39.28 ± 0.29</b>	<b>19.65 ± 0.06</b>	<b>60.32 ± 0.13</b>	12.07 ± 0.58	<b>39.27 ± 0.09</b>	<b>85.71 ± 8.25</b>
0.6 (OWL)	Wanda	99.38 ± 1.37	73.00 ± 0.37	111.24 ± 1.83	61.28 ± 0.28	29.88 ± 0.11	<b>52.07 ± 0.82</b>	40.22 ± 0.09	<b>20.82 ± 0.05</b>	60.03 ± 0.20	<b>14.80 ± 0.12</b>	<b>39.87 ± 0.17</b>	33.33 ± 4.76
	MOONSHOT-Wanda	<b>73.97 ± 0.19</b>	<b>58.81 ± 0.28</b>	<b>100.07 ± 1.96</b>	<b>61.81 ± 0.08</b>	<b>30.47 ± 0.07</b>	51.51 ± 0.37	<b>40.92 ± 0.22</b>	20.71 ± 0.15	<b>60.36 ± 0.19</b>	13.27 ± 0.35	39.86 ± 0.15	<b>66.67 ± 4.76</b>
2:4	Wanda	164.32 ± 2.37	114.73 ± 2.32	190.58 ± 1.50	57.28 ± 1.00	28.32 ± 0.03	<b>51.51 ± 0.15</b>	35.76 ± 0.26	18.34 ± 0.09	58.41 ± 0.35	<b>13.67 ± 0.07</b>	37.61 ± 0.10	28.57 ± 0.00
	MOONSHOT-Wanda	<b>110.74 ± 1.17</b>	<b>78.55 ± 1.14</b>	<b>126.91 ± 1.66</b>	<b>61.08 ± 0.57</b>	<b>28.51 ± 0.05</b>	50.62 ± 0.25	<b>38.41 ± 0.26</b>	<b>19.43 ± 0.08</b>	<b>59.36 ± 0.33</b>	12.67 ± 0.37	<b>38.58 ± 0.14</b>	<b>71.43 ± 0.00</b>

(b) Llama-3.2-3B

Sparsity	Method	C4 ↓	WikiText2 ↓	PTB ↓	BoolQ ↑	HellaSwag ↑	WinoGrande ↑	ARC-e ↑	ARC-c ↑	PIQA ↑	OBQA ↑	Mean ↑	Win Rate ↑
Dense	-	11.34	7.81	13.54	72.72	55.30	69.22	74.37	42.41	76.71	31.20	60.28	-
0.1 (structured)	OSSCAR	16.80 ± 0.46	13.76 ± 0.85	22.34 ± 0.43	64.38 ± 0.86	49.03 ± 0.86	61.30 ± 0.74	65.82 ± 0.64	34.10 ± 0.53	<b>75.70 ± 0.78</b>	<b>27.33 ± 0.41</b>	<b>53.95 ± 0.08</b>	<b>33.33 ± 4.76</b>
	MOONSHOT-OSSCAR	<b>16.56 ± 0.38</b>	<b>13.28 ± 0.55</b>	<b>22.10 ± 0.28</b>	<b>64.91 ± 1.00</b>	<b>49.50 ± 0.62</b>	<b>61.43 ± 0.47</b>	<b>66.39 ± 0.83</b>	<b>34.41 ± 0.33</b>	<b>75.54 ± 0.45</b>	<b>26.80 ± 0.76</b>	<b>54.14 ± 0.20</b>	<b>66.67 ± 4.76</b>
0.6	SparseGPT	33.63 ± 0.14	26.12 ± 0.23	42.69 ± 0.73	66.82 ± 0.60	38.14 ± 0.14	60.91 ± 0.65	53.89 ± 0.11	26.28 ± 0.18	67.75 ± 0.34	18.47 ± 0.44	47.47 ± 0.08	4.76 ± 4.76
	MOONSHOT-SparseGPT	<b>28.23 ± 0.11</b>	<b>22.46 ± 0.17</b>	<b>35.63 ± 0.68</b>	<b>67.76 ± 0.35</b>	<b>39.13 ± 0.07</b>	<b>61.01 ± 0.23</b>	<b>57.59 ± 0.92</b>	<b>27.79 ± 0.71</b>	<b>69.44 ± 0.13</b>	<b>20.00 ± 0.53</b>	<b>48.96 ± 0.12</b>	<b>95.24 ± 4.76</b>
0.6 (Alpha-Pruning)	SparseGPT	34.19 ± 0.40	26.06 ± 0.64	42.82 ± 0.13	68.17 ± 0.43	38.62 ± 0.16	61.98 ± 0.81	53.37 ± 0.70	26.00 ± 0.23	68.06 ± 0.38	19.80 ± 0.90	48.00 ± 0.38	14.29 ± 8.25
	MOONSHOT-SparseGPT	<b>28.64 ± 0.25</b>	<b>22.17 ± 0.19</b>	<b>35.48 ± 1.52</b>	<b>68.73 ± 0.11</b>	<b>39.34 ± 0.18</b>	<b>62.27 ± 0.52</b>	<b>56.52 ± 0.25</b>	<b>26.93 ± 0.37</b>	<b>69.13 ± 0.32</b>	<b>20.33 ± 0.24</b>	<b>49.04 ± 0.06</b>	<b>85.71 ± 8.25</b>
0.6 (OWL)	SparseGPT	29.15 ± 0.31	23.58 ± 0.40	36.58 ± 0.78	66.64 ± 0.77	39.89 ± 0.16	<b>61.96 ± 0.40</b>	55.57 ± 0.41	27.53 ± 0.12	68.72 ± 0.14	<b>21.20 ± 0.35</b>	48.79 ± 0.11	28.57 ± 8.25
	MOONSHOT-SparseGPT	<b>25.31 ± 0.12</b>	<b>20.85 ± 0.11</b>	<b>31.39 ± 0.69</b>	<b>67.41 ± 0.47</b>	<b>40.64 ± 0.13</b>	61.62 ± 0.09	<b>57.39 ± 0.78</b>	<b>28.16 ± 0.62</b>	<b>69.73 ± 0.29</b>	20.60 ± 0.50	<b>49.36 ± 0.17</b>	<b>71.43 ± 8.25</b>
2:4	SparseGPT	30.00 ± 0.29	24.40 ± 0.23	38.32 ± 0.74	<b>65.64 ± 0.70</b>	<b>38.31 ± 0.08</b>	<b>59.93 ± 0.13</b>	55.63 ± 1.10	26.14 ± 0.42	<b>68.34 ± 0.33</b>	20.87 ± 0.35	<b>47.83 ± 0.18</b>	<b>52.38 ± 9.52</b>
	MOONSHOT-SparseGPT	<b>28.79 ± 0.29</b>	<b>23.21 ± 0.32</b>	<b>35.94 ± 0.73</b>	65.58 ± 0.08	38.06 ± 0.13	59.30 ± 0.41	<b>55.99 ± 0.63</b>	<b>26.56 ± 0.37</b>	67.94 ± 0.21	<b>20.93 ± 0.55</b>	47.77 ± 0.18	47.62 ± 9.52
0.6	Wanda	41.98 ± 0.40	30.56 ± 0.32	51.00 ± 0.45	<b>64.82 ± 0.35</b>	35.12 ± 0.07	<b>56.56 ± 0.46</b>	50.58 ± 0.41	23.83 ± 0.12	65.58 ± 0.19	<b>16.93 ± 0.07</b>	<b>44.77 ± 0.10</b>	38.10 ± 4.76
	MOONSHOT-Wanda	<b>37.73 ± 0.19</b>	<b>27.71 ± 0.26</b>	<b>46.47 ± 0.08</b>	61.33 ± 0.91	<b>35.53 ± 0.08</b>	54.83 ± 0.14	<b>52.53 ± 0.29</b>	<b>24.69 ± 0.21</b>	<b>66.81 ± 0.03</b>	16.60 ± 0.23	44.62 ± 0.12	<b>61.90 ± 4.76</b>
0.6 (Alpha-Pruning)	Wanda	40.03 ± 0.04	29.19 ± 0.20	50.24 ± 0.27	<b>65.93 ± 0.24</b>	35.95 ± 0.01	<b>57.51 ± 0.09</b>	51.09 ± 0.16	24.32 ± 0.36	66.03 ± 0.15	<b>16.87 ± 0.24</b>	45.39 ± 0.03	38.10 ± 4.76
	MOONSHOT-Wanda	<b>37.37 ± 0.21</b>	<b>27.08 ± 0.16</b>	<b>47.01 ± 0.29</b>	63.38 ± 0.65	<b>36.05 ± 0.10</b>	57.51 ± 0.56	<b>54.15 ± 0.27</b>	<b>25.43 ± 0.18</b>	<b>66.96 ± 0.36</b>	16.27 ± 0.18	<b>45.68 ± 0.10</b>	<b>61.90 ± 4.76</b>
0.6 (OWL)	Wanda	37.35 ± 0.14	27.93 ± 0.23	44.25 ± 0.74	<b>67.26 ± 0.07</b>	37.18 ± 0.13	<b>59.06 ± 0.73</b>	51.89 ± 0.34	25.54 ± 0.23	66.47 ± 0.10	<b>17.20 ± 0.12</b>	<b>46.37 ± 0.18</b>	33.33 ± 9.52
	MOONSHOT-Wanda	<b>34.56 ± 0.11</b>	<b>25.59 ± 0.06</b>	<b>40.41 ± 0.78</b>	63.73 ± 0.92	<b>37.40 ± 0.12</b>	58.93 ± 0.35	<b>54.31 ± 0.07</b>	<b>25.68 ± 0.13</b>	<b>67.05 ± 0.25</b>	17.00 ± 0.23	46.30 ± 0.13	<b>66.67 ± 9.52</b>
2:4	Wanda	49.79 ± 0.26	35.90 ± 0.37	68.16 ± 0.29	<b>64.29 ± 0.13</b>	34.16 ± 0.06	<b>55.88 ± 0.36</b>	50.88 ± 0.23	25.28 ± 0.03	65.25 ± 0.10	17.13 ± 0.24	<b>44.70 ± 0.06</b>	38.10 ± 12.60
	MOONSHOT-Wanda	<b>45.10 ± 0.26</b>	<b>32.47 ± 0.50</b>	<b>61.19 ± 0.23</b>	61.60 ± 0.86	<b>34.41 ± 0.16</b>	55.51 ± 0.66	<b>52.53 ± 0.18</b>	<b>25.60 ± 0.38</b>	<b>65.40 ± 0.17</b>	<b>17.20 ± 0.61</b>	44.61 ± 0.09	<b>61.90 ± 12.60</b>

(c) Llama-2-13b-chat-hf

Sparsity	Method	C4 ↓	WikiText2 ↓	PTB ↓	BoolQ ↑	HellaSwag ↑	WinoGrande ↑	ARC-e ↑	ARC-c ↑	PIQA ↑	OBQA ↑	Mean ↑	Win Rate ↑
Dense	-	8.49	6.11	50.36	81.65	60.71	71.11	77.53	46.16	77.91	35.20	64.32	-
0.7	SparseGPT	27.62 ± 0.31	24.87 ± 0.41	460.75 ± 16.25	71.86 ± 1.18	38.59 ± 0.1	59.93 ± 0.27	53.65 ± 1.11	<b>28.21 ± 0.16</b>	67.16 ± 0.54	<b>22.73 ± 0.07</b>	48.88 ± 0.28	19.05 ± 4.76
	MOONSHOT-SparseGPT	<b>23.67 ± 0.06</b>	<b>20.95 ± 0.59</b>	<b>368.57 ± 5.75</b>	<b>75.71 ± 0.5</b>	<b>40.01 ± 0.27</b>	<b>61.12 ± 0.42</b>	<b>57.41 ± 0.72</b>	28.16 ± 0.23	<b>68.64 ± 0.28</b>	21.73 ± 0.77	<b>50.4 ± 0.3</b>	<b>80.95 ± 4.76</b>
0.7	Wanda	46.15 ± 0.16	47.85 ± 1.0	629.11 ± 9.28	64.05 ± 0.04	32.17 ± 0.13	54.22 ± 0.28	43.95 ± 0.27	20.71 ± 0.17	61.35 ± 0.3	<b>17.0 ± 0.12</b>	41.92 ± 0.07	19.05 ± 4.76
	MOONSHOT-Wanda	<b>41.0 ± 0.25</b>	<b>38.38 ± 1.31</b>	<b>607.22 ± 8.48</b>	<b>66.68 ± 0.13</b>	<b>34.11 ± 0.08</b>	<b>54.75 ± 0.38</b>	<b>48.22 ± 0.27</b>	<b>21.16 ± 0.09</b>	<b>64.4 ± 0.1</b>	14.53 ± 0.52	<b>43.41 ± 0.1</b>	<b>80.95 ± 4.76</b>

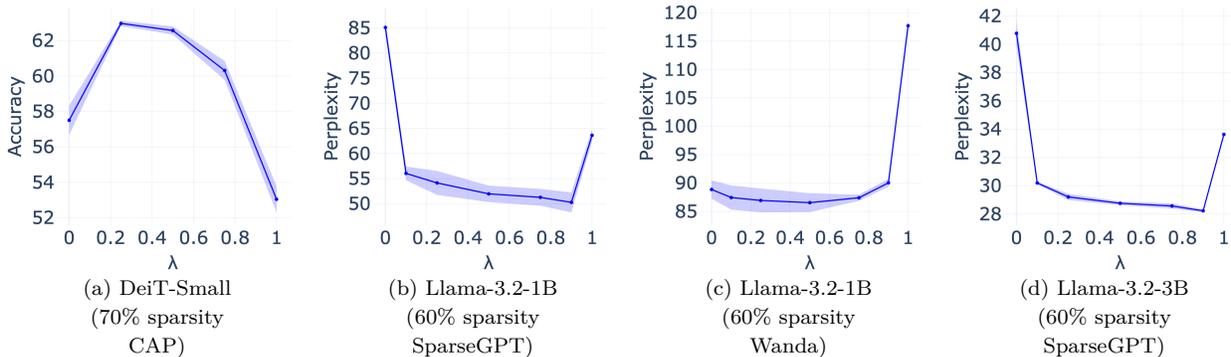


Figure 3: Performance of *MOONSHOT* across values of  $\lambda$  on DeiT-Small using CAP (70% sparsity) and Llama-3.2 models using SparseGPT/Wanda (60% and 2:4 sparsity). ImageNet-1k accuracies for DeiT-Small and C4 perplexities for the Llama-3.2 models are averaged over 3 seeds with standard errors.

tuned sparsity, combining *MOONSHOT* with OWL leads to a further 13-point reduction in C4 perplexity and a 0.4-point increase in mean downstream accuracy.

In the case of structured pruning, the gains are particularly high, with up to 30% lower perplexity on WikiText2, 22% on PTB, and 11% on C4, together with a +4.9 point improvement in mean accuracy.

Finally, in terms of win rate, *MOONSHOT* outperforms the baseline on most benchmarks across sparsity regimes, architectures, and pruning baselines.

## 4 Ablation studies

### 4.1 Selecting $\lambda$

To demonstrate the efficacy of our proposed multi-objective formulation, we try different values of  $\lambda$  in equation 10.  $\lambda$  determines the balance between the layer-wise reconstruction error and the Fisher loss.

Figure 3 below, and Figure 5 in Appendix A.10, illustrate that setting neither  $\lambda = 0.0$  nor  $\lambda = 1.0$  achieves the best results on ResNet-50, DeiT-Base the Llama-3.2 models. The results are striking for the Llama models, for which the test perplexity on C4 is substantially lower in the multi-objective regime than the single objective regime. An intermediate value of  $\lambda$  that leverages the advantages of both loss functions is more effective. Furthermore,  $\lambda$  seems to require minimal tuning, as a value of  $\lambda$  other than 0 and 1 is often sufficient to achieve near-optimal performance.  $\lambda = 0.5$  for vision models and  $\lambda = 0.9$  for LLMs are relatively good choices across all architectures and sparsity levels tested.

### 4.2 Performance of *MOONSHOT* Across Sparsity Regimes

Figure 4 shows that *MOONSHOT* consistently outperforms the baselines across all tested sparsity levels. Performance gains are especially pronounced at higher sparsity levels, where preserving original performance is increasingly difficult. Additional results can be found in Appendix 6.

## 5 Conclusion

We present *MOONSHOT*, a framework that replaces the traditional single-objective formulation in one-shot pruning algorithms with a multi-objective approach. By incorporating both the layer-wise reconstruction loss (a local objective) and the second-order Taylor approximation of the training loss (a global objective), *MOONSHOT* significantly enhances the performance of state-of-the-art single-objective algorithms. Beyond these performance improvements, our work shows that generalizing existing pruning algorithms to a multi-

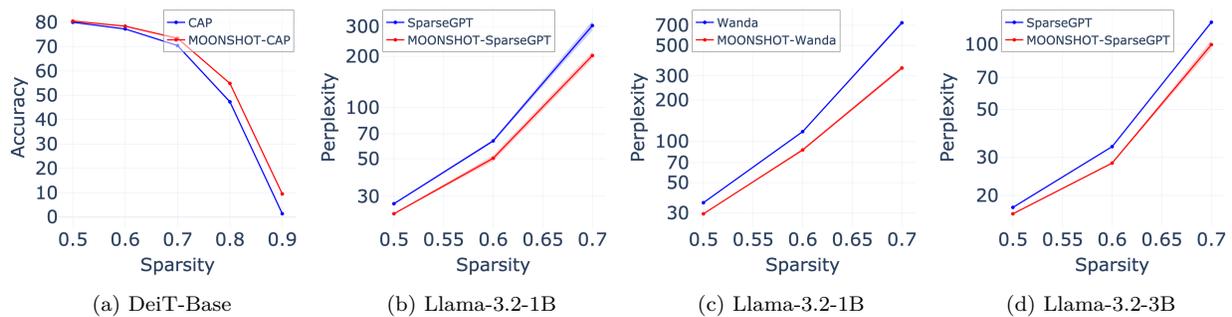


Figure 4: Impact of *MOONSHOT* across sparsity levels on CAP for DeiT-Base, and SparseGPT/Wanda on the Llama-3.2 models. ImageNet-1k accuracies for DeiT-Base and C4 perplexities for the Llama-3.2 models are averaged over 3 seeds with standard errors.

objective framework can be done efficiently to scale to modern large language models, making it a compelling approach for real-world applications.

## References

- Riade Benbaki, Wenyu Chen, Xiang Meng, Hussein Hazimeh, Natalia Ponomareva, Zhe Zhao, and Rahul Mazumder. Fast as CHITA: Neural network pruning with combinatorial optimization. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 2031–2049. PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/benbaki23a.html>.
- Yonatan Bisk, Rowan Zellers, Ronan Le bras, Jianfeng Gao, and Yejin Choi. Piqa: Reasoning about physical commonsense in natural language. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(05): 7432–7439, Apr. 2020. doi: 10.1609/aaai.v34i05.6239. URL <https://ojs.aaai.org/index.php/AAAI/article/view/6239>.
- Davis Blalock, Jose Javier Gonzalez Ortiz, Jonathan Frankle, and John Gutttag. What is the state of neural network pruning? *Proceedings of machine learning and systems*, 2:129–146, 2020.
- Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. BoolQ: Exploring the surprising difficulty of natural yes/no questions. In Jill Burstein, Christy Doran, and Thamar Solorio (eds.), *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 2924–2936, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1300. URL <https://aclanthology.org/N19-1300>.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge, 2018. URL <https://arxiv.org/abs/1803.05457>.
- Shail Dave, Riyadh Baghdadi, Tony Nowatzki, Sasikanth Avancha, Aviral Shrivastava, and Baoxin Li. Hardware acceleration of sparse and irregular tensor computations of ml models: A survey and insights. *Proceedings of the IEEE*, 109(10):1706–1752, 2021. doi: 10.1109/JPROC.2021.3098483.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.

- Xin Dong, Shangyu Chen, and Sinno Pan. Learning to prune deep neural networks via layer-wise optimal brain surgeon. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL [https://proceedings.neurips.cc/paper\\_files/paper/2017/file/c5dc3e08849bec07e33ca353de62ea04-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/c5dc3e08849bec07e33ca353de62ea04-Paper.pdf).
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=YicbFdNTTy>.
- Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=rJl-b3RcF7>.
- Jonathan Frankle, Gintare Karolina Dziugaite, Daniel Roy, and Michael Carbin. Pruning neural networks at initialization: Why are we missing the mark? In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=Ig-VyQc-MLK>.
- Elias Frantar and Dan Alistarh. SPDY: accurate pruning with speedup guarantees. *CoRR*, abs/2201.13096, 2022. URL <https://arxiv.org/abs/2201.13096>.
- Elias Frantar and Dan Alistarh. Sparsegpt: Massive language models can be accurately pruned in one-shot, 2023. URL <https://arxiv.org/abs/2301.00774>.
- Elias Frantar, Eldar Kurtic, and Dan Alistarh. M-fac: Efficient matrix-free approximations of second-order information. *Advances in Neural Information Processing Systems*, 34:14873–14886, 2021.
- Elias Frantar, Sidak Pal Singh, and Dan Alistarh. Optimal Brain Compression: a framework for accurate post-training quantization and pruning. *Advances in Neural Information Processing Systems*, 36, 2022.
- Trevor Gale, Erich Elsen, and Sara Hooker. The state of sparsity in deep neural networks, 2019.
- Trevor Gale, Matei Zaharia, Cliff Young, and Erich Elsen. Sparse gpu kernels for deep learning. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis, SC '20*. IEEE Press, 2020. ISBN 9781728199986.
- Mitchell A Gordon, Kevin Duh, and Nicholas Andrews. Compressing bert: Studying the effects of weight pruning on transfer learning. *arXiv preprint arXiv:2002.08307*, 2020.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, et al. The llama 3 herd of models, 2024. URL <https://arxiv.org/abs/2407.21783>.
- Song Han, Jeff Pool, John Tran, and William Dally. Learning both weights and connections for efficient neural network. *Advances in neural information processing systems*, 28, 2015.
- Stephen Hanson and Lorien Pratt. Comparing biases for minimal network construction with back-propagation. *Advances in neural information processing systems*, 1, 1988.
- Babak Hassibi and David Stork. Second order derivatives for network pruning: Optimal brain surgeon. *Advances in neural information processing systems*, 5, 1992a.
- Babak Hassibi and David Stork. Second order derivatives for network pruning: Optimal brain surgeon. In S. Hanson, J. Cowan, and C. Giles (eds.), *Advances in Neural Information Processing Systems*, volume 5. Morgan-Kaufmann, 1992b. URL [https://proceedings.neurips.cc/paper\\_files/paper/1992/file/303ed4c69846ab36c2904d3ba8573050-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/1992/file/303ed4c69846ab36c2904d3ba8573050-Paper.pdf).
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition, 2015.

- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Yihui He, Xiangyu Zhang, and Jian Sun. Channel pruning for accelerating very deep neural networks. In *Proceedings of the IEEE international conference on computer vision*, pp. 1389–1397, 2017.
- Itay Hubara, Yury Nahshan, Yair Hanani, Ron Banner, and Daniel Soudry. Accurate post training quantization with small calibration sets. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 4466–4475. PMLR, 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/hubara21a.html>.
- Eldar Kurtic, Daniel Campos, Tuan Nguyen, Elias Frantar, Mark Kurtz, Benjamin Fineran, Michael Goin, and Dan Alistarh. The optimal BERT surgeon: Scalable and accurate second-order pruning for large language models. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 4163–4181, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.279. URL <https://aclanthology.org/2022.emnlp-main.279/>.
- Eldar Kurtic, Elias Frantar, and Dan Alistarh. ZipLM: Inference-aware structured pruning of language models. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=d8j31sBwPv>.
- Denis Kuznedelev, Eldar Kurtic, Elias Frantar, and Dan Alistarh. CAP: Correlation-aware pruning for highly-accurate sparse vision models. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=Xy7DoWSNZX>.
- Yann LeCun, John Denker, and Sara Solla. Optimal brain damage. *Advances in neural information processing systems*, 2, 1989a.
- Yann LeCun, John Denker, and Sara Solla. Optimal brain damage. In D. Touretzky (ed.), *Advances in Neural Information Processing Systems*, volume 2. Morgan-Kaufmann, 1989b. URL [https://proceedings.neurips.cc/paper\\_files/paper/1989/file/6c9882bbac1c7093bd25041881277658-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/1989/file/6c9882bbac1c7093bd25041881277658-Paper.pdf).
- Namhoon Lee, Thalaisyasingam Ajanthan, and Philip Torr. SNIP: SINGLE-SHOT NETWORK PRUNING BASED ON CONNECTION SENSITIVITY. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=B1VZqjAcYX>.
- Namhoon Lee, Thalaisyasingam Ajanthan, Stephen Gould, and Philip H. S. Torr. A signal propagation perspective for pruning neural networks at initialization, 2020.
- Zhuang Liu, Mingjie Sun, Tinghui Zhou, Gao Huang, and Trevor Darrell. Rethinking the value of network pruning. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=rJlnB3C5Ym>.
- Christos Louizos, Karen Ullrich, and Max Welling. Bayesian compression for deep learning, 2017.
- Haiquan Lu, Yefan Zhou, Shiwei Liu, Zhangyang Wang, Michael W. Mahoney, and Yaoqing Yang. Alphapruning: Using heavy-tailed self regularization theory for improved layer-wise pruning of large language models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=fHq4x2YXVv>.
- Xinyin Ma, Gongfan Fang, and Xinchao Wang. LLM-pruner: On the structural pruning of large language models. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=J8Ajf9WfXP>.
- Mitchell Marcus, Grace Kim, Mary Ann Marcinkiewicz, Robert MacIntyre, Ann Bies, Mark Ferguson, Karen Katz, and Britta Schasberger. The Penn Treebank: Annotating predicate argument structure. In *Human Language Technology: Proceedings of a Workshop held at Plainsboro, New Jersey, March 8-11, 1994*, 1994. URL <https://aclanthology.org/H94-1020>.

- Xiang Meng, Wenyu Chen, Riade Benbaki, and Rahul Mazumder. FALCON: FLOP-aware combinatorial optimization for neural network pruning. In Sanjoy Dasgupta, Stephan Mandt, and Yingzhen Li (eds.), *Proceedings of The 27th International Conference on Artificial Intelligence and Statistics*, volume 238 of *Proceedings of Machine Learning Research*, pp. 4384–4392. PMLR, 02–04 May 2024a. URL <https://proceedings.mlr.press/v238/meng24a.html>.
- Xiang Meng, Shibal Ibrahim, Kayhan Behdin, Hussein Hazimeh, Natalia Ponomareva, and Rahul Mazumder. OSSCAR: One-shot structured pruning in vision and language models with combinatorial optimization. In *Forty-first International Conference on Machine Learning*, 2024b. URL <https://openreview.net/forum?id=ZctlF8R1V4>.
- Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models, 2016. URL <https://arxiv.org/abs/1609.07843>.
- Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct electricity? a new dataset for open book question answering. In Ellen Riloff, David Chiang, Julia Hockenmaier, and Jun'ichi Tsujii (eds.), *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 2381–2391, Brussels, Belgium, October–November 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-1260. URL <https://aclanthology.org/D18-1260>.
- Asit Mishra, Jorge Albericio Latorre, Jeff Pool, Darko Stosic, Dusan Stosic, Ganesh Venkatesh, Chong Yu, and Paulius Micekevicius. Accelerating sparse deep neural networks, 2021. URL <https://arxiv.org/abs/2104.08378>.
- Michael C Mozer and Paul Smolensky. Using relevance to reduce network size automatically. *Connection Science*, 1(1):3–16, 1989.
- Markus Nagel, Rana Ali Amjad, Mart Van Baalen, Christos Louizos, and Tijmen Blankevoort. Up or down? Adaptive rounding for post-training quantization. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 7197–7206. PMLR, 13–18 Jul 2020. URL <https://proceedings.mlr.press/v119/nagel20a.html>.
- NeuralMagic. DeepSparse. <https://github.com/neuralmagic/deepsparse>, 2021. Accessed: 2025-08-11.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL [https://proceedings.neurips.cc/paper\\_files/paper/2019/file/bdbca288fee7f92f2bfa9f7012727740-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2019/file/bdbca288fee7f92f2bfa9f7012727740-Paper.pdf).
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67, 2020. URL <http://jmlr.org/papers/v21/20-074.html>.
- Alex Renda, Jonathan Frankle, and Michael Carbin. Comparing rewinding and fine-tuning in neural network pruning, 2020.
- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: an adversarial winograd schema challenge at scale. *Commun. ACM*, 64(9):99–106, August 2021. ISSN 0001-0782. doi: 10.1145/3474381. URL <https://doi.org/10.1145/3474381>.
- Sidak Pal Singh and Dan Alistarh. Woodfisher: Efficient second-order approximation for neural network compression. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 18098–18109. Curran

- Associates, Inc., 2020a. URL [https://proceedings.neurips.cc/paper\\_files/paper/2020/file/d1ff1ec86b62cd5f3903ff19c3a326b2-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/d1ff1ec86b62cd5f3903ff19c3a326b2-Paper.pdf).
- Sidak Pal Singh and Dan Alistarh. Woodfisher: Efficient second-order approximation for neural network compression. *Advances in Neural Information Processing Systems*, 33:18098–18109, 2020b.
- Yang Sui, Miao Yin, Yi Xie, Huy Phan, Saman Aliari Zonouz, and Bo Yuan. CHIP: CHannel independence-based pruning for compact neural networks. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*, 2021. URL <https://openreview.net/forum?id=EmeWbcWORRg>.
- Mingjie Sun, Zhuang Liu, Anna Bair, and J. Zico Kolter. A simple and effective pruning approach for large language models, 2024. URL <https://arxiv.org/abs/2306.11695>.
- Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Herve Jegou. Training data-efficient image transformers and distillation through attention. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 10347–10357. PMLR, 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/touvron21a.html>.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrutu Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models, 2023. URL <https://arxiv.org/abs/2307.09288>.
- Chaoqi Wang, Guodong Zhang, and Roger Grosse. Picking winning tickets before training by preserving gradient flow, 2020.
- M.A. Woodbury and Princeton University. Department of Statistics. *Inverting Modified Matrices*. Memorandum Report / Statistical Research Group, Princeton. Department of Statistics, Princeton University, 1950. URL [https://books.google.com/books?id=\\_zAnzGECAAJ](https://books.google.com/books?id=_zAnzGECAAJ).
- Lu Yin, You Wu, Zhenyu Zhang, Cheng-Yu Hsieh, Yaqing Wang, Yiling Jia, Gen Li, AJAY KUMAR JAISWAL, Mykola Pechenizkiy, Yi Liang, Michael Bendersky, Zhangyang Wang, and Shiwei Liu. Outlier weighed layerwise sparsity (OWL): A missing secret sauce for pruning LLMs to high sparsity. In *Forty-first International Conference on Machine Learning*, 2024. URL <https://openreview.net/forum?id=ahEm312P6w>.
- Xin Yu, Thiago Serra, Srikumar Ramalingam, and Shandian Zhe. The combinatorial brain surgeon: Pruning weights that cancel one another in neural networks. In *International Conference on Machine Learning*, pp. 25668–25683. PMLR, 2022.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. HellaSwag: Can a machine really finish your sentence? In Anna Korhonen, David Traum, and Lluís Màrquez (eds.), *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 4791–4800, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1472. URL <https://aclanthology.org/P19-1472>.
- Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt Shuster, Daniel

Simig, Punit Singh Koura, Anjali Sridhar, Tianlu Wang, and Luke Zettlemoyer. Opt: Open pre-trained transformer language models, 2022. URL <https://arxiv.org/abs/2205.01068>.

Zhenyu Zhang, Xuxi Chen, Tianlong Chen, and Zhangyang Wang. Efficient lottery ticket finding: Less data is more. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 12380–12390. PMLR, 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/zhang21c.html>.

## A Appendix

### A.1 Related Work

Many techniques have been proposed in order to prune a neural network to a desired sparsity. While some methods emphasize the use of gradual pruning to recover accuracy (Han et al., 2015; Gale et al., 2019; Singh & Alistarh, 2020a; Blalock et al., 2020; Benbaki et al., 2023), others attempt to prune during training or at initialization (Louizos et al., 2017; Frankle & Carbin, 2019; Lee et al., 2019; Liu et al., 2019; Lee et al., 2020; Wang et al., 2020; Renda et al., 2020; Frankle et al., 2021; Sui et al., 2021; Zhang et al., 2021). While effective, these methods require extensive retraining and are often too costly or impractical in resource-constrained settings with large models. Therefore, we focus on post-training one-shot pruning of large models in this work.

**Post-training One-Shot Pruning.** In the post-training one-shot pruning literature, we identify three main types of approaches that are generally proposed: (i) Magnitude-based methods (Hanson & Pratt, 1988; Mozer & Smolensky, 1989; Gordon et al., 2020) use the weight magnitudes to determine their importances and whether or not they should be pruned. Since magnitude alone may not be the best proxy for weight relevance, alternatives have been proposed. (ii) Second-order approaches, such as OBD and OBS (Optimal Brain Damage/Surgeon) (LeCun et al., 1989a; Hassibi & Stork, 1992a), consider a local quadratic approximation of the loss around the pre-trained weights. These methods employ impact-based pruning, removing weights based on the estimated effect of their removal on the loss function. This line of work uses a second-order Taylor approximation of the training loss and the empirical Fisher information as a proxy for the Hessian. Singh & Alistarh (2020b) proposed a block-diagonal approximation of the empirical Fisher matrix for scaling the OBS framework to modern vision model sizes. Yu et al. (2022) propose to select weights to prune based on their joint rather than individual impact on the loss. (iii) Layer-wise pruning methods adapt the OBS framework to the layer-wise reconstruction objective. Dong et al. (2017) prunes each layer independently to overcome the computational challenge of computing the per-sample gradients needed in OBS. With the Optimal Brain Compression (OBC) framework, Frantar et al. (2022) adapts the OBS algorithm (Hassibi & Stork, 1992b) to the layer-wise reconstruction error and proposes rank-1 updates of the Hessian for efficient pruning.

**Pruning Vision Models.** While the pruning techniques mentioned previously are generally applicable, they have been investigated primarily in the context of Convolutional Neural Networks (CNNs). For instance, in the case of the ResNet architecture (He et al., 2015), OBC (Frantar et al., 2022) represents a state-of-the-art post-training one-shot pruning approach. Kuznedelev et al. (2023), with their Correlation Aware Pruner (CAP), adapted the greedy OBS algorithm used in OBC (Frantar et al., 2022) to Vision Transformers. This approach is a state-of-the-art method for post-training one-shot pruning in Vision Transformers.

**Pruning Large Language Models.** Pruning is particularly important for Large Language Models, which can have billions of parameters. In the case of unstructured and semi-structured pruning, state-of-the-art methods include SparseGPT (Frantar & Alistarh, 2023) and Wanda (Sun et al., 2024), both of which build on the OBC/OBS (Frantar et al., 2022; Hassibi & Stork, 1992b) framework. To make pruning more scalable, SparseGPT prunes the weight matrix by groups of columns. This enables to reduce the number of computations during pruning by considering only the Hessian for the weights within the support. However, SparseGPT still relies on a block-diagonal approximation of the Hessian. Wanda simplifies this further by using a diagonal matrix instead. Both methods focus solely on minimizing the layer-wise reconstruction loss. In the case of structured pruning, ZipLM (Kurtic et al., 2023) uses the layer-wise reconstruction error to guide the pruning of attention heads and projection layers. Building on a careful reformulation of this objective, Meng et al. (2024b) introduce OSSCAR, a more scalable and stronger baseline for one-shot

structured pruning. OSSCAR greedily selects columns for removal to reduce the reconstruction objective, and optionally applies a local-search refinement step. OSSCAR is a state-of-the-art one-shot structured pruning method for LLMs.

**Non-Uniform Sparsity.** LLMs typically suffer significant degradation beyond 50% uniform sparsity. This has motivated recent work on non-uniform sparsity for LLMs, which assigns different sparsity levels to different layers to better preserve model quality under the same global sparsity constraint. OWL (Yin et al., 2024) improves pruning by allocating layer-wise sparsity based on the distribution of outlier weights, leading to notable gains for methods like SparseGPT and Wanda. AlphaPruning (Lu et al., 2024) takes a more principled approach and uses heavy-tailed self-regularization theory to measure how well each layer is trained. After quantifying the heavy-tail distribution of the weights, it assigns lower sparsity to the better trained layers.

## A.2 MOONSHOT-SparseGPT Algorithm

We present below our adaptation of the SparseGPT algorithm (Frantar & Alistarh, 2023), denoted MOONSHOT-SparseGPT.

---

### Algorithm 2 MOONSHOT-SparseGPT

---

**Input:** Layer weight matrix  $\widehat{W}^{(l)} \in \mathbb{R}^{d_{\text{out}}^{(l)} \times d_{\text{in}}^{(l)}}$ , layer input matrix  $X^{(l)} \in \mathbb{R}^{d_{\text{in}}^{(l)} \times N}$ , per-sample gradients  $A_k^{(l)} \in \mathbb{R}^{d_{\text{in}}^{(l)} \times N}$  for each block  $k = 1, \dots, d_{\text{out}}^{(l)}$ , multi-objective weight  $\lambda \in [0, 1]$ , lazy batch-update block-size  $B$ , adaptive mask selection blocksize  $B_s$ , number of blocks to prune in parallel  $K_P$  and sparsity level  $p$ .

- 1: Initialize Pruned Weights:  $W^{(l)} \leftarrow \widehat{W}^{(l)}$
- 2: Initialize Binary Pruning Mask:  $M \leftarrow \mathbb{1}_{d_{\text{out}}^{(l)} \times d_{\text{in}}^{(l)}}$
- 3: Initialize Block Errors:  $E \leftarrow \mathbb{0}_{d_{\text{out}}^{(l)} \times B}$
- 4: **for**  $k_p = 0, K_P, 2K_P, \dots$  **do**
- 5:   Compute Hessian Inverse  $\{G_k^{(l)}\}_{k=k_p}^{k_p+K_P} = \{F_k^{(l)-1}\}_{k=k_p}^{k_p+K_P}$  using Algorithm 1 (tensor format  $G^{(l)} \in \mathbb{R}^{K_P \times d_{\text{in}}^{(l)} \times d_{\text{in}}^{(l)}}$ )
- 6:   Compute Cholesky Decomposition for each block  $G_k^{(l)} \leftarrow \text{Cholesky}(G_k^{(l)})^T$
- 7:   **for**  $i = 0, B, 2B, \dots$  **do**
- 8:     **for**  $j = i, \dots, i + B - 1$  **do**
- 9:       **if**  $j \bmod B_s = 0$  **then**
- 10:           $M_{k_p:k_p+K_P, j:(j+B_s)} \leftarrow$  mask of  $(1-p)\%$  weights  $w_c \in W_{k_p:k_p+K_P, j:(j+B_s)}^{(l)}$
- 11:          with largest  $w_c^2 / [G^{(l)}]_{k_p:k_p+K_P, c, c}^2$
- 12:          Compute Pruning Errors:  $E_{k_p:k_p+K_P, j-i} \leftarrow W_{k_p:k_p+K_P, j}^{(l)} / [G^{(l)}]_{k_p:k_p+K_P, j, j}$
- 13:          Freeze Weights:  $E_{k_p:k_p+K_P, i, j-i} \leftarrow (\mathbb{1}_{d_{\text{out}}^{(l)} \times d_{\text{in}}^{(l)}} - M_{k_p:k_p+K_P, j}) \cdot E_{k_p:k_p+K_P, j-i}$
- 14:          Weights update:  $W_{k_p:k_p+K_P, j:(i+B)}^{(l)} \leftarrow W_{k_p:k_p+K_P, j:(i+B)}^{(l)} - E_{k_p:k_p+K_P, j-i} \cdot G_{k_p:k_p+K_P, j, j:(i+B)}^{(l)}$
- 15:          Weights update:  $W_{k_p:k_p+K_P, (i+B)}^{(l)} \leftarrow W_{k_p:k_p+K_P, (i+B)}^{(l)} - E_{k_p:k_p+K_P, :} \cdot G_{k_p:k_p+K_P, i:(i+B), (i+B)}^{(l)}$
- 16:    Set pruned weights to 0:  $W_{k_p:k_p+K_P, :}^{(l)} \leftarrow W_{k_p:k_p+K_P, :}^{(l)} \cdot M_{k_p:k_p+K_P, :}$

**Output:** Pruned weights  $W^{(l)}$

---

A key component of SparseGPT is the use of the Hessian inverse. In the multi-objective setting, however, directly computing matrix inverses for every block is computationally infeasible. To overcome this, we replace Step 2 in Algorithm 2 with the more efficient procedure outlined in Algorithm 1. We also note that the larger Hessian size makes it challenging to prune all blocks simultaneously. To address this, we perform pruning in parallel over  $K_P$  blocks at a time. This means that the uniform sparsity budget is applied within each group of  $K_P$  blocks rather than across the entire weight matrix. While this introduces a more local form of sparsity allocation, choosing  $K_P$  sufficiently large (all the blocks for Llama-3.2-1B and at least 50% of them for Llama-3.2-3B) ensures that the effect is minimal in practice.

### A.3 Woodbury Update in 1

Let  $A \in \mathbb{R}^{n \times n}$ ,  $C \in \mathbb{R}^{k \times k}$ ,  $U \in \mathbb{R}^{n \times k}$  and  $V \in \mathbb{R}^{k \times n}$ . The Woodbury matrix identity (Woodbury & of Statistics, 1950) gives us that:

$$(A + UCV)^{-1} = A^{-1} - A^{-1}U(C + VA^{-1}U)^{-1}VA^{-1} \quad (12)$$

In our case, we want to compute  $G_k^{(l)} = (F_k^{(l)})^{-1} = \left( \frac{\lambda}{\mathcal{L}_R^{(l)}(0)} L_k^{(l)} + \frac{1-\lambda}{\mathcal{L}_F^{(l)}(0)} H_k^{(l)} \right)^{-1}$

We focus in Algorithm 1 on methods like SparseGPT and OSSCAR, which use the layer-wise reconstruction error objective to scale to billion-parameter-LLMs (without further diagonal approximation on the exact block-diagonal Hessian). For these pruning baselines,  $L_1^{(l)} = \dots = L_K^{(l)} = X^{(l)}(X^{(l)})^T$  (with  $K = d_{\text{out}}^{(l)}$ ). In addition, as seen in Section 2.3, we can write  $H_k^{(l)} = \frac{1}{N} A_k^{(l)} A_k^{(l)T}$ . Therefore,  $G_k^{(l)}$  can be rewritten as:

$$G_k^{(l)} = \left( \frac{\lambda}{\mathcal{L}_R^{(l)}(0)} X^{(l)}(X^{(l)})^T + \frac{1-\lambda}{N\mathcal{L}_F^{(l)}(0)} A_k^{(l)} A_k^{(l)T} \right)^{-1} \quad (13)$$

This is the same form as equation 12 with:

$$A = \frac{\lambda}{\mathcal{L}_R^{(l)}(0)} X^{(l)}(X^{(l)})^T = J_0^{-1}, \quad U = \frac{1-\lambda}{N\mathcal{L}_F^{(l)}(0)} A_k^{(l)}, \quad C = I_N, \quad V = (A_k^{(l)})^T$$

Therefore, using equation 12, equation 13 becomes:

$$\begin{aligned} G_k^{(l)} &= \left( \frac{\lambda}{\mathcal{L}_R^{(l)}(0)} X^{(l)}(X^{(l)})^T + \frac{1-\lambda}{N\mathcal{L}_F^{(l)}(0)} A_k^{(l)} A_k^{(l)T} \right)^{-1} \\ &= J_0 - J_0 \left( \frac{1-\lambda}{N\mathcal{L}_F^{(l)}(0)} A_k^{(l)} \right) \left( I_N + A_k^{(l)T} J_0 \left( \frac{1-\lambda}{N\mathcal{L}_F^{(l)}(0)} A_k^{(l)} \right) \right)^{-1} A_k^{(l)T} J_0 \end{aligned}$$

This expression of  $G_k^{(l)}$  is the same as the one used in Algorithm 1.

### A.4 Efficient backsolve for OSSCAR

After greedy selection, OSSCAR performs a backsolve to update the weights on the support of remaining columns  $S$ . With the original single-objective, the Hessian has repeated blocks, yielding the closed form

$$W_{S,:}^* = [XX^\top]_{S,S}^{-1} [XX^\top]_{S,:} \widehat{W}.$$

With *MOONSHOT* (when  $\lambda \neq 1$ ), the Hessian becomes block-diagonal with row-dependent blocks,  $H = \text{Diag}(F_1, \dots, F_K)$ , where each  $F_k$  has the same shape as  $XX^\top$ . A direct extension would require inverting all  $K$  matrices, which is unnecessarily expensive. Instead, we compute the Cholesky decomposition for each block and use an efficient solver in PyTorch for the system:

$$\text{Diag}([F_1]_{S,S}, \dots, [F_K]_{S,S}) \text{vec}(W_{S,:}^*) = H_{S,:} \text{vec}(\widehat{W}),$$

where  $\text{vec}(W)$  denotes the vector form of the weight matrix  $W$ .

### A.5 Pruning Time

The pruning times obtained when applying *MOONSHOT* on the different baselines are described in Table 4. OBC and CAP allow pruning at multiple sparsity levels in a single run, with minimal additional cost compared to pruning at a single level. For these methods, we report the total pruning time needed to obtain

Table 4: Pruning times (in seconds) of *MOONSHOT* over 3 seeds, for different values of  $\lambda$ , architectures, and pruning baselines.

		OBC	CAP			Wanda		SparseGPT	
Sparsity	$\lambda$	ResNet-50	DeiT-Tiny	DeiT-Small	DeiT-Base	Llama-3.2-1B	Llama-3.2-3B	Llama-3.2-1B	Llama-3.2-3B
Unstructured	0.00	7360.65 $\pm$ 8.35	71.89 $\pm$ 7.91	268.58 $\pm$ 23.98	1303.13 $\pm$ 22.99	69.59 $\pm$ 0.41	403.51 $\pm$ 37.9	473.03 $\pm$ 1.95	2181.89 $\pm$ 128.72
	0.25	7392.68 $\pm$ 23.96	65.63 $\pm$ 1.12	240.06 $\pm$ 1.73	1186.93 $\pm$ 9.11	68.94 $\pm$ 0.47	407.2 $\pm$ 38.62	474.86 $\pm$ 1.48	2269.66 $\pm$ 94.02
	0.50	7370.94 $\pm$ 1.0	66.77 $\pm$ 1.51	226.0 $\pm$ 6.68	1185.66 $\pm$ 14.01	69.63 $\pm$ 0.64	464.06 $\pm$ 48.64	472.11 $\pm$ 1.07	2287.58 $\pm$ 88.63
	0.75	7389.37 $\pm$ 9.43	67.58 $\pm$ 2.02	227.61 $\pm$ 0.77	1257.79 $\pm$ 87.92	70.23 $\pm$ 0.83	420.58 $\pm$ 47.86	471.77 $\pm$ 0.19	2165.25 $\pm$ 171.46
	1.00	<b>7146.09</b> $\pm$ 5.34	<b>21.62</b> $\pm$ 1.15	<b>43.07</b> $\pm$ 0.89	<b>247.25</b> $\pm$ 3.04	<b>21.0</b> $\pm$ 0.36	<b>70.66</b> $\pm$ 7.23	<b>73.9</b> $\pm$ 0.29	<b>307.73</b> $\pm$ 11.31
Semi-Structured (2:4)	0.00	3965.42 $\pm$ 13.45	60.19 $\pm$ 0.9	229.22 $\pm$ 2.39	1236.61 $\pm$ 8.72	71.0 $\pm$ 0.54	413.24 $\pm$ 53.15	478.6 $\pm$ 1.14	2197.32 $\pm$ 142.89
	0.25	4033.31 $\pm$ 85.1	62.06 $\pm$ 1.11	217.56 $\pm$ 1.78	1123.35 $\pm$ 13.82	71.62 $\pm$ 0.39	442.05 $\pm$ 35.45	477.58 $\pm$ 0.25	2345.26 $\pm$ 86.07
	0.50	3956.29 $\pm$ 4.2	62.35 $\pm$ 1.41	218.5 $\pm$ 2.5	1110.06 $\pm$ 9.58	72.62 $\pm$ 1.04	474.96 $\pm$ 6.32	476.66 $\pm$ 2.13	2362.17 $\pm$ 59.99
	0.75	3974.67 $\pm$ 5.72	65.18 $\pm$ 4.0	213.7 $\pm$ 2.33	1102.55 $\pm$ 3.82	73.87 $\pm$ 0.33	433.84 $\pm$ 77.23	477.92 $\pm$ 0.51	2237.57 $\pm$ 216.16
	1.00	<b>3737.4</b> $\pm$ 5.68	<b>16.31</b> $\pm$ 1.21	<b>29.38</b> $\pm$ 0.1	<b>178.98</b> $\pm$ 2.99	<b>22.74</b> $\pm$ 0.11	<b>80.78</b> $\pm$ 6.02	<b>79.01</b> $\pm$ 0.13	<b>333.86</b> $\pm$ 13.01

pruned weights for all target sparsities: 0.5, 0.6, 0.7, 0.8, and 0.9. SparseGPT and Wanda only support pruning one sparsity level at a time, but their runtime is not affected by the chosen sparsity. We therefore report their pruning time at sparsity 0.6.

We note that the baselines correspond to  $\lambda = 0$  for CAP and  $\lambda = 1$  for OBC, Wanda and SparseGPT. We observe that *MOONSHOT* incur no to almost no additional computational overhead to CAP and OBC for the DeiT models and ResNet-50 respectively. With *MOONSHOT*-SparseGPT, we achieve pruning times of under 40 minutes for Llama-3.2-3B and under 8 minutes for Llama-3.2-1B. While this is a slowdown compared to SparseGPT’s original pruning times of 5 and 1.2 minutes respectively, the increase is reasonable given the significant performance improvements. Moreover, since the ultimate goal is to deploy a more compact and efficient model post-pruning, the slightly longer pruning time, viewed as a one-time cost, is a worthwhile trade-off. For Wanda, pruning times are similarly manageable, with 8 minutes for Llama-3.2-3B and 1.2 minutes for Llama-3.2-1B, compared to 1.2 minutes and 20 seconds respectively for the single-objective version.

## A.6 Additional Hyperparameters

**Dampening Term.** As noted in previous work (Frantar et al., 2022; Frantar & Alistarh, 2023; Kuznedelev et al., 2023), the Hessian is not always invertible in practice. To address this issue, a common approach is to add a dampening factor  $\mu$  to the diagonal of the Hessian, ensuring it is positive definite. However, selecting  $\mu$  is non-trivial: if  $\mu$  is too small, numerical instabilities can persist, while a large  $\mu$  may degrade algorithm performance. Following SparseGPT (Frantar & Alistarh, 2023), we set  $\mu$  to 1% of the mean of the diagonal elements of the Hessian. For OBC (Frantar et al., 2022) and CAP (Kuznedelev et al., 2023), we use the diagonal of each block  $F_k^{(l)}$  in the multi-objective formulation equation 10, while for SparseGPT, we use the diagonal of  $X^{(l)}(X^{(l)})^T$  only to maintain the efficiency of the inverse Hessian computation described in Section 2.3. For Wanda (Sun et al., 2024) only,  $\mu$  is set to 0 as recommended by the authors (no inverse Hessian is required).

**OWL and AlphaPruning.** Following the optimal parameters found by the authors of OWL (Yin et al., 2024) on Llama-7B, we use  $\lambda^{(\text{OWL})} = 0.08$  and  $M = 5$ . For AlphaPruning, we fixed  $\tau = 0.05$  for Llama-3.2-1B and  $\tau = 0.1$  for Llama-3.2-3B.

## A.7 Hessian Recomputation

In this section, we report additional results with Hessian recomputation after each block of layers for SparseGPT and Wanda. Following prior Fisher-based work (Benbaki et al., 2023), we denote these variants as *MOONSHOT*-SparseGPT++ and *MOONSHOT*-Wanda++. Although more costly than computing the Hessian once, our implementation remains tractable: we compute per-sample gradients for one block of layers at a time and update the dataset after pruning each block to reduce further gradient costs (using only the remaining dense layers). With this approach, *MOONSHOT*-SparseGPT++ prunes Llama-3.2-1B in under 11 minutes (vs. under 2 minutes for standard SparseGPT), and *MOONSHOT*-Wanda++ prunes Llama-3.2-1B in under 4 minutes (vs. around 1 minute for standard Wanda) on a single L40 GPU (40 GB). Comprehensive results are presented in Tables 5 and 6.

Table 5: Test perplexity on C4, WikiText2 and PTB and zero-shot accuracy of Llama-3.2-1B using *MOONSHOT-SparseGPT++*. Zero-shot **mean performance** and **win rate** are computed over the 7 classification tasks. Results are averaged over 3 seeds with standard errors.

Sparsity	$\lambda$	C4 ↓	WikiText2 ↓	PTB ↓	BoolQ ↑	HellaSwag ↑	WinoGrande ↑	ARC-e ↑	ARC-c ↑	PIQA ↑	OBQA ↑	Mean ↑	Win Rate ↑
Dense	-	14.02	9.75	17.59	64.01	47.73	60.14	65.15	31.23	74.32	26.4	52.71	-
0.5	0.0	27.37 ± 0.02	20.3 ± 0.03	33.19 ± 0.31	62.25 ± 0.19	37.0 ± 0.15	55.01 ± 0.88	53.47 ± 0.4	24.43 ± 0.8	67.43 ± 0.08	18.0 ± 0.12	45.37 ± 0.14	9.52 ± 9.52
	0.1	24.02 ± 0.1	17.71 ± 0.11	28.99 ± 0.26	62.5 ± 0.26	38.37 ± 0.03	55.3 ± 0.37	56.17 ± 0.3	25.63 ± 0.19	<b>68.57</b> ± 0.19	19.2 ± 0.12	46.53 ± 0.04	28.57 ± 8.25
	0.25	23.8 ± 0.12	17.57 ± 0.01	28.86 ± 0.02	<b>63.03</b> ± 0.11	38.71 ± 0.03	54.67 ± 0.37	55.2 ± 0.26	26.65 ± 0.08	68.17 ± 0.17	20.33 ± 0.44	46.68 ± 0.02	33.33 ± 12.6
	0.5	23.6 ± 0.06	17.38 ± 0.02	28.72 ± 0.12	62.8 ± 0.46	38.91 ± 0.05	55.14 ± 0.16	55.61 ± 0.52	25.97 ± 0.33	68.41 ± 0.22	19.73 ± 0.57	46.65 ± 0.07	23.81 ± 12.6
	0.75	23.45 ± 0.06	<b>17.23</b> ± 0.03	28.43 ± 0.11	62.94 ± 0.35	38.93 ± 0.05	55.96 ± 0.34	55.56 ± 0.06	<b>27.05</b> ± 0.34	67.85 ± 0.19	19.8 ± 0.46	46.87 ± 0.07	42.86 ± 8.25
	1.0	<b>23.39</b> ± 0.05	<b>17.23</b> ± 0.02	<b>28.16</b> ± 0.16	62.4 ± 0.14	39.11 ± 0.12	55.72 ± 0.39	<b>56.44</b> ± 0.26	26.22 ± 0.16	68.55 ± 0.42	20.73 ± 0.93	47.03 ± 0.1	<b>47.62</b> ± 26.51
0.6	0.0	59.74 ± 0.72	48.97 ± 0.26	72.66 ± 3.54	60.47 ± 0.82	30.38 ± 0.19	51.43 ± 0.62	41.68 ± 0.58	20.73 ± 0.57	60.57 ± 0.13	14.13 ± 0.48	39.92 ± 0.34	4.76 ± 4.76
	0.1	47.21 ± 0.16	37.08 ± 0.66	53.82 ± 0.38	62.1 ± 0.1	31.97 ± 0.12	52.22 ± 0.25	45.86 ± 0.98	21.56 ± 0.54	63.02 ± 0.09	15.6 ± 0.6	41.76 ± 0.25	52.38 ± 9.52
	0.25	46.93 ± 0.41	36.73 ± 0.2	52.18 ± 0.65	61.89 ± 0.27	32.34 ± 0.11	52.04 ± 0.13	45.9 ± 0.62	21.9 ± 0.5	63.44 ± 0.65	15.73 ± 0.29	41.89 ± 0.09	<b>52.38</b> ± 4.76
	0.5	46.31 ± 0.42	35.07 ± 0.29	50.84 ± 0.64	62.16 ± 0.15	32.5 ± 0.13	<b>53.01</b> ± 0.55	45.66 ± 0.71	21.53 ± 0.2	<b>63.89</b> ± 0.13	16.47 ± 0.66	42.17 ± 0.3	<b>71.43</b> ± 14.29
	0.75	46.19 ± 0.4	<b>34.85</b> ± 0.78	51.53 ± 1.17	62.05 ± 0.07	<b>32.68</b> ± 0.11	52.64 ± 0.32	<b>46.65</b> ± 0.38	<b>21.96</b> ± 0.33	63.51 ± 0.53	16.2 ± 0.35	<b>42.24</b> ± 0.13	57.14 ± 8.25
	1.0	<b>46.08</b> ± 0.79	35.6 ± 0.44	<b>50.36</b> ± 1.23	61.67 ± 0.35	32.55 ± 0.16	52.67 ± 0.64	46.34 ± 0.55	21.53 ± 0.52	63.46 ± 0.64	16.87 ± 0.18	42.16 ± 0.14	66.67 ± 12.6
0.7	0.0	168.52 ± 1.54	156.79 ± 3.73	190.61 ± 10.89	54.75 ± 3.4	27.18 ± 0.16	50.07 ± 0.16	32.72 ± 0.72	18.66 ± 0.59	56.4 ± 0.07	11.73 ± 0.33	35.93 ± 0.52	28.57 ± 0.0
	0.1	130.06 ± 1.93	107.57 ± 2.09	135.44 ± 4.14	61.88 ± 0.1	27.76 ± 0.06	<b>50.83</b> ± 0.46	34.09 ± 0.18	18.69 ± 0.6	57.13 ± 0.2	11.53 ± 0.37	37.41 ± 0.2	47.62 ± 17.17
	0.25	127.78 ± 1.31	103.16 ± 1.48	<b>125.71</b> ± 4.58	61.47 ± 0.26	27.78 ± 0.1	49.3 ± 0.46	<b>34.83</b> ± 0.42	17.55 ± 0.51	56.71 ± 0.15	12.2 ± 0.35	37.12 ± 0.08	23.81 ± 4.76
	0.5	<b>127.22</b> ± 1.21	<b>101.14</b> ± 0.6	130.32 ± 3.42	60.84 ± 0.58	27.73 ± 0.06	49.46 ± 0.55	34.39 ± 0.68	18.43 ± 0.31	<b>57.25</b> ± 0.33	11.6 ± 0.6	37.1 ± 0.1	38.1 ± 12.6
	0.75	132.13 ± 0.92	103.3 ± 2.55	130.82 ± 5.53	<b>62.04</b> ± 0.12	<b>27.94</b> ± 0.09	49.8 ± 0.27	34.5 ± 0.07	<b>18.8</b> ± 0.37	56.64 ± 0.35	11.47 ± 0.37	37.31 ± 0.13	<b>47.62</b> ± 9.52
	1.0	135.42 ± 1.28	105.49 ± 1.42	134.09 ± 2.7	61.31 ± 0.59	27.92 ± 0.06	50.78 ± 0.57	34.57 ± 0.19	18.52 ± 0.21	57.18 ± 0.23	11.27 ± 0.27	37.36 ± 0.12	<b>52.14</b> ± 16.5
2.4	0.0	159.35 ± 11.39	143.16 ± 4.52	174.23 ± 10.38	61.44 ± 0.45	27.86 ± 0.04	<b>50.83</b> ± 0.41	34.53 ± 0.68	18.71 ± 0.47	56.0 ± 0.41	<b>13.67</b> ± 0.18	<b>37.58</b> ± 0.08	-
	0.1	53.47 ± 0.38	40.62 ± 0.39	64.24 ± 0.82	62.09 ± 0.22	30.43 ± 0.1	52.59 ± 0.59	42.61 ± 0.2	19.94 ± 0.28	60.75 ± 0.35	14.2 ± 0.7	40.37 ± 0.04	19.05 ± 4.76
	0.1	43.81 ± 0.17	32.47 ± 0.2	51.62 ± 0.17	62.16 ± 0.1	31.67 ± 0.11	53.14 ± 0.55	44.39 ± 0.87	20.16 ± 0.38	62.3 ± 0.32	14.93 ± 0.48	41.25 ± 0.2	23.81 ± 4.76
	0.25	43.39 ± 0.09	<b>32.08</b> ± 0.31	51.18 ± 0.48	<b>62.21</b> ± 0.12	31.72 ± 0.03	53.62 ± 0.89	46.04 ± 0.17	19.94 ± 0.33	62.01 ± 0.09	15.2 ± 0.95	41.53 ± 0.09	<b>47.62</b> ± 4.76
	0.5	<b>43.07</b> ± 0.12	32.2 ± 0.48	50.95 ± 0.2	62.11 ± 0.12	31.83 ± 0.05	53.88 ± 0.43	46.07 ± 0.48	21.02 ± 0.35	61.95 ± 0.17	15.13 ± 0.18	41.71 ± 0.11	<b>47.62</b> ± 9.52
	0.75	43.19 ± 0.22	32.17 ± 0.07	50.55 ± 0.57	61.95 ± 0.14	31.86 ± 0.06	53.67 ± 0.12	45.62 ± 0.61	20.56 ± 0.3	<b>62.42</b> ± 0.17	14.67 ± 0.57	41.54 ± 0.11	42.86 ± 14.29
0.9	0.0	43.27 ± 0.1	<b>32.08</b> ± 0.29	<b>50.51</b> ± 0.65	61.75 ± 0.48	31.8 ± 0.08	54.72 ± 0.65	<b>46.34</b> ± 0.17	20.42 ± 0.58	61.64 ± 0.21	14.87 ± 0.44	41.65 ± 0.07	33.33 ± 4.76
	0.1	45.73 ± 0.77	34.12 ± 0.55	52.81 ± 0.61	61.64 ± 0.39	<b>32.18</b> ± 0.19	<b>54.93</b> ± 0.57	45.74 ± 0.91	<b>21.5</b> ± 0.69	61.81 ± 0.17	<b>16.2</b> ± 0.92	<b>42.0</b> ± 0.4	-

Table 6: Test perplexity on C4, WikiText2 and PTB and zero-shot accuracy of Llama-3.2-1B using *MOONSHOT-Wanda++*. Zero-shot **mean performance** and **win rate** are computed over the 7 classification tasks. Results are averaged over 3 seeds with standard errors.

Sparsity	$\lambda$	C4 ↓	WikiText2 ↓	PTB ↓	BoolQ ↑	HellaSwag ↑	WinoGrande ↑	ARC-e ↑	ARC-c ↑	PIQA ↑	OBQA ↑	Mean ↑	Win Rate ↑
Dense	-	14.02	9.75	17.59	64.01	47.73	60.14	65.15	31.23	74.32	26.4	52.71	-
0.5	0.0	31.13 ± 0.14	21.86 ± 0.1	38.27 ± 0.35	61.99 ± 0.08	35.08 ± 0.06	54.35 ± 0.09	54.45 ± 0.16	24.09 ± 0.38	66.09 ± 0.23	16.6 ± 0.23	44.66 ± 0.1	66.67 ± 4.76
	0.1	30.56 ± 0.26	21.46 ± 0.15	37.5 ± 0.53	62.12 ± 0.12	35.3 ± 0.07	54.09 ± 0.41	54.91 ± 0.43	24.29 ± 0.14	66.74 ± 0.13	16.47 ± 0.24	44.84 ± 0.05	66.67 ± 4.76
	0.25	30.37 ± 0.09	21.13 ± 0.06	37.4 ± 0.32	61.96 ± 0.12	35.37 ± 0.01	54.43 ± 0.53	<b>54.97</b> ± 0.46	24.6 ± 0.25	<b>66.79</b> ± 0.19	17.53 ± 0.47	45.09 ± 0.07	80.95 ± 4.76
	0.5	30.08 ± 0.21	20.83 ± 0.06	36.74 ± 0.46	62.08 ± 0.12	35.59 ± 0.02	<b>54.93</b> ± 0.41	54.77 ± 0.27	25.06 ± 0.2	66.36 ± 0.23	<b>17.67</b> ± 0.29	45.21 ± 0.07	<b>95.24</b> ± 4.76
	0.75	29.91 ± 0.26	20.61 ± 0.17	<b>36.63</b> ± 0.45	62.13 ± 0.14	35.73 ± 0.07	54.62 ± 0.43	54.45 ± 0.19	<b>25.26</b> ± 0.13	66.61 ± 0.13	17.6 ± 0.0	45.2 ± 0.09	90.48 ± 9.52
	1.0	<b>29.78</b> ± 0.2	<b>20.54</b> ± 0.1	36.69 ± 0.3	<b>62.22</b> ± 0.06	<b>35.85</b> ± 0.03	54.91 ± 0.33	54.95 ± 0.2	24.94 ± 0.15	66.63 ± 0.13	17.13 ± 0.57	<b>45.23</b> ± 0.09	90.48 ± 9.52
0.6	0.0	93.22 ± 0.76	65.58 ± 0.41	91.86 ± 1.37	61.99 ± 0.1	28.55 ± 0.06	51.07 ± 0.43	39.63 ± 0.04	18.83 ± 0.52	58.96 ± 0.38	12.87 ± 0.41	38.84 ± 0.13	76.19 ± 9.52
	0.1	90.14 ± 1.95	62.95 ± 0.88	90.3 ± 0.91	62.16 ± 0.04	28.75 ± 0.08	50.99 ± 0.25	40.28 ± 0.21	19.37 ± 0.18	59.7 ± 0.5	12.6 ± 0.2	39.12 ± 0.1	85.71 ± 8.25
	0.25	87.37 ± 1.18	<b>62.97</b> ± 0.47	87.48 ± 1.14	<b>62.17</b> ± 0.06	28.85 ± 0.12	51.51 ± 0.89	40.17 ± 0.18	18.94 ± 0.05	59.5 ± 0.18	12.73 ± 0.13	39.13 ± 0.14	80.95 ± 4.76
	0.5	<b>87.1</b> ± 1.68	62.86 ± 0.3	88.14 ± 1.31	62.08 ± 0.06	28.89 ± 0.06	51.38 ± 0.3	39.76 ± 0.22	19.28 ± 0.05	59.9 ± 0.11	13.47 ± 0.29	39.25 ± 0.06	<b>90.48</b> ± 4.76
	0.75	87.58 ± 0.79	62.94 ± 1.14	<b>85.97</b> ± 0.77	62.06 ± 0.13	29.0 ± 0.05	<b>51.67</b> ± 0.07	40.01 ± 0.2	19.25 ± 0.25	59.7 ± 0.12	13.2 ± 0.2	39.27 ± 0.01	<b>90.48</b> ± 4.76
	1.0	89.06 ± 0.44	63.87 ± 0.25	89.52 ± 0.28	61.62 ± 0.34	<b>29.13</b> ± 0.03	50.41 ± 0.21	<b>40.71</b> ± 0.07	<b>19.54</b> ± 0.1	<b>60.08</b> ± 0.14	<b>13.6</b> ± 0.12	<b>39.3</b> ± 0.11	<b>90.48</b> ± 4.76
0.7	0.0	109.62 ± 1.07	79.05 ± 0.37	107.92 ± 1.11	58.75 ± 1.04	28.43 ± 0.07	50.49 ± 0.48	37.61 ± 0.24	19.31 ± 0.19	58.03 ± 0.13	12.73 ± 0.18	37.91 ± 0.2	-
	0.1	382.23 ± 8.94	274.94 ± 10.62	298.74 ± 11.35	42.25 ± 0.96	26.9 ± 0.04	48.67 ± 0.56	<b>30.35</b> ± 0.36	18.43 ± 0.3	54.95 ± 0.13	11.4 ± 0.12	33.28 ± 0.05	57.14 ± 8.25
	0.1	352.6 ± 7.68	254.96 ± 2.39	274.48 ± 14.3	39.62 ± 0.51	26.9 ± 0.01	49.2 ± 0.33	30.29 ± 0.2	18.32 ± 0.27	55.11 ± 0.14	<b>12.13</b> ± 0.18	33.08 ± 0.18	57.14 ± 8.25
	0.25	368.59 ± 38.77	251.25 ± 22.88	277.37 ± 16.15	<b>45.3</b> ± 3.24	<b>27.02</b> ± 0.02	49.78 ± 0.29	29.97 ± 0.32	19.08 ± 0.12	54.9 ± 0.22	11.93 ± 0.27	<b>34.0</b> ± 0.43	52.38 ± 4.76
	0.5	<b>349.39</b> ± 17.94	<b>246.09</b> ± 15.9	<b>273.31</b> ± 18.89	39.83 ± 0.55	26.89 ± 0.06	<b>50.28</b> ± 0.66	29.66 ± 0.48	17.95 ± 0.16	54.82 ± 0.22	11.8 ± 0.4	33.03 ± 0.23	<b>33.33</b> ± 9.52
	0.75	368.59 ± 5.73	252.62 ± 7.82	292.4 ± 7.11	44.65 ± 2.05	26.82 ± 0.01	49.91 ± 0.26	29.76 ± 0.27	18.43 ± 0.27	55.22 ± 0.24	12.0 ± 0.31	33.83 ± 0.35	<b>66.67</b> ± 12.6
2.4	0.0	352.34 ± 3.1	249.04 ± 2.44	287.06 ± 4.69	40.15 ± 0.93	26.84 ± 0.04	48.88 ± 0.16	29.94 ± 0.26	18.12 ± 0.32	<b>55.24</b> ± 0.02	11.8 ± 0.64	33.0 ± 0.23	52.38 ± 9.52
	0.1	467.13 ± 14.94	436.87 ± 28.98	507.38 ± 37.85	40.84 ± 1.23	26.45 ± 0.1	<b>50.28</b> ± 0.37	29.31 ± 0.12	<b>19.28</b> ± 0.44	55.17 ± 0.3	<b>12.13</b> ± 0.44	33.35 ± 0.09	-
	0.1	96.63 ± 1.95	69.45 ± 0.97	97.33 ± 1.97	61.39 ± 0.43	28.54 ± 0.06	48.7 ± 0.43	38.43 ± 0.41	<b>18.89</b> ± 0.57	59.48 ± 0.13	12.47 ± 0.13	38.27 ± 0.03	71.43 ± 0.0
	0.1	96.35 ± 1.78	69.13 ± 0.85	98.94 ± 1.31	61.51 ± 0.04	28.56 ± 0.12	51.3 ± 0.16	39.13 ± 0.26	18.46 ± 0.5	59.39 ± 0.16	12.0 ± 0.42	38.62 ± 0.08	<b>76.19</b> ± 9.52
	0.25	<b>94.01</b> ± 1.72	<b>65.79</b> ± 0.34	<b>93.54</b> ± 0.96	62.17 ± 0.02	28.56 ± 0.09	<b>51.41</b> ± 0.28	39.52 ± 0.12	18.63 ± 0.19	59.07 ± 0.38	12.27 ± 0.07	38.8 ± 0.04	<b>76.19</b> ± 9.52
	0.5	96.97 ± 2.48	67.35 ± 1.33	99.48 ± 4.12	60.93 ± 0.85	28.54 ± 0.13	50.3 ± 0.93	<b>39.6</b> ± 0.62	18.34 ± 0.1	<b>59.54</b> ± 0.24	12.87 ± 0.52		

terms of computation time and performance gains. As mentioned in Section 3,  $K_p$  (see Algorithm 2) needs to be reduced in order to prune the larger projection layers. In particular, we use the following values of  $K_p$ :

Table 7: Number  $K_p$  of rows (Algorithm 2) pruned at the same time for *MOONSHOT*-SparseGPT ( $\lambda \neq 1$ ).

Layer	Llama-3.2-1B		Llama-3.2-3B	
	$K_p$	$n_{\text{blocks}}$	$K_p$	$n_{\text{blocks}}$
q_proj	2048	2048	1536	3072
k_proj	512	512	512	512
v_proj	512	512	512	512
o_proj	2048	2048	1536	3072
gate_proj	1024	8192	1024	8192
up_proj	1024	8192	1024	8192
down_proj	64	2048	192	3072

Llama-3.2-1B is pruned using a single L40 (40GB) and Llama-3.2-3B using a single A100 (80GB). Empirically, applying *MOONSHOT* to projection layers too (`gate_proj`, `up_proj`, `down_proj`) can yield additional gains:

Table 8: Impact of *MOONSHOT* on SparseGPT/Wanda on the LLaMA-3.2 models at 60% unstructured sparsity. The perplexities on C4, WikiText2 and PTB, as well as the zero-shot accuracies are averaged over 3 seeds with standard errors. Mean performance and win rate are computed over the 7 zero-shot downstream classification tasks.

(a) Llama-3.2-1B

Sparsity	Method	C4 ↓	WikiText2 ↓	PTB ↓	BoolQ ↑	HellaSwag ↑	WinoGrande ↑	ARC-e ↑	ARC-c ↑	PIQA ↑	OBQA ↑	Mean ↑	Win Rate ↑
0.6	SparseGPT	63.63 ± 1.18	54.60 ± 1.00	81.11 ± 3.99	60.07 ± 0.59	32.16 ± 0.20	54.46 ± 0.53	44.94 ± 0.11	21.47 ± 0.48	62.21 ± 0.20	17.07 ± 0.41	41.85 ± 0.20	-
	<i>MOONSHOT</i> -SparseGPT	50.28 ± 1.99	39.13 ± 1.54	60.14 ± 2.90	<b>62.36 ± 0.12</b>	32.49 ± 0.13	53.09 ± 0.18	46.49 ± 0.38	21.30 ± 0.23	63.22 ± 0.17	15.73 ± 0.55	42.10 ± 0.11	57.14 ± 8.25
	<i>MOONSHOT</i> (all)-SparseGPT	<b>42.96 ± 0.27</b>	<b>34.29 ± 0.46</b>	<b>54.92 ± 1.22</b>	62.09 ± 0.13	<b>32.88 ± 0.08</b>	53.14 ± 0.78	<b>46.87 ± 0.14</b>	<b>21.50 ± 0.57</b>	<b>63.73 ± 0.43</b>	16.93 ± 0.44	<b>42.45 ± 0.22</b>	<b>71.43 ± 8.25</b>
0.6	Wanda	117.71 ± 0.87	84.73 ± 0.73	119.64 ± 1.00	58.96 ± 1.39	28.86 ± 0.03	51.35 ± 0.49	38.82 ± 0.32	18.94 ± 0.26	59.05 ± 0.18	<b>13.93 ± 0.24</b>	38.56 ± 0.12	-
	<i>MOONSHOT</i> -Wanda	86.55 ± 1.67	63.57 ± 1.61	98.44 ± 3.98	61.56 ± 0.29	29.53 ± 0.06	51.64 ± 0.23	40.40 ± 0.17	<b>19.60 ± 0.06</b>	61.12 ± 0.08	13.40 ± 0.20	39.61 ± 0.07	80.95 ± 4.76
	<i>MOONSHOT</i> (all)-Wanda	<b>85.66 ± 0.92</b>	<b>63.12 ± 0.29</b>	<b>94.39 ± 1.78</b>	<b>61.81 ± 0.18</b>	<b>29.62 ± 0.05</b>	<b>51.88 ± 0.38</b>	<b>41.27 ± 0.22</b>	19.11 ± 0.13	<b>61.19 ± 0.17</b>	12.73 ± 0.47	<b>39.66 ± 0.06</b>	<b>85.71 ± 8.25</b>

(b) Llama-3.2-3B

Sparsity	Method	C4 ↓	WikiText2 ↓	PTB ↓	BoolQ ↑	HellaSwag ↑	WinoGrande ↑	ARC-e ↑	ARC-c ↑	PIQA ↑	OBQA ↑	Mean ↑	Win Rate ↑
0.6	SparseGPT	33.63 ± 0.14	26.12 ± 0.23	42.69 ± 0.73	66.82 ± 0.60	38.14 ± 0.14	60.91 ± 0.65	53.89 ± 0.11	26.28 ± 0.18	67.75 ± 0.34	18.47 ± 0.44	47.47 ± 0.08	-
	<i>MOONSHOT</i> -SparseGPT	28.23 ± 0.11	22.46 ± 0.17	35.63 ± 0.68	<b>67.76 ± 0.35</b>	39.13 ± 0.07	61.01 ± 0.23	57.59 ± 0.92	<b>27.79 ± 0.71</b>	69.44 ± 0.13	<b>20.00 ± 0.53</b>	<b>48.96 ± 0.12</b>	<b>95.24 ± 4.76</b>
	<i>MOONSHOT</i> (all)-SparseGPT	<b>26.26 ± 0.12</b>	<b>20.67 ± 0.16</b>	<b>33.01 ± 0.86</b>	65.66 ± 1.76	<b>39.60 ± 0.06</b>	<b>61.19 ± 0.43</b>	<b>58.25 ± 0.79</b>	27.45 ± 1.13	<b>69.57 ± 0.07</b>	<b>20.00 ± 0.31</b>	48.82 ± 0.38	80.95 ± 9.52
0.6	Wanda	41.98 ± 0.40	30.56 ± 0.32	51.00 ± 0.45	<b>64.82 ± 0.35</b>	35.12 ± 0.07	<b>56.56 ± 0.46</b>	50.58 ± 0.41	23.83 ± 0.12	65.58 ± 0.19	16.93 ± 0.07	<b>44.77 ± 0.10</b>	-
	<i>MOONSHOT</i> -Wanda	<b>37.73 ± 0.19</b>	<b>27.71 ± 0.26</b>	<b>46.47 ± 0.08</b>	61.33 ± 0.91	35.53 ± 0.08	54.83 ± 0.14	<b>52.53 ± 0.29</b>	24.69 ± 0.21	66.81 ± 0.03	16.60 ± 0.23	44.62 ± 0.12	61.90 ± 4.76
	<i>MOONSHOT</i> (all)-Wanda	37.83 ± 0.04	27.84 ± 0.18	47.37 ± 0.57	60.40 ± 0.51	<b>35.54 ± 0.07</b>	56.30 ± 0.21	52.24 ± 0.17	<b>24.72 ± 0.16</b>	<b>66.85 ± 0.35</b>	<b>17.20 ± 0.23</b>	44.75 ± 0.14	<b>76.19 ± 4.76</b>

With the exception of Wanda on Llama-3.2-3B, we observe consistent additional improvements when extending *MOONSHOT* to projection layers, indicating that the multi-objective formulation is beneficial beyond the attention blocks. However, pruning time increases by  $8\times$  for Wanda and  $12\times$  for SparseGPT on Llama-3.2-1B. For Llama-3.2-3B, pruning time increases by  $8\times$  for Wanda and  $6\times$  for SparseGPT. This increase is due to the smaller feasible  $K_p$  and the increased computations with the larger Hessian.

Pruning projection layers with *MOONSHOT* is thus a practical compute/memory trade-off: depending on available resources, users may choose to apply *MOONSHOT* only to attention layers for efficiency, or extend it to projection layers for additional performance gains. Since pruning is typically performed once as an offline step, the extra runtime can be justified when resources permit, given the consistent improvements we observe.

## A.9 Pruning Llama-3.2 Instruct Models

We provide in this section additional results for the Instruct version of Llama-3.2-1B and Llama-3.2-3B.

Table 9: Test perplexity on C4, WikiText2 and PTB and zero-shot accuracy of Llama-3.2-1B-Instruct using *MOONSHOT*-SparseGPT. Zero-shot mean performance and win rate are computed over the 7 classification tasks. Results are averaged over 3 seeds with standard errors.

Sparsity	Method	C4 ↓	WikiText2 ↓	PTB ↓	BoolQ ↑	HellaSwag ↑	WinoGrande ↑	ARC-e ↑	ARC-e ↑	PIQA ↑	OBQA ↑	Mean ↑	Win Rate ↑
Dense	-	21.31	13.16	25.69	69.36	45.11	59.67	68.31	35.75	74.16	24.80	53.88	-
0.5	0.0	35.63 ± 0.15	26.56 ± 0.26	47.38 ± 0.53	62.87 ± 0.32	36.55 ± 0.18	54.33 ± 0.37	55.99 ± 0.12	25.91 ± 0.47	67.16 ± 0.25	20.33 ± 0.27	46.16 ± 0.09	0.0 ± 0.0
	0.1	31.72 ± 0.06	22.93 ± 0.07	42.37 ± 0.55	63.22 ± 0.18	38.0 ± 0.05	54.46 ± 0.43	58.57 ± 0.16	27.9 ± 0.44	68.82 ± 0.17	20.8 ± 0.53	47.4 ± 0.08	38.1 ± 4.76
	0.25	31.3 ± 0.07	22.52 ± 0.05	41.28 ± 0.48	63.19 ± 0.12	38.23 ± 0.09	54.91 ± 0.03	58.52 ± 0.33	28.04 ± 0.24	68.7 ± 0.21	20.2 ± 0.4	47.4 ± 0.08	42.86 ± 0.0
	0.5	31.17 ± 0.04	22.43 ± 0.1	41.27 ± 0.61	63.2 ± 0.08	38.28 ± 0.03	54.78 ± 0.73	58.61 ± 0.15	<b>28.73</b> ± 0.21	<b>69.04</b> ± 0.14	21.2 ± 0.35	47.64 ± 0.07	38.1 ± 4.76
	0.75	31.01 ± 0.03	<b>22.26</b> ± 0.06	40.84 ± 0.78	63.82 ± 0.21	38.37 ± 0.07	54.78 ± 0.51	58.95 ± 0.18	27.7 ± 0.12	68.48 ± 0.13	20.8 ± 0.46	47.56 ± 0.07	47.62 ± 4.76
	0.9	<b>30.97</b> ± 0.03	<b>22.26</b> ± 0.09	<b>40.2</b> ± 0.38	63.74 ± 0.15	<b>38.45</b> ± 0.04	55.12 ± 0.78	<b>59.15</b> ± 0.11	28.58 ± 0.52	68.97 ± 0.11	20.93 ± 0.27	<b>47.85</b> ± 0.15	<b>61.9</b> ± 12.6
	1.0	33.94 ± 0.18	24.56 ± 0.15	44.76 ± 0.34	<b>63.98</b> ± 0.35	38.39 ± 0.06	<b>55.93</b> ± 0.13	57.48 ± 0.14	27.53 ± 0.48	68.35 ± 0.43	<b>22.27</b> ± 0.13	47.7 ± 0.07	-
0.6	0.0	73.52 ± 0.28	63.14 ± 0.29	97.29 ± 0.62	62.21 ± 0.15	30.95 ± 0.17	51.99 ± 0.66	43.2 ± 0.68	21.05 ± 0.3	60.75 ± 0.04	15.47 ± 0.55	40.8 ± 0.14	14.29 ± 8.25
	0.1	54.01 ± 0.28	43.44 ± 0.9	69.89 ± 1.4	62.27 ± 0.06	32.69 ± 0.07	51.75 ± 0.5	47.95 ± 0.69	<b>23.04</b> ± 0.49	62.88 ± 0.32	<b>17.33</b> ± 0.41	42.56 ± 0.13	42.86 ± 8.25
	0.25	52.92 ± 0.29	42.57 ± 0.84	70.33 ± 1.27	62.4 ± 0.08	32.97 ± 0.06	51.64 ± 0.43	48.74 ± 0.67	22.53 ± 0.26	63.13 ± 0.33	16.47 ± 0.75	42.55 ± 0.15	47.62 ± 4.76
	0.5	52.01 ± 0.14	41.49 ± 0.78	67.8 ± 1.32	62.26 ± 0.05	33.06 ± 0.05	<b>52.72</b> ± 0.16	49.06 ± 0.83	22.87 ± 0.2	63.62 ± 0.1	16.27 ± 0.57	42.84 ± 0.13	57.14 ± 0.0
	0.75	51.13 ± 0.18	40.71 ± 0.74	66.39 ± 0.56	62.35 ± 0.15	33.11 ± 0.19	52.57 ± 0.25	49.45 ± 0.79	22.84 ± 0.27	63.76 ± 0.27	17.07 ± 0.41	43.02 ± 0.3	<b>76.19</b> ± 12.6
	0.9	<b>50.87</b> ± 0.22	<b>40.47</b> ± 0.43	<b>66.03</b> ± 0.58	62.46 ± 0.13	<b>33.29</b> ± 0.11	52.07 ± 0.03	<b>49.54</b> ± 0.8	<b>23.04</b> ± 0.34	<b>64.15</b> ± 0.3	16.67 ± 0.07	<b>43.03</b> ± 0.22	71.43 ± 16.5
	1.0	61.34 ± 0.64	51.01 ± 1.18	77.49 ± 1.38	<b>62.57</b> ± 0.05	32.81 ± 0.09	51.96 ± 0.44	47.98 ± 0.53	22.75 ± 0.23	63.22 ± 0.65	16.27 ± 0.41	42.51 ± 0.06	-
0.7	0.0	339.24 ± 24.3	347.52 ± 25.13	561.79 ± 42.06	54.19 ± 1.67	27.34 ± 0.12	49.41 ± 0.85	31.72 ± 0.27	18.52 ± 0.47	55.42 ± 0.54	13.0 ± 0.42	35.66 ± 0.29	4.76 ± 4.76
	0.1	169.63 ± 5.41	151.44 ± 4.14	263.45 ± 20.91	61.02 ± 0.5	28.14 ± 0.16	<b>51.93</b> ± 1.22	36.43 ± 0.36	<b>19.45</b> ± 0.2	57.54 ± 0.51	14.47 ± 0.93	48.56 ± 0.27	80.95 ± 4.76
	0.25	158.23 ± 2.87	<b>147.81</b> ± 4.46	247.05 ± 8.59	<b>61.97</b> ± 0.11	28.27 ± 0.12	50.96 ± 0.37	36.25 ± 0.44	19.11 ± 0.3	57.47 ± 0.21	14.27 ± 0.41	38.33 ± 0.2	80.95 ± 12.6
	0.5	155.75 ± 2.38	154.07 ± 8.74	245.63 ± 12.03	61.62 ± 0.11	28.4 ± 0.18	51.33 ± 0.54	36.55 ± 0.65	19.23 ± 0.32	57.94 ± 0.11	14.93 ± 0.29	38.57 ± 0.13	80.95 ± 12.6
	0.75	<b>154.39</b> ± 2.31	156.53 ± 7.65	249.79 ± 14.84	61.35 ± 0.11	28.4 ± 0.08	51.7 ± 0.55	37.37 ± 1.03	19.31 ± 0.12	<b>58.31</b> ± 0.24	<b>15.07</b> ± 0.37	<b>38.79</b> ± 0.13	85.71 ± 8.25
	0.9	<b>152.47</b> ± 1.6	159.61 ± 7.77	<b>235.26</b> ± 9.45	61.53 ± 0.04	<b>28.49</b> ± 0.03	<b>51.38</b> ± 0.24	<b>37.85</b> ± 0.66	19.08 ± 0.16	58.23 ± 0.31	14.2 ± 0.12	38.68 ± 0.15	<b>90.48</b> ± 9.52
	1.0	220.95 ± 7.06	278.28 ± 18.56	393.47 ± 37.16	60.43 ± 0.48	27.92 ± 0.1	50.51 ± 0.83	33.26 ± 0.48	19.11 ± 0.26	56.67 ± 0.07	13.67 ± 0.18	37.37 ± 0.21	-

Table 10: Test perplexity on C4, WikiText2 and PTB and zero-shot accuracy of Llama-3.2-1B-Instruct using *MOONSHOT*-Wanda. Zero-shot mean performance and win rate are computed over the 7 classification tasks. Results are averaged over 3 seeds with standard errors.

Sparsity	Method	C4 ↓	WikiText2 ↓	PTB ↓	BoolQ ↑	HellaSwag ↑	WinoGrande ↑	ARC-e ↑	ARC-e ↑	PIQA ↑	OBQA ↑	Mean ↑	Win Rate ↑
Dense	-	21.31	13.16	25.69	69.36	45.11	59.67	68.31	35.75	74.16	24.80	53.88	-
0.5	0.0	41.06 ± 0.24	29.22 ± 0.09	53.47 ± 0.42	62.19 ± 0.11	35.7 ± 0.01	53.67 ± 0.25	56.87 ± 0.43	26.14 ± 0.22	66.52 ± 0.16	16.8 ± 0.23	45.41 ± 0.1	42.86 ± 0.0
	0.1	40.41 ± 0.11	28.61 ± 0.06	52.63 ± 0.41	62.85 ± 0.37	35.92 ± 0.05	53.64 ± 0.27	56.69 ± 0.42	26.11 ± 0.44	66.56 ± 0.1	16.87 ± 0.29	45.52 ± 0.07	61.9 ± 12.6
	0.25	39.73 ± 0.1	<b>27.87</b> ± 0.09	51.39 ± 0.4	62.61 ± 0.31	35.99 ± 0.06	54.2 ± 0.38	57.1 ± 0.28	26.22 ± 0.32	66.65 ± 0.16	17.0 ± 0.31	45.68 ± 0.06	52.38 ± 9.52
	0.5	<b>39.58</b> ± 0.28	<b>27.87</b> ± 0.19	<b>51.26</b> ± 0.5	62.75 ± 0.23	36.12 ± 0.01	53.96 ± 0.09	<b>57.25</b> ± 0.33	26.02 ± 0.52	66.72 ± 0.16	17.07 ± 0.33	45.7 ± 0.1	61.9 ± 9.52
	0.75	39.91 ± 0.27	28.13 ± 0.1	51.4 ± 0.12	62.95 ± 0.59	36.22 ± 0.04	54.06 ± 0.68	56.9 ± 0.15	<b>26.39</b> ± 0.35	<b>67.05</b> ± 0.3	17.27 ± 0.24	45.83 ± 0.22	<b>66.67</b> ± 17.17
	0.9	40.16 ± 0.17	28.28 ± 0.03	52.16 ± 0.18	<b>63.06</b> ± 0.27	<b>36.34</b> ± 0.0	54.72 ± 0.48	56.96 ± 0.05	26.28 ± 0.2	66.74 ± 0.13	17.87 ± 0.07	<b>46.0</b> ± 0.1	<b>66.67</b> ± 12.6
	1.0	46.5 ± 0.11	33.63 ± 0.04	59.43 ± 0.11	62.62 ± 0.18	35.77 ± 0.08	<b>55.99</b> ± 0.52	54.05 ± 0.24	25.97 ± 0.34	65.89 ± 0.06	<b>18.07</b> ± 0.44	45.48 ± 0.2	-
0.6	0.0	106.89 ± 2.02	85.13 ± 2.17	115.83 ± 3.22	62.2 ± 0.02	29.38 ± 0.1	51.54 ± 0.36	41.41 ± 0.37	18.66 ± 0.06	59.47 ± 0.03	14.0 ± 0.2	39.52 ± 0.01	47.62 ± 4.76
	0.1	100.94 ± 1.72	78.78 ± 1.76	109.13 ± 2.96	62.15 ± 0.02	29.53 ± 0.12	52.14 ± 0.46	41.46 ± 0.37	19.25 ± 0.23	60.08 ± 0.14	14.6 ± 0.46	39.89 ± 0.03	66.67 ± 4.76
	0.25	96.71 ± 1.01	74.61 ± 1.17	<b>103.93</b> ± 2.76	62.15 ± 0.01	29.71 ± 0.06	52.54 ± 0.59	42.19 ± 0.46	19.48 ± 0.27	60.37 ± 0.34	14.8 ± 0.12	40.18 ± 0.23	76.19 ± 12.6
	0.5	<b>96.18</b> ± 1.15	<b>72.81</b> ± 1.32	104.85 ± 2.75	62.27 ± 0.01	29.91 ± 0.09	52.83 ± 0.37	42.68 ± 0.33	19.8 ± 0.26	60.45 ± 0.17	<b>15.2</b> ± 0.23	40.45 ± 0.13	<b>90.48</b> ± 4.76
	0.75	97.72 ± 0.56	74.27 ± 1.04	106.11 ± 2.83	62.2 ± 0.02	30.21 ± 0.1	<b>53.01</b> ± 0.32	43.17 ± 0.23	19.97 ± 0.34	60.43 ± 0.27	14.6 ± 0.23	40.51 ± 0.1	76.19 ± 9.52
	0.9	99.91 ± 1.39	77.19 ± 1.06	105.98 ± 1.66	<b>62.31</b> ± 0.06	<b>30.29</b> ± 0.07	52.2 ± 0.38	<b>43.74</b> ± 0.21	<b>20.42</b> ± 0.21	<b>60.57</b> ± 0.07	14.2 ± 0.12	<b>40.53</b> ± 0.06	76.19 ± 4.76
	1.0	154.55 ± 0.94	129.18 ± 1.52	141.38 ± 1.56	61.73 ± 0.09	29.62 ± 0.02	52.28 ± 0.23	39.69 ± 0.29	19.8 ± 0.64	59.16 ± 0.26	15.0 ± 0.12	39.61 ± 0.21	-
0.7	0.0	455.59 ± 8.85	467.43 ± 20.46	477.51 ± 15.39	49.42 ± 3.01	26.75 ± 0.08	49.33 ± 0.2	30.09 ± 0.33	<b>18.34</b> ± 0.27	55.02 ± 0.1	12.27 ± 0.27	34.46 ± 0.45	71.43 ± 8.25
	0.1	444.16 ± 9.19	437.37 ± 24.56	450.97 ± 22.14	53.84 ± 2.97	26.93 ± 0.02	50.51 ± 0.3	30.22 ± 0.26	18.0 ± 0.21	55.24 ± 0.11	12.0 ± 0.5	35.25 ± 0.34	71.43 ± 0.0
	0.25	<b>438.43</b> ± 21.47	416.16 ± 37.44	428.54 ± 38.04	55.31 ± 1.87	27.02 ± 0.03	50.91 ± 0.6	30.65 ± 0.14	18.15 ± 0.14	55.55 ± 0.24	12.4 ± 0.12	35.71 ± 0.18	80.95 ± 9.52
	0.5	454.78 ± 12.84	<b>413.61</b> ± 22.88	<b>425.15</b> ± 18.8	<b>58.4</b> ± 0.18	27.09 ± 0.04	51.04 ± 0.74	30.64 ± 0.11	17.92 ± 0.34	55.57 ± 0.24	12.93 ± 0.18	36.23 ± 0.08	85.71 ± 8.25
	0.75	481.74 ± 11.64	466.41 ± 26.96	426.24 ± 10.96	57.66 ± 1.42	27.0 ± 0.05	<b>52.67</b> ± 0.59	<b>31.02</b> ± 0.18	18.2 ± 0.11	55.77 ± 0.14	12.67 ± 0.33	<b>36.43</b> ± 0.24	<b>90.48</b> ± 9.52
	0.9	583.01 ± 11.37	559.07 ± 16.93	513.8 ± 22.06	56.83 ± 0.84	<b>27.15</b> ± 0.04	51.91 ± 0.25	31.0 ± 0.13	17.72 ± 0.11	<b>55.89</b> ± 0.52	12.73 ± 0.33	36.18 ± 0.08	76.19 ± 4.76
	1.0	1797.15 ± 176.66	2120.64 ± 117.11	2225.24 ± 246.18	52.1 ± 3.48	26.62 ± 0.06	49.25 ± 0.33	28.3 ± 0.2	18.26 ± 0.2	54.39 ± 0.15	<b>13.0</b> ± 0.23	34.56 ± 0.45	-

Table 11: Test perplexity on C4, WikiText2 and PTB and zero-shot accuracy of Llama-3.2-3B-Instruct using *MOONSHOT*-SparseGPT. Zero-shot mean performance and win rate are computed over the 7 classification tasks. Results are averaged over 3 seeds with standard errors.

Sparsity	Method	C4 ↓	WikiText2 ↓	PTB ↓	BoolQ ↑	HellaSwag ↑	WinoGrande ↑	ARC-e ↑	ARC-e ↑	PIQA ↑	OBQA ↑	Mean ↑	Win Rate ↑
Dense	-	16.49	11.04	20.42	78.47	52.24	67.40	73.99	43.43	75.79	27.40	59.82	-
0.5	0.0	24.55 ± 0.14	19.22 ± 0.44	33.43 ± 1.03	74.55 ± 0.01	43.43 ± 0.16	61.64 ± 0.25	64.37 ± 0.31	33.25 ± 0.6	70.51 ± 0.22	21.73 ± 0.48	52.78 ± 0.18	4.76 ± 4.76
	0.1	22.61 ± 0.02	17.45 ± 0.28	29.93 ± 0.2	75.52 ± 0.32	44.79 ± 0.1	62.88 ± 0.37	65.42 ± 0.27	34.3 ± 0.34	<b>71.85</b> ± 0.1	22.0 ± 0.2	53.82 ± 0.11	38.1 ± 4.76
	0.25	22.45 ± 0.01	17.21 ± 0.14	29.67 ± 0.09	75.77 ± 0.4	44.81 ± 0.02	63.27 ± 0.31	65.5 ± 0.56	34.3 ± 0.44	71.8 ± 0.32	22.13 ± 0.29	53.94 ± 0.2	38.1 ± 12.6
	0.5	22.41 ± 0.03	17.12 ± 0.18	29.42 ± 0.03	<b>76.13</b> ± 0.04	44.87 ± 0.05	63.04 ± 0.52	65.74 ± 0.53	<b>34.33</b> ± 0.37	71.84 ± 0.16	21.93 ± 0.64	<b>53.98</b> ± 0.16	<b>42.86</b> ± 14.29
	0.75	22.38 ± 0.04	17.09 ± 0.17	<b>29.21</b> ± 0.24	75.95 ± 0.15	44.88 ± 0.02	62.93 ± 0.21	66.11 ± 0.09	34.07 ± 0.31	71.47 ± 0.15	22.13 ± 0.59	53.94 ± 0.14	38.1 ± 4.76
	0.9	<b>22.37</b> ± 0.02	<b>17.04</b> ± 0.14	29.24 ± 0.27	75.81 ± 0.13	44.89 ± 0.01	63.35 ± 0.13	<b>66.16</b> ± 0.06	<b>34.33</b> ± 0.14	71.65 ± 0.05	22.33 ± 0.24	54.08 ± 0.02	<b>42.86</b> ± 8.25
	1.0	23.11 ± 0.											

Table 12: Test perplexity on C4, WikiText2 and PTB and zero-shot accuracy of Llama-3.2-3B-Instruct using *MOONSHOT*-Wanda. Zero-shot mean performance and win rate are computed over the 7 classification tasks. Results are averaged over 3 seeds with standard errors.

Sparsity	Method	C4 ↓	WikiText2 ↓	PTB ↓	BoolQ ↑	HellaSwag ↑	WinoGrande ↑	ARC-e ↑	ARC-e ↑	PIQA ↑	OBQA ↑	Mean ↑	Win Rate ↑
Dense	-	16.49	11.04	20.42	78.47	52.24	67.40	73.99	43.43	75.79	27.40	59.82	-
0.5	0.0	25.54 ± 0.11	19.79 ± 0.17	34.17 ± 0.36	73.04 ± 0.34	42.38 ± 0.08	60.8 ± 0.59	66.26 ± 0.23	33.73 ± 0.28	70.98 ± 0.05	21.47 ± 0.24	52.67 ± 0.22	38.1 ± 4.76
	0.1	25.15 ± 0.07	19.26 ± 0.1	33.32 ± 0.19	73.21 ± 0.33	42.57 ± 0.05	62.04 ± 0.21	66.02 ± 0.09	33.79 ± 0.27	70.96 ± 0.1	<b>22.33 ± 0.07</b>	52.99 ± 0.12	42.86 ± 0.0
	0.25	24.86 ± 0.07	19.06 ± 0.2	32.71 ± 0.21	73.79 ± 0.27	42.7 ± 0.05	61.88 ± 0.28	66.22 ± 0.15	34.19 ± 0.29	<b>71.42 ± 0.02</b>	21.93 ± 0.24	53.16 ± 0.13	47.62 ± 4.76
	0.5	24.53 ± 0.06	18.73 ± 0.14	32.3 ± 0.13	73.37 ± 0.38	42.88 ± 0.06	62.77 ± 0.47	66.33 ± 0.15	34.7 ± 0.27	71.27 ± 0.11	22.07 ± 0.13	53.34 ± 0.07	57.14 ± 0.0
	0.75	24.48 ± 0.13	18.57 ± 0.15	31.99 ± 0.23	73.86 ± 0.51	42.9 ± 0.07	62.67 ± 0.18	66.22 ± 0.13	35.01 ± 0.2	71.4 ± 0.13	21.93 ± 0.18	53.43 ± 0.08	<b>61.9 ± 4.76</b>
	0.9	<b>24.35 ± 0.08</b>	<b>18.46 ± 0.01</b>	<b>31.62 ± 0.17</b>	73.86 ± 0.26	43.2 ± 0.01	63.59 ± 0.09	<b>66.4 ± 0.1</b>	<b>35.32 ± 0.21</b>	71.25 ± 0.02	21.6 ± 0.23	<b>53.6 ± 0.04</b>	52.38 ± 4.76
1.0	24.76 ± 0.01	18.7 ± 0.04	32.27 ± 0.13	<b>74.28 ± 0.34</b>	<b>43.43 ± 0.06</b>	<b>63.85 ± 0.16</b>	64.66 ± 0.12	34.5 ± 0.22	70.26 ± 0.16	21.27 ± 0.35	53.18 ± 0.04	-	
0.6	0.0	73.41 ± 4.15	66.23 ± 4.51	101.71 ± 3.95	63.9 ± 0.4	32.57 ± 0.25	55.67 ± 0.43	49.89 ± 0.66	23.83 ± 0.52	63.64 ± 0.56	14.47 ± 0.24	43.42 ± 0.31	19.05 ± 12.6
	0.1	67.82 ± 1.91	61.5 ± 2.44	94.64 ± 1.55	64.48 ± 0.19	33.09 ± 0.18	55.62 ± 0.26	50.74 ± 0.95	24.94 ± 0.38	64.18 ± 0.34	15.0 ± 0.5	44.01 ± 0.29	23.81 ± 9.52
	0.25	63.87 ± 0.75	57.93 ± 1.4	90.5 ± 1.03	64.69 ± 0.04	33.55 ± 0.18	56.75 ± 0.21	51.32 ± 0.51	24.91 ± 0.45	64.24 ± 0.19	15.33 ± 0.07	44.4 ± 0.2	33.33 ± 4.76
	0.5	60.41 ± 0.83	54.03 ± 1.15	85.55 ± 0.97	65.38 ± 0.13	33.92 ± 0.17	56.91 ± 0.41	51.02 ± 0.58	25.17 ± 0.32	64.25 ± 0.17	15.87 ± 0.18	44.65 ± 0.09	47.62 ± 4.76
	0.75	57.64 ± 0.24	50.62 ± 0.77	82.07 ± 0.91	65.89 ± 0.24	34.25 ± 0.06	56.99 ± 0.16	<b>51.5 ± 0.26</b>	<b>25.74 ± 0.23</b>	64.33 ± 0.15	<b>16.67 ± 0.27</b>	45.05 ± 0.08	76.19 ± 9.52
	0.9	<b>56.47 ± 0.27</b>	<b>49.09 ± 0.69</b>	<b>80.34 ± 0.05</b>	<b>66.41 ± 0.4</b>	<b>34.49 ± 0.06</b>	<b>57.51 ± 0.27</b>	51.46 ± 0.09	25.65 ± 0.19	<b>64.54 ± 0.21</b>	16.2 ± 0.35	<b>45.18 ± 0.13</b>	<b>80.95 ± 12.6</b>
1.0	57.74 ± 0.74	50.82 ± 0.81	82.33 ± 0.32	65.54 ± 0.49	34.33 ± 0.07	57.14 ± 0.24	48.72 ± 0.28	25.06 ± 0.24	63.98 ± 0.23	<b>16.67 ± 0.07</b>	44.49 ± 0.06	-	
0.7	0.0	323.12 ± 19.39	340.56 ± 32.95	302.69 ± 10.61	39.43 ± 0.25	26.81 ± 0.07	49.46 ± 0.73	31.14 ± 0.19	<b>18.54 ± 0.28</b>	55.82 ± 0.41	11.67 ± 0.18	33.27 ± 0.18	42.86 ± 8.25
	0.1	301.09 ± 12.36	311.9 ± 24.51	288.06 ± 3.24	41.1 ± 0.2	26.86 ± 0.1	49.7 ± 0.6	31.1 ± 0.11	18.15 ± 0.2	56.33 ± 0.16	12.0 ± 0.2	33.6 ± 0.14	47.62 ± 12.6
	0.25	290.05 ± 12.47	296.65 ± 19.83	286.03 ± 6.45	45.08 ± 0.64	26.93 ± 0.05	49.83 ± 0.17	31.14 ± 0.29	17.95 ± 0.22	56.31 ± 0.17	11.4 ± 0.31	34.09 ± 0.15	38.1 ± 9.52
	0.5	277.25 ± 10.52	279.64 ± 12.43	282.93 ± 6.27	46.76 ± 1.4	27.1 ± 0.05	49.96 ± 0.05	31.59 ± 0.12	17.61 ± 0.08	56.57 ± 0.11	11.87 ± 0.35	34.49 ± 0.26	42.86 ± 8.25
	0.75	264.81 ± 3.25	266.54 ± 8.53	280.08 ± 6.69	47.18 ± 0.96	27.26 ± 0.03	<b>50.3 ± 0.25</b>	31.75 ± 0.11	18.0 ± 0.3	56.06 ± 0.33	11.73 ± 0.13	34.7 ± 0.21	<b>52.38 ± 4.76</b>
	0.9	255.35 ± 1.85	<b>258.39 ± 3.53</b>	268.83 ± 4.87	45.51 ± 1.63	27.38 ± 0.04	<b>50.3 ± 0.15</b>	<b>31.8 ± 0.2</b>	17.78 ± 0.19	<b>57.05 ± 0.3</b>	11.47 ± 0.13	34.47 ± 0.29	<b>52.38 ± 12.6</b>
1.0	<b>234.03 ± 3.39</b>	271.85 ± 7.51	<b>264.91 ± 5.09</b>	<b>55.83 ± 0.19</b>	<b>27.42 ± 0.04</b>	48.8 ± 0.18	31.0 ± 0.18	17.83 ± 0.25	56.84 ± 0.18	<b>12.2 ± 0.2</b>	<b>35.7 ± 0.07</b>	-	

## A.10 Additional Ablations on $\lambda$

In Section 3, we analyze the sensitivity of  $\lambda$  and find that the optimum almost never occurs at the extremes  $\lambda \in (0, 1)$ . Across all architectures, pruning baselines, and sparsity regimes we tested, intermediate values consistently outperform the single-objective endpoints, as shown in Figure 5. This pattern supports the idea that balancing the reconstruction and Fisher terms yields a more robust pruning criterion than either alone.

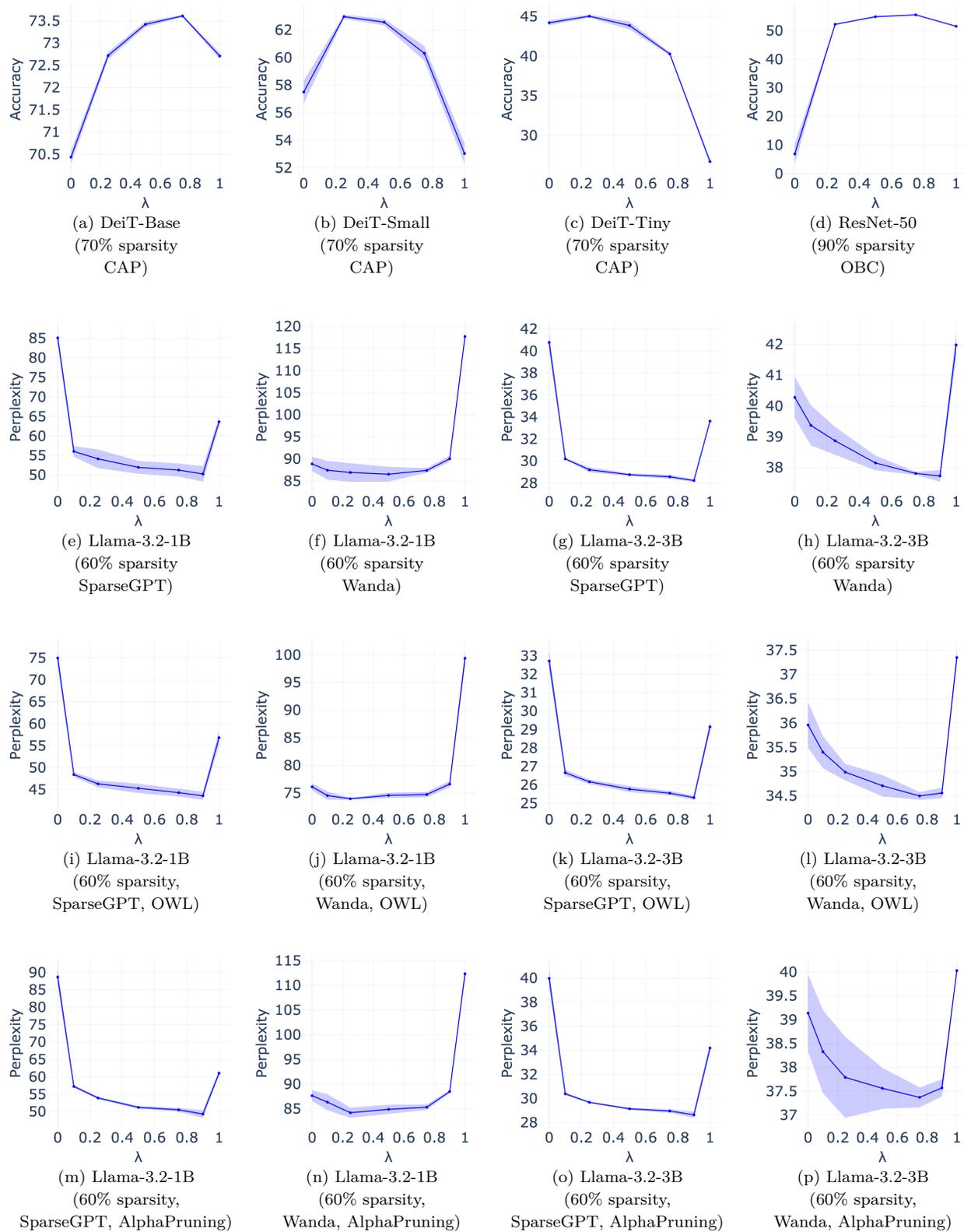


Figure 5: Performance of *MOONSHOT* across different values of  $\lambda$  on the DeiT models (70% sparsity), ResNet-50 (90% sparsity), and Llama-3.2 models (60% and 2:4 sparsity), using CAP, OBC, and SparseGPT as base methods respectively. Accuracy is reported for vision models and perplexity on C4 for LLMs.

### A.11 Further Evaluation of *MOONSHOT* Across Sparsity Regimes

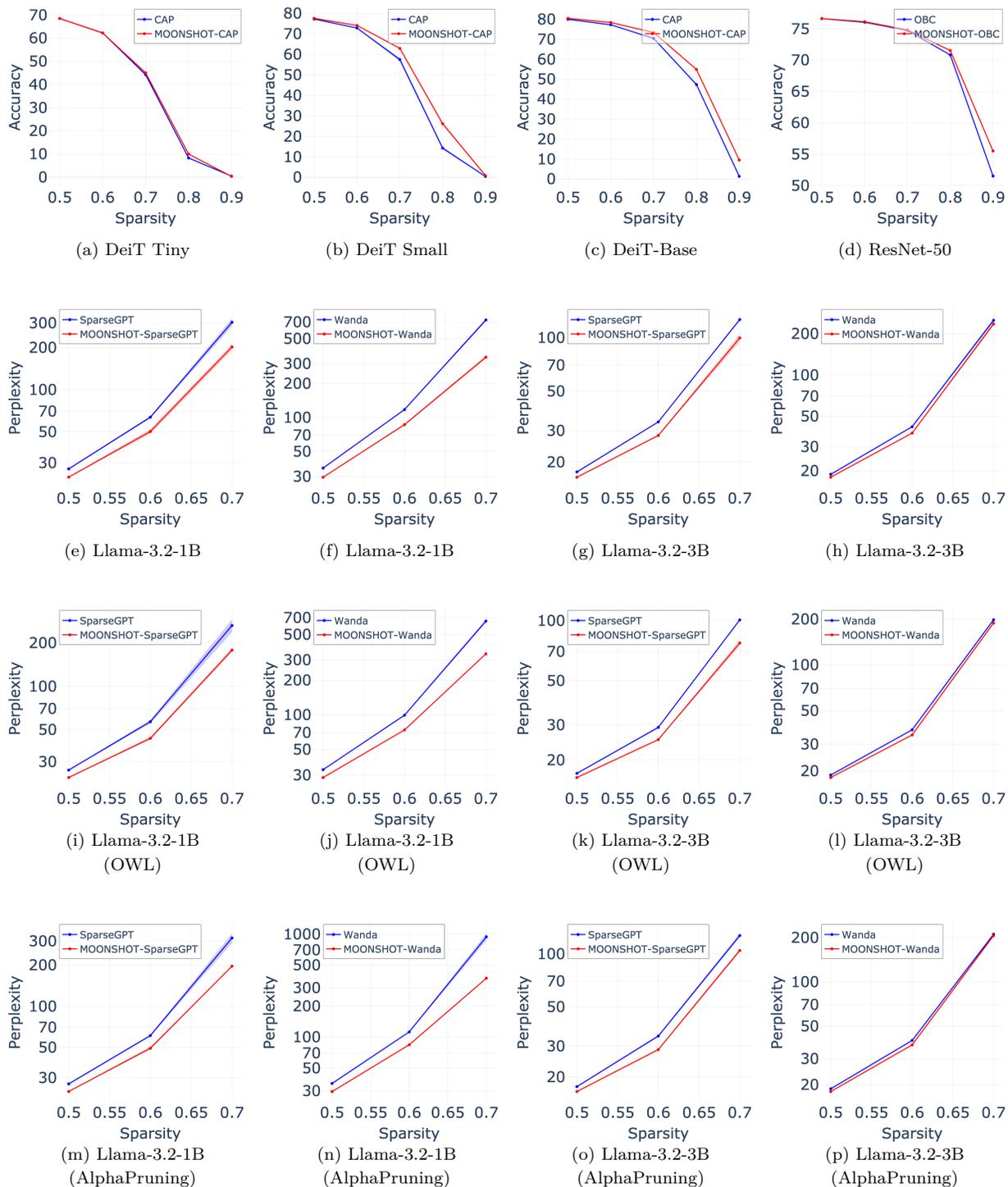


Figure 6: Impact of *MOONSHOT* across sparsity levels on CAP for the DeiT models, OBC on ResNet-50 and SparseGPT/Wanda on the Llama-3.2 models.

Figure 6 expands the sparsity sweeps for all architectures, pruning algorithms and pruning methods, and shows a consistent trend: *MOONSHOT* consistently yields better performance–sparsity tradeoff than their single-objective counterparts, with the gap widening in the high-sparsity regime where baselines degrade most. Moreover, when combined with non-uniform sparsity allocation methods (OWL and AlphaPruning), *MOONSHOT*’s gains are additive: curves shift upward at nearly all sparsity regimes, indicating that our multi-objective signal complements allocation strategies rather than replacing them.

## A.12 Comprehensive Experimental Results

Tables 2 and 3, as well as Figure 1 show the results of *MOONSHOT* for the optimal value of  $\lambda$  across a selection of sparsity regimes. In this section, we provide the results of *MOONSHOT* precisely for each value of  $\lambda$  that was tried, across all sparsity regimes.

Table 13: Test Accuracy for the DeiT models and ResNet-50 using *MOONSHOT*-CAP and *MOONSHOT*-OBC respectively

Sparsity	$\lambda$	DeiT Tiny	DeiT Small	DeiT Base	ResNet-50
Dense	-	72.14	79.83	81.80	77.11
0.5	0.00	68.49 $\pm$ 0.1	77.27 $\pm$ 0.03	80.01 $\pm$ 0.01	50.88 $\pm$ 25.39
	0.25	<b>68.62</b> $\pm$ 0.03	77.63 $\pm$ 0.02	80.5 $\pm$ 0.02	76.56 $\pm$ 0.03
	0.50	68.35 $\pm$ 0.05	<b>77.67</b> $\pm$ 0.01	80.56 $\pm$ 0.02	76.61 $\pm$ 0.03
	0.75	68.02 $\pm$ 0.07	77.49 $\pm$ 0.03	<b>80.6</b> $\pm$ 0.02	<b>76.63</b> $\pm$ 0.04
	1.00	65.49 $\pm$ 0.02	76.56 $\pm$ 0.04	80.58 $\pm$ 0.01	<b>76.63</b> $\pm$ 0.05
0.6	0.00	<b>62.28</b> $\pm$ 0.05	72.89 $\pm$ 0.04	77.27 $\pm$ 0.1	50.37 $\pm$ 25.13
	0.25	62.22 $\pm$ 0.09	<b>74.16</b> $\pm$ 0.04	78.41 $\pm$ 0.06	76.04 $\pm$ 0.01
	0.50	61.76 $\pm$ 0.1	74.14 $\pm$ 0.02	78.62 $\pm$ 0.04	<b>76.13</b> $\pm$ 0.02
	0.75	60.7 $\pm$ 0.2	73.76 $\pm$ 0.09	<b>78.81</b> $\pm$ 0.04	76.11 $\pm$ 0.0
	1.00	54.18 $\pm$ 0.15	71.31 $\pm$ 0.19	78.67 $\pm$ 0.01	76.04 $\pm$ 0.02
0.7	0.00	44.22 $\pm$ 0.32	57.5 $\pm$ 0.83	70.44 $\pm$ 0.15	48.94 $\pm$ 24.42
	0.25	<b>45.05</b> $\pm$ 0.2	<b>62.97</b> $\pm$ 0.15	72.72 $\pm$ 0.08	74.7 $\pm$ 0.04
	0.50	43.87 $\pm$ 0.49	62.57 $\pm$ 0.23	73.42 $\pm$ 0.06	<b>74.82</b> $\pm$ 0.05
	0.75	40.29 $\pm$ 0.19	60.31 $\pm$ 0.54	<b>73.61</b> $\pm$ 0.03	<b>74.82</b> $\pm$ 0.04
	1.00	26.71 $\pm$ 0.27	53.04 $\pm$ 0.74	72.71 $\pm$ 0.07	74.73 $\pm$ 0.03
0.8	0.00	8.28 $\pm$ 0.32	14.28 $\pm$ 0.65	47.32 $\pm$ 0.48	44.4 $\pm$ 22.15
	0.25	<b>9.97</b> $\pm$ 0.26	<b>26.2</b> $\pm$ 0.59	53.56 $\pm$ 0.12	71.02 $\pm$ 0.08
	0.50	8.58 $\pm$ 0.26	24.36 $\pm$ 0.56	54.9 $\pm$ 0.26	71.39 $\pm$ 0.03
	0.75	6.23 $\pm$ 0.15	19.27 $\pm$ 0.88	<b>55.37</b> $\pm$ 0.2	<b>71.54</b> $\pm$ 0.02
	1.00	2.14 $\pm$ 0.16	9.84 $\pm$ 0.62	49.2 $\pm$ 0.18	70.83 $\pm$ 0.04
0.9	0.00	0.43 $\pm$ 0.06	0.43 $\pm$ 0.03	1.37 $\pm$ 0.09	6.89 $\pm$ 3.42
	0.25	0.41 $\pm$ 0.05	<b>0.91</b> $\pm$ 0.07	7.79 $\pm$ 0.29	52.2 $\pm$ 0.06
	0.50	0.43 $\pm$ 0.09	0.86 $\pm$ 0.09	<b>9.49</b> $\pm$ 0.37	54.86 $\pm$ 0.28
	0.75	<b>0.5</b> $\pm$ 0.04	0.72 $\pm$ 0.08	9.4 $\pm$ 0.28	<b>55.52</b> $\pm$ 0.09
	1.00	0.3 $\pm$ 0.03	0.46 $\pm$ 0.08	3.07 $\pm$ 0.07	51.52 $\pm$ 0.07
2:4	0.00	52.28 $\pm$ 0.04	69.65 $\pm$ 0.02	76.21 $\pm$ 0.07	0.1 $\pm$ 0.0
	0.25	<b>54.23</b> $\pm$ 0.1	71.1 $\pm$ 0.07	77.3 $\pm$ 0.04	75.37 $\pm$ 0.02
	0.50	54.2 $\pm$ 0.15	<b>71.54</b> $\pm$ 0.08	77.67 $\pm$ 0.07	<b>75.5</b> $\pm$ 0.04
	0.75	53.78 $\pm$ 0.03	71.52 $\pm$ 0.05	<b>77.88</b> $\pm$ 0.05	<b>75.5</b> $\pm$ 0.03
	1.00	47.65 $\pm$ 0.11	70.25 $\pm$ 0.04	77.74 $\pm$ 0.04	75.46 $\pm$ 0.03

Table 14: Test perplexity on C4, WikiText2 and PTB and zero-shot accuracy of Llama-3.2-1B using *MOON-SHOT-OSSCAR*. Zero-shot mean performance and win rate are computed over the 7 classification tasks. Results are averaged over 3 seeds with standard errors.

Sparsity	$\lambda$	C4 $\downarrow$	WikiText2 $\downarrow$	PTB $\downarrow$	BoolQ $\uparrow$	HellaSwag $\uparrow$	WinoGrande $\uparrow$	ARC-e $\uparrow$	ARC-c $\uparrow$	PIQA $\uparrow$	OBQA $\uparrow$	Mean $\uparrow$	Win Rate $\uparrow$
Dense	-	14.02	9.75	17.59	64.01	47.73	60.14	65.15	31.23	74.32	26.4	52.71	-
0.1	0.0	315.6 $\pm$ 149.49	242.7 $\pm$ 114.26	4120.96 $\pm$ 2133.48	51.86 $\pm$ 1.12	30.82 $\pm$ 3.98	50.51 $\pm$ 1.02	45.45 $\pm$ 4.51	22.41 $\pm$ 2.1	60.99 $\pm$ 3.65	15.27 $\pm$ 1.17	39.62 $\pm$ 2.5	28.57 $\pm$ 28.57
	0.1	<b>38.04</b> $\pm$ 0.58	30.93 $\pm$ 0.92	<b>86.73</b> $\pm$ 9.88	<b>56.55</b> $\pm$ 0.9	37.89 $\pm$ 1.11	<b>55.51</b> $\pm$ 0.37	<b>58.84</b> $\pm$ 0.38	28.58 $\pm$ 0.51	<b>71.49</b> $\pm$ 0.06	<b>18.93</b> $\pm$ 0.48	<b>46.83</b> $\pm$ 0.42	80.95 $\pm$ 12.6
	0.25	41.23 $\pm$ 1.45	33.54 $\pm$ 1.31	102.9 $\pm$ 6.89	56.02 $\pm$ 1.61	37.44 $\pm$ 1.22	54.12 $\pm$ 0.66	58.52 $\pm$ 0.51	<b>28.64</b> $\pm$ 0.21	71.33 $\pm$ 0.22	17.8 $\pm$ 0.42	46.27 $\pm$ 0.53	76.19 $\pm$ 4.76
	0.5	39.35 $\pm$ 2.63	<b>30.62</b> $\pm$ 1.23	91.55 $\pm$ 16.49	54.31 $\pm$ 3.38	<b>38.75</b> $\pm$ 0.98	55.22 $\pm$ 0.61	57.46 $\pm$ 2.72	27.9 $\pm$ 0.51	71.2 $\pm$ 0.81	17.67 $\pm$ 1.77	46.07 $\pm$ 1.49	80.95 $\pm$ 12.6
	0.75	38.46 $\pm$ 1.14	33.27 $\pm$ 0.24	97.7 $\pm$ 8.46	52.1 $\pm$ 2.0	36.75 $\pm$ 0.38	53.64 $\pm$ 0.5	55.2 $\pm$ 2.72	27.05 $\pm$ 0.53	70.57 $\pm$ 1.2	17.47 $\pm$ 0.74	44.68 $\pm$ 0.94	80.95 $\pm$ 9.52
	0.9	39.31 $\pm$ 1.06	37.59 $\pm$ 0.66	97.6 $\pm$ 9.05	52.57 $\pm$ 1.72	36.59 $\pm$ 0.45	55.2 $\pm$ 1.05	48.13 $\pm$ 6.52	24.91 $\pm$ 1.96	67.01 $\pm$ 3.23	17.4 $\pm$ 0.42	43.12 $\pm$ 1.85	<b>85.71</b> $\pm$ 8.25
1.0	43.0 $\pm$ 1.28	43.91 $\pm$ 0.6	106.62 $\pm$ 8.57	52.61 $\pm$ 1.63	35.82 $\pm$ 0.29	54.17 $\pm$ 1.11	45.19 $\pm$ 7.83	24.23 $\pm$ 2.04	65.65 $\pm$ 3.68	15.93 $\pm$ 0.33	41.94 $\pm$ 1.98	0.0 $\pm$ 0.0	
0.15	0.0	379.28 $\pm$ 177.85	302.81 $\pm$ 141.13	1378.15 $\pm$ 633.62	50.31 $\pm$ 1.0	29.86 $\pm$ 3.4	51.3 $\pm$ 1.22	42.82 $\pm$ 5.83	21.76 $\pm$ 1.67	60.95 $\pm$ 4.05	15.87 $\pm$ 1.18	38.98 $\pm$ 2.56	38.1 $\pm$ 31.23
	0.1	148.98 $\pm$ 3.3	113.79 $\pm$ 4.23	392.84 $\pm$ 35.98	<b>54.37</b> $\pm$ 1.47	30.91 $\pm$ 0.41	54.12 $\pm$ 0.69	51.37 $\pm$ 0.27	<b>25.74</b> $\pm$ 0.25	64.54 $\pm$ 0.48	<b>16.27</b> $\pm$ 0.52	42.48 $\pm$ 0.35	61.9 $\pm$ 12.6
	0.25	87.21 $\pm$ 20.52	79.97 $\pm$ 20.05	365.34 $\pm$ 138.38	53.16 $\pm$ 1.55	32.68 $\pm$ 1.47	53.43 $\pm$ 0.57	<b>52.62</b> $\pm$ 1.92	25.48 $\pm$ 0.75	66.54 $\pm$ 1.46	16.13 $\pm$ 0.58	<b>42.86</b> $\pm$ 1.09	71.43 $\pm$ 8.25
	0.5	64.68 $\pm$ 9.44	61.36 $\pm$ 15.64	282.32 $\pm$ 156.56	51.96 $\pm$ 2.46	33.42 $\pm$ 1.79	<b>54.7</b> $\pm$ 0.66	46.83 $\pm$ 9.01	24.46 $\pm$ 2.89	64.91 $\pm$ 5.71	15.93 $\pm$ 1.76	41.74 $\pm$ 3.34	47.62 $\pm$ 17.17
	0.75	48.07 $\pm$ 2.62	<b>42.59</b> $\pm$ 0.5	117.8 $\pm$ 4.46	47.82 $\pm$ 1.19	<b>36.05</b> $\pm$ 0.43	54.06 $\pm$ 0.46	49.51 $\pm$ 5.8	25.11 $\pm$ 1.77	<b>68.03</b> $\pm$ 2.48	15.8 $\pm$ 0.58	42.34 $\pm$ 1.68	66.67 $\pm$ 4.76
	0.9	<b>47.48</b> $\pm$ 1.15	46.86 $\pm$ 1.09	<b>102.39</b> $\pm$ 3.55	48.88 $\pm$ 2.07	35.65 $\pm$ 0.5	54.22 $\pm$ 0.32	46.49 $\pm$ 6.13	25.17 $\pm$ 1.7	66.74 $\pm$ 2.73	15.93 $\pm$ 0.74	41.87 $\pm$ 1.83	<b>76.19</b> $\pm$ 9.52
1.0	52.65 $\pm$ 1.85	54.51 $\pm$ 1.35	132.62 $\pm$ 9.26	51.23 $\pm$ 2.07	35.22 $\pm$ 0.32	52.83 $\pm$ 0.43	45.61 $\pm$ 6.45	24.23 $\pm$ 1.72	66.12 $\pm$ 2.69	14.53 $\pm$ 0.55	41.4 $\pm$ 1.74	0.0 $\pm$ 0.0	

Table 15: Test perplexity on C4, WikiText2 and PTB and zero-shot accuracy of Llama-3.2-1B using *MOON-SHOT*-SparseGPT. Zero-shot **mean performance** and **win rate** are computed over the 7 classification tasks. Results are averaged over 3 seeds with standard errors.

Sparsity	$\lambda$	C4 $\downarrow$	WikiText2 $\downarrow$	PTB $\downarrow$	BoolQ $\uparrow$	HellaSwag $\uparrow$	WinoGrande $\uparrow$	ARC-e $\uparrow$	ARC-c $\uparrow$	PIQA $\uparrow$	OBQA $\uparrow$	Mean $\uparrow$	Win Rate $\uparrow$
Dense	-	14.02	9.75	17.59	64.01	47.73	60.14	65.15	31.23	74.32	26.4	52.71	-
0.5	0.0	29.14 $\pm$ 0.16	22.03 $\pm$ 0.15	36.04 $\pm$ 0.38	60.08 $\pm$ 0.67	36.28 $\pm$ 0.07	54.06 $\pm$ 0.67	52.83 $\pm$ 0.16	24.97 $\pm$ 0.5	66.92 $\pm$ 0.27	17.8 $\pm$ 0.12	44.71 $\pm$ 0.1	9.52 $\pm$ 4.76
	0.1	24.77 $\pm$ 0.09	18.43 $\pm$ 0.09	30.72 $\pm$ 0.35	62.18 $\pm$ 0.67	38.25 $\pm$ 0.03	55.38 $\pm$ 0.73	55.39 $\pm$ 0.49	26.68 $\pm$ 0.3	68.12 $\pm$ 0.05	19.6 $\pm$ 0.5	46.51 $\pm$ 0.16	33.33 $\pm$ 4.76
	0.25	24.31 $\pm$ 0.04	17.93 $\pm$ 0.1	30.25 $\pm$ 0.41	62.33 $\pm$ 0.77	38.52 $\pm$ 0.14	55.01 $\pm$ 0.14	55.46 $\pm$ 0.4	26.25 $\pm$ 0.21	<b>68.81</b> $\pm$ 0.07	18.87 $\pm$ 0.47	46.46 $\pm$ 0.13	38.1 $\pm$ 9.52
	0.5	24.01 $\pm$ 0.08	17.74 $\pm$ 0.14	29.72 $\pm$ 0.34	62.6 $\pm$ 0.18	38.61 $\pm$ 0.07	<b>55.54</b> $\pm$ 0.34	<b>56.0</b> $\pm$ 0.36	26.05 $\pm$ 0.12	68.3 $\pm$ 0.15	19.07 $\pm$ 0.41	46.6 $\pm$ 0.12	28.57 $\pm$ 14.29
	0.75	23.77 $\pm$ 0.1	17.53 $\pm$ 0.1	29.17 $\pm$ 0.14	<b>63.16</b> $\pm$ 0.11	38.71 $\pm$ 0.16	54.78 $\pm$ 0.56	55.82 $\pm$ 0.49	26.05 $\pm$ 0.37	68.63 $\pm$ 0.33	20.0 $\pm$ 0.5	46.74 $\pm$ 0.2	38.1 $\pm$ 9.52
	0.9	<b>23.7</b> $\pm$ 0.06	<b>17.49</b> $\pm$ 0.12	<b>29.05</b> $\pm$ 0.16	63.14 $\pm$ 0.31	38.86 $\pm$ 0.12	55.38 $\pm$ 0.14	55.6 $\pm$ 0.42	<b>26.96</b> $\pm$ 0.34	68.59 $\pm$ 0.42	20.27 $\pm$ 0.58	<b>46.97</b> $\pm$ 0.22	<b>47.62</b> $\pm$ 9.52
1.0	27.15 $\pm$ 0.23	19.98 $\pm$ 0.17	33.63 $\pm$ 0.1	61.67 $\pm$ 1.48	<b>39.19</b> $\pm$ 0.08	55.51 $\pm$ 0.38	55.16 $\pm$ 0.53	26.51 $\pm$ 0.53	68.59 $\pm$ 0.18	<b>21.87</b> $\pm$ 0.29	46.93 $\pm$ 0.32	-	
0.5 (Alpha-Pruning)	0.0	29.64 $\pm$ 0.28	22.33 $\pm$ 0.14	36.34 $\pm$ 0.53	61.12 $\pm$ 0.36	36.36 $\pm$ 0.19	55.96 $\pm$ 0.57	51.95 $\pm$ 0.49	25.63 $\pm$ 1.05	66.92 $\pm$ 0.27	18.67 $\pm$ 0.37	45.23 $\pm$ 0.35	14.29 $\pm$ 8.25
	0.1	24.9 $\pm$ 0.1	18.44 $\pm$ 0.09	30.33 $\pm$ 0.33	61.56 $\pm$ 0.98	38.33 $\pm$ 0.09	<b>56.06</b> $\pm$ 0.53	54.8 $\pm$ 0.4	26.34 $\pm$ 0.08	68.73 $\pm$ 0.26	19.27 $\pm$ 0.37	46.44 $\pm$ 0.13	28.57 $\pm$ 8.25
	0.25	24.47 $\pm$ 0.13	18.05 $\pm$ 0.16	29.97 $\pm$ 0.24	62.6 $\pm$ 0.32	38.62 $\pm$ 0.06	55.33 $\pm$ 0.24	54.98 $\pm$ 0.27	<b>26.91</b> $\pm$ 0.28	<b>68.81</b> $\pm$ 0.47	19.87 $\pm$ 0.94	46.73 $\pm$ 0.17	<b>42.86</b> $\pm$ 8.25
	0.5	24.17 $\pm$ 0.11	17.79 $\pm$ 0.19	29.57 $\pm$ 0.35	62.01 $\pm$ 1.3	38.79 $\pm$ 0.08	55.3 $\pm$ 0.38	55.12 $\pm$ 0.12	26.48 $\pm$ 0.64	68.72 $\pm$ 0.51	19.0 $\pm$ 0.42	46.48 $\pm$ 0.1	33.33 $\pm$ 14.29
	0.75	23.93 $\pm$ 0.1	17.58 $\pm$ 0.15	29.26 $\pm$ 0.38	<b>63.13</b> $\pm$ 0.22	38.89 $\pm$ 0.07	54.49 $\pm$ 0.74	55.43 $\pm$ 0.06	26.54 $\pm$ 0.5	68.72 $\pm$ 0.51	19.8 $\pm$ 0.5	46.71 $\pm$ 0.12	33.33 $\pm$ 17.17
	0.9	<b>23.77</b> $\pm$ 0.06	<b>17.48</b> $\pm$ 0.17	<b>29.11</b> $\pm$ 0.38	62.81 $\pm$ 0.4	39.1 $\pm$ 0.1	55.59 $\pm$ 0.27	<b>55.71</b> $\pm$ 0.28	26.31 $\pm$ 0.32	<b>68.81</b> $\pm$ 0.29	20.13 $\pm$ 0.57	46.92 $\pm$ 0.09	38.1 $\pm$ 19.05
1.0	27.03 $\pm$ 0.29	19.85 $\pm$ 0.21	32.85 $\pm$ 0.24	62.35 $\pm$ 1.39	<b>39.12</b> $\pm$ 0.17	55.75 $\pm$ 0.43	55.2 $\pm$ 0.6	26.54 $\pm$ 0.3	<b>68.81</b> $\pm$ 0.23	<b>21.67</b> $\pm$ 0.29	<b>47.06</b> $\pm$ 0.14	-	
0.5 (OWL)	0.0	21.44 $\pm$ 0.21	21.25 $\pm$ 0.1	34.47 $\pm$ 0.26	61.41 $\pm$ 0.56	37.54 $\pm$ 0.08	55.3 $\pm$ 0.34	51.54 $\pm$ 0.43	25.31 $\pm$ 0.4	67.46 $\pm$ 0.11	18.6 $\pm$ 1.03	45.31 $\pm$ 0.2	0.0 $\pm$ 0.0
	0.1	21.13 $\pm$ 0.07	18.29 $\pm$ 0.11	30.68 $\pm$ 0.46	62.68 $\pm$ 0.3	39.11 $\pm$ 0.07	56.35 $\pm$ 0.61	55.39 $\pm$ 0.82	26.25 $\pm$ 0.36	67.92 $\pm$ 0.15	20.53 $\pm$ 0.68	46.89 $\pm$ 0.29	23.81 $\pm$ 17.17
	0.25	23.84 $\pm$ 0.08	18.02 $\pm$ 0.15	29.81 $\pm$ 0.27	62.76 $\pm$ 0.2	39.15 $\pm$ 0.06	56.38 $\pm$ 0.16	55.32 $\pm$ 0.67	26.96 $\pm$ 0.09	68.46 $\pm$ 0.14	20.8 $\pm$ 0.2	47.12 $\pm$ 0.11	38.1 $\pm$ 12.6
	0.5	23.57 $\pm$ 0.07	17.77 $\pm$ 0.1	29.69 $\pm$ 0.13	62.77 $\pm$ 0.45	39.4 $\pm$ 0.11	55.99 $\pm$ 0.47	55.44 $\pm$ 0.59	26.93 $\pm$ 0.49	<b>68.63</b> $\pm$ 0.28	20.27 $\pm$ 0.37	47.06 $\pm$ 0.27	42.86 $\pm$ 14.29
	0.75	23.41 $\pm$ 0.05	17.69 $\pm$ 0.09	29.38 $\pm$ 0.19	<b>63.12</b> $\pm$ 0.5	39.46 $\pm$ 0.05	56.8 $\pm$ 0.73	55.88 $\pm$ 0.23	<b>27.45</b> $\pm$ 0.28	68.48 $\pm$ 0.27	20.67 $\pm$ 0.47	<b>47.41</b> $\pm$ 0.22	<b>57.14</b> $\pm$ 14.29
	0.9	<b>23.31</b> $\pm$ 0.04	<b>17.58</b> $\pm$ 0.11	<b>29.16</b> $\pm$ 0.31	62.88 $\pm$ 0.34	39.56 $\pm$ 0.1	<b>56.96</b> $\pm$ 0.78	<b>56.05</b> $\pm$ 0.35	27.08 $\pm$ 0.12	68.53 $\pm$ 0.13	20.27 $\pm$ 0.35	47.33 $\pm$ 0.16	52.38 $\pm$ 12.6
1.0	26.25 $\pm$ 0.05	19.74 $\pm$ 0.08	32.2 $\pm$ 0.26	62.86 $\pm$ 0.31	<b>39.94</b> $\pm$ 0.08	56.59 $\pm$ 0.55	55.36 $\pm$ 0.32	26.42 $\pm$ 0.36	68.32 $\pm$ 0.22	<b>21.93</b> $\pm$ 0.37	47.37 $\pm$ 0.15	-	
0.6	0.0	85.09 $\pm$ 1.34	72.05 $\pm$ 0.81	104.4 $\pm$ 1.31	60.89 $\pm$ 1.43	29.47 $\pm$ 0.16	53.17 $\pm$ 1.07	39.46 $\pm$ 0.42	19.51 $\pm$ 0.33	59.97 $\pm$ 0.11	13.93 $\pm$ 0.74	43.99 $\pm$ 0.36	14.29 $\pm$ 0.0
	0.1	56.05 $\pm$ 1.38	44.2 $\pm$ 0.94	67.85 $\pm$ 2.09	61.47 $\pm$ 0.23	31.68 $\pm$ 0.11	52.64 $\pm$ 0.39	44.99 $\pm$ 0.81	21.42 $\pm$ 0.52	62.44 $\pm$ 0.37	15.47 $\pm$ 0.84	41.44 $\pm$ 0.18	38.1 $\pm$ 12.6
	0.25	54.15 $\pm$ 2.38	42.06 $\pm$ 1.34	63.21 $\pm$ 1.81	62.15 $\pm$ 0.2	31.99 $\pm$ 0.12	52.57 $\pm$ 0.39	45.19 $\pm$ 0.52	21.33 $\pm$ 0.2	62.6 $\pm$ 0.18	14.87 $\pm$ 0.59	41.53 $\pm$ 0.1	47.62 $\pm$ 12.6
	0.5	51.97 $\pm$ 1.64	40.32 $\pm$ 1.13	61.81 $\pm$ 2.66	62.19 $\pm$ 0.07	32.13 $\pm$ 0.18	52.91 $\pm$ 0.45	46.34 $\pm$ 0.47	21.3 $\pm$ 0.23	62.88 $\pm$ 0.09	15.87 $\pm$ 0.24	41.94 $\pm$ 0.16	52.38 $\pm$ 9.52
	0.75	51.28 $\pm$ 1.68	39.69 $\pm$ 1.27	60.91 $\pm$ 2.28	62.16 $\pm$ 0.04	32.38 $\pm$ 0.19	53.38 $\pm$ 0.32	46.25 $\pm$ 0.34	<b>21.53</b> $\pm$ 0.21	<b>63.31</b> $\pm$ 0.13	15.67 $\pm$ 0.29	<b>42.1</b> $\pm$ 0.13	<b>57.14</b> $\pm$ 8.25
	0.9	<b>50.28</b> $\pm$ 1.99	<b>39.13</b> $\pm$ 1.54	<b>60.14</b> $\pm$ 2.9	<b>62.36</b> $\pm$ 0.12	<b>32.49</b> $\pm$ 0.13	53.09 $\pm$ 0.18	<b>46.49</b> $\pm$ 0.38	21.3 $\pm$ 0.24	63.22 $\pm$ 0.17	15.73 $\pm$ 0.55	<b>42.1</b> $\pm$ 0.11	<b>57.14</b> $\pm$ 8.25
1.0	63.63 $\pm$ 1.18	54.16 $\pm$ 1.0	81.11 $\pm$ 3.99	60.67 $\pm$ 0.59	32.16 $\pm$ 0.2	<b>54.46</b> $\pm$ 0.53	<b>44.94</b> $\pm$ 0.11	21.47 $\pm$ 0.48	62.21 $\pm$ 0.2	<b>17.07</b> $\pm$ 0.41	41.85 $\pm$ 0.2	-	
0.6 (Alpha-Pruning)	0.0	88.67 $\pm$ 2.38	75.35 $\pm$ 3.09	109.66 $\pm$ 5.07	61.52 $\pm$ 0.51	29.5 $\pm$ 0.03	51.64 $\pm$ 0.47	38.61 $\pm$ 0.43	19.48 $\pm$ 0.45	59.5 $\pm$ 0.27	12.67 $\pm$ 0.47	38.99 $\pm$ 0.28	4.76 $\pm$ 4.76
	0.1	57.24 $\pm$ 0.41	45.06 $\pm$ 0.23	68.58 $\pm$ 0.08	61.88 $\pm$ 0.24	31.64 $\pm$ 0.07	51.8 $\pm$ 0.38	44.89 $\pm$ 0.87	21.13 $\pm$ 0.34	61.29 $\pm$ 0.41	15.73 $\pm$ 0.59	41.29 $\pm$ 0.33	19.05 $\pm$ 9.52
	0.25	53.93 $\pm$ 0.37	42.36 $\pm$ 0.2	64.63 $\pm$ 1.4	62.13 $\pm$ 0.04	31.89 $\pm$ 0.07	52.7 $\pm$ 0.6	45.23 $\pm$ 0.44	21.76 $\pm$ 0.15	61.82 $\pm$ 0.19	15.87 $\pm$ 0.64	41.63 $\pm$ 0.22	28.57 $\pm$ 0.0
	0.5	51.22 $\pm$ 0.44	40.66 $\pm$ 0.49	62.51 $\pm$ 1.04	62.22 $\pm$ 0.08	32.16 $\pm$ 0.06	53.38 $\pm$ 0.74	45.74 $\pm$ 0.58	22.07 $\pm$ 0.37	62.79 $\pm$ 0.17	15.93 $\pm$ 0.13	42.04 $\pm$ 0.13	57.14 $\pm$ 14.29
	0.75	50.52 $\pm$ 0.55	39.44 $\pm$ 0.3	61.78 $\pm$ 0.85	62.06 $\pm$ 0.22	32.38 $\pm$ 0.07	53.33 $\pm$ 0.39	<b>46.68</b> $\pm$ 0.6	<b>22.21</b> $\pm$ 0.22	62.66 $\pm$ 0.16	16.13 $\pm$ 0.18	42.21 $\pm$ 0.21	61.9 $\pm$ 4.76
	0.9	<b>49.31</b> $\pm$ 1.1	<b>38.44</b> $\pm$ 0.52	<b>60.32</b> $\pm$ 1.2	<b>62.29</b> $\pm$ 0.05	<b>32.53</b> $\pm$ 0.06	53.7 $\pm$ 0.55	46.3 $\pm$ 0.3	21.99 $\pm$ 0.15	<b>63.13</b> $\pm$ 0.1	16.6 $\pm$ 0.12	<b>42.36</b> $\pm$ 0.14	<b>66.67</b> $\pm$ 9.52
1.0	61.05 $\pm$ 0.77	52.8 $\pm$ 0.43	78.27 $\pm$ 3.64	62.08 $\pm$ 0.26	32.0 $\pm$ 0.08	<b>53.88</b> $\pm$ 0.25	45.29 $\pm$ 0.49	22.01 $\pm$ 0.59	62.02 $\pm$ 0.42	<b>17.6</b> $\pm$ 0.42	42.13 $\pm$ 0.06	-	
0.6 (OWL)	0.0	74.97 $\pm$ 1.68	64.7 $\pm$ 1.38	94.32 $\pm$ 3.39	62.11 $\pm$ 0.11	30.75 $\pm$ 0.2	52.07 $\pm$ 0.57	40.39 $\pm$ 0.1	21.05 $\pm$ 0.37	61.08 $\pm$ 0.18	14.2 $\pm$ 0.5	40.23 $\pm$ 0.08	4.76 $\pm$ 4.76
	0.1	48.43 $\pm$ 0.56	40.09 $\pm$ 0.31	59.81 $\pm$ 1.44	62.22 $\pm$ 0.05	32.92 $\pm$ 0.08	52.59 $\pm$ 0.22	45.69 $\pm$ 0.51	22.04 $\pm$ 0.15	63.58 $\pm$ 0.13	16.13 $\pm$ 0.52	42.17 $\pm$ 0.13	42.86 $\pm$ 8.25
	0.25	46.31 $\pm$ 0.78	38.41 $\pm$ 0.17	56.18 $\pm$ 0.88	<b>62.27</b> $\pm$ 0.1	33.06 $\pm$ 0.06	52.99 $\pm$ 0.51	45.17 $\pm$ 0.32	22.64 $\pm$ 0.42	64.02 $\pm$ 0.12	15.8 $\pm$ 0.61	42.28 $\pm$ 0.11	61.9 $\pm$ 9.52
	0.5	45.28 $\pm$ 1.06	37.36 $\pm$ 0.54	54.54 $\pm$ 1.37	62.21 $\pm$ 0.02	33.25 $\pm$ 0.2	53.62 $\pm$ 0.38	45.99 $\pm$ 0.5	23.07 $\pm$ 0.79	<b>64.2</b> $\pm$ 0.14	16.6 $\pm$ 0.12	42.71 $\pm$ 0.26	<b>71.43</b> $\pm$ 16.5
	0.75	44.3 $\pm$ 0.89	36.24 $\pm$ 0.24	53.87 $\pm$ 1.29	<b>62.27</b> $\pm$ 0.09	33.48 $\pm$ 0.16	<b>54.06</b> $\pm$ 0.44	<b>46.89</b> $\pm$ 0.15					

Table 16: Test perplexity on C4, WikiText2 and PTB and zero-shot accuracy of Llama-3.2-1B using *MOON-SHOT*-Wanda. Zero-shot mean performance and win rate are computed over the 7 classification tasks. Results are averaged over 3 seeds with standard errors.

Sparsity	$\lambda$	C4 ↓	WikiText2 ↓	PTB ↓	BoolQ ↑	HellaSwag ↑	WinoGrande ↑	ARC-c ↑	ARC-c ↑	PIQA ↑	OBQA ↑	Mean ↑	Win Rate ↑
Dense	-	14.02	9.75	17.59	64.01	47.73	60.14	65.15	31.23	74.32	26.4	52.71	-
0.5	0.0	30.48 ± 0.14	21.48 ± 0.08	39.41 ± 0.22	61.59 ± 0.37	35.57 ± 0.08	54.14 ± 0.6	54.32 ± 0.28	24.94 ± 0.12	65.96 ± 0.21	17.8 ± 0.7	44.9 ± 0.21	80.95 ± 4.76
	0.1	30.26 ± 0.17	21.32 ± 0.12	38.76 ± 0.18	61.85 ± 0.22	35.59 ± 0.13	54.51 ± 0.11	54.66 ± 0.28	24.94 ± 0.1	66.36 ± 0.02	17.93 ± 0.7	45.12 ± 0.17	80.95 ± 4.76
	0.25	30.08 ± 0.09	21.1 ± 0.07	38.19 ± 0.04	61.98 ± 0.19	35.68 ± 0.13	53.96 ± 0.3	54.7 ± 0.44	25.2 ± 0.15	66.29 ± 0.1	18.0 ± 1.01	45.11 ± 0.14	80.95 ± 4.76
	0.5	29.82 ± 0.04	20.81 ± 0.04	37.68 ± 0.25	61.99 ± 0.27	35.91 ± 0.11	54.14 ± 0.6	54.69 ± 0.18	25.2 ± 0.41	<b>66.56 ± 0.13</b>	<b>18.27 ± 0.44</b>	<b>45.25 ± 0.19</b>	<b>85.71 ± 8.25</b>
	0.75	<b>29.63 ± 0.04</b>	<b>20.69 ± 0.08</b>	37.32 ± 0.16	62.06 ± 0.18	35.95 ± 0.06	54.46 ± 0.05	54.45 ± 0.21	<b>25.31 ± 0.12</b>	66.2 ± 0.19	17.8 ± 0.23	45.18 ± 0.08	71.43 ± 0.0
	1.0	29.7 ± 0.08	<b>20.69 ± 0.05</b>	<b>37.23 ± 0.1</b>	<b>62.13 ± 0.13</b>	<b>36.16 ± 0.06</b>	<b>54.59 ± 0.21</b>	54.31 ± 0.15	25.26 ± 0.31	66.34 ± 0.11	17.4 ± 0.42	45.17 ± 0.1	80.95 ± 4.76
0.9	35.71 ± 0.21	24.43 ± 0.04	43.15 ± 0.26	60.23 ± 0.49	35.26 ± 0.04	54.56 ± 0.03	51.59 ± 0.29	24.69 ± 0.21	65.51 ± 0.13	18.2 ± 0.2	44.29 ± 0.07	-	
0.5 (Alpha-Pruning)	0.0	30.46 ± 0.16	21.25 ± 0.11	38.88 ± 0.19	61.15 ± 0.09	35.64 ± 0.07	54.91 ± 0.39	53.62 ± 0.09	24.32 ± 0.13	66.29 ± 0.23	17.73 ± 0.24	44.81 ± 0.09	80.95 ± 4.76
	0.1	30.28 ± 0.13	21.1 ± 0.08	38.24 ± 0.19	<b>61.78 ± 0.08</b>	35.67 ± 0.09	54.33 ± 0.32	53.87 ± 0.17	24.12 ± 0.25	66.3 ± 0.07	17.8 ± 0.72	44.84 ± 0.18	71.43 ± 8.25
	0.25	30.24 ± 0.16	21.01 ± 0.05	38.06 ± 0.08	61.67 ± 0.15	35.74 ± 0.11	54.35 ± 0.26	<b>54.12 ± 0.22</b>	24.32 ± 0.13	66.23 ± 0.07	<b>18.0 ± 0.7</b>	44.92 ± 0.15	80.95 ± 4.76
	0.5	29.93 ± 0.18	20.74 ± 0.09	37.38 ± 0.08	61.58 ± 0.2	35.79 ± 0.08	54.75 ± 0.13	54.05 ± 0.09	24.26 ± 0.16	66.61 ± 0.19	<b>18.0 ± 0.61</b>	45.01 ± 0.09	80.95 ± 4.76
	0.75	<b>29.57 ± 0.03</b>	20.5 ± 0.04	<b>37.03 ± 0.1</b>	61.64 ± 0.13	35.92 ± 0.01	<b>55.46 ± 0.48</b>	53.98 ± 0.15	24.63 ± 0.21	66.65 ± 0.11	<b>18.0 ± 0.31</b>	<b>45.18 ± 0.11</b>	<b>85.71 ± 8.25</b>
	1.0	29.64 ± 0.06	<b>20.48 ± 0.07</b>	37.13 ± 0.14	61.41 ± 0.13	<b>36.08 ± 0.03</b>	55.38 ± 0.35	53.7 ± 0.08	<b>24.66 ± 0.21</b>	<b>66.83 ± 0.04</b>	17.93 ± 0.13	45.14 ± 0.03	<b>90.48 ± 4.76</b>
0.9	35.47 ± 0.21	23.99 ± 0.11	42.65 ± 0.3	58.79 ± 0.37	35.31 ± 0.13	55.09 ± 0.32	51.07 ± 0.08	24.57 ± 0.31	65.14 ± 0.19	17.27 ± 0.24	43.89 ± 0.12	-	
0.5 (OWL)	0.0	28.88 ± 0.09	21.1 ± 0.1	38.01 ± 0.37	61.68 ± 0.12	36.48 ± 0.04	54.54 ± 0.51	52.95 ± 0.43	24.32 ± 0.13	66.16 ± 0.05	19.2 ± 0.2	45.05 ± 0.14	57.14 ± 8.25
	0.1	28.8 ± 0.09	20.98 ± 0.09	37.79 ± 0.41	61.82 ± 0.2	36.52 ± 0.11	55.12 ± 0.43	52.92 ± 0.56	24.43 ± 0.32	66.23 ± 0.29	19.27 ± 0.24	45.19 ± 0.17	76.19 ± 4.76
	0.25	28.82 ± 0.12	20.9 ± 0.1	37.48 ± 0.36	<b>62.08 ± 0.09</b>	36.66 ± 0.06	54.75 ± 0.07	<b>53.3 ± 0.56</b>	24.74 ± 0.09	<b>66.52 ± 0.16</b>	19.53 ± 0.29	45.37 ± 0.09	76.19 ± 12.6
	0.5	28.7 ± 0.13	20.71 ± 0.13	36.85 ± 0.42	61.96 ± 0.2	36.79 ± 0.05	54.78 ± 0.36	53.04 ± 0.21	<b>25.11 ± 0.21</b>	66.38 ± 0.25	19.93 ± 0.07	45.43 ± 0.06	85.71 ± 8.25
	0.75	28.51 ± 0.04	<b>20.58 ± 0.07</b>	36.22 ± 0.29	61.98 ± 0.02	36.76 ± 0.04	55.25 ± 0.44	53.24 ± 0.13	24.69 ± 0.38	66.49 ± 0.14	19.87 ± 0.07	45.47 ± 0.11	85.71 ± 8.25
	1.0	28.52 ± 0.03	20.59 ± 0.02	<b>36.18 ± 0.12</b>	61.82 ± 0.32	<b>37.0 ± 0.06</b>	<b>55.56 ± 0.45</b>	52.89 ± 0.18	24.89 ± 0.29	66.36 ± 0.2	<b>20.27 ± 0.29</b>	<b>45.54 ± 0.15</b>	<b>95.24 ± 4.76</b>
0.9	33.38 ± 0.16	23.16 ± 0.06	40.06 ± 0.08	60.01 ± 0.77	36.63 ± 0.05	54.83 ± 0.4	50.53 ± 0.26	24.66 ± 0.38	65.42 ± 0.1	19.47 ± 0.29	44.51 ± 0.12	-	
0.6	0.0	88.87 ± 1.62	66.78 ± 1.37	108.78 ± 4.4	<b>61.99 ± 0.11</b>	29.39 ± 0.07	51.14 ± 0.3	39.86 ± 0.39	19.0 ± 0.42	60.52 ± 0.14	13.27 ± 0.35	39.31 ± 0.12	80.95 ± 12.6
	0.1	87.45 ± 2.13	65.41 ± 1.82	106.65 ± 4.21	61.78 ± 0.22	29.47 ± 0.03	<b>52.35 ± 0.53</b>	40.31 ± 0.35	19.2 ± 0.15	60.23 ± 0.33	13.67 ± 0.24	39.57 ± 0.06	80.95 ± 12.6
	0.25	86.94 ± 2.11	64.92 ± 2.09	103.22 ± 3.32	61.9 ± 0.11	29.56 ± 0.04	51.72 ± 0.59	40.19 ± 0.13	19.37 ± 0.13	60.66 ± 0.16	13.73 ± 0.37	39.59 ± 0.06	<b>85.71 ± 8.25</b>
	0.5	<b>86.55 ± 1.67</b>	<b>63.57 ± 1.61</b>	98.44 ± 3.98	61.56 ± 0.29	29.53 ± 0.06	51.64 ± 0.23	40.4 ± 0.17	<b>19.6 ± 0.06</b>	<b>61.12 ± 0.08</b>	13.4 ± 0.2	39.61 ± 0.07	80.95 ± 4.76
	0.75	87.41 ± 0.53	63.83 ± 0.69	<b>96.72 ± 2.62</b>	61.55 ± 0.37	29.57 ± 0.04	52.33 ± 0.76	40.95 ± 0.15	19.03 ± 0.09	60.9 ± 0.07	13.4 ± 0.23	39.67 ± 0.12	76.19 ± 9.52
	1.0	90.04 ± 0.69	64.72 ± 0.13	97.73 ± 0.94	61.64 ± 0.21	<b>29.64 ± 0.04</b>	52.17 ± 0.44	<b>41.11 ± 0.38</b>	19.06 ± 0.1	61.1 ± 0.11	13.33 ± 0.37	<b>39.72 ± 0.07</b>	80.95 ± 12.6
0.9	117.71 ± 0.87	84.73 ± 0.73	119.64 ± 1.0	58.96 ± 1.39	28.86 ± 0.03	51.35 ± 0.49	38.82 ± 0.32	18.94 ± 0.26	59.05 ± 0.18	<b>13.93 ± 0.24</b>	38.56 ± 0.12	-	
0.6 (Alpha-Pruning)	0.0	87.66 ± 1.09	66.36 ± 0.7	107.17 ± 2.93	<b>61.94 ± 0.15</b>	29.34 ± 0.1	51.49 ± 0.39	39.06 ± 0.38	19.62 ± 0.31	59.99 ± 0.07	12.4 ± 0.53	39.12 ± 0.08	71.43 ± 0.0
	0.1	86.35 ± 1.64	65.13 ± 1.22	104.51 ± 3.58	61.79 ± 0.29	29.49 ± 0.12	51.7 ± 0.34	39.16 ± 0.7	19.74 ± 0.08	60.28 ± 0.06	12.4 ± 0.5	39.22 ± 0.13	<b>85.71 ± 8.25</b>
	0.25	<b>84.21 ± 1.04</b>	63.3 ± 0.92	101.96 ± 2.98	61.69 ± 0.24	29.56 ± 0.08	<b>52.33 ± 0.37</b>	39.28 ± 0.29	19.65 ± 0.06	60.32 ± 0.13	12.07 ± 0.58	39.27 ± 0.09	<b>85.71 ± 8.25</b>
	0.5	84.9 ± 0.94	62.75 ± 1.05	98.54 ± 2.09	61.86 ± 0.21	<b>29.67 ± 0.09</b>	51.93 ± 0.18	39.96 ± 0.26	19.74 ± 0.4	60.39 ± 0.27	12.13 ± 0.37	39.38 ± 0.04	76.19 ± 4.76
	0.75	85.33 ± 0.55	<b>62.02 ± 0.85</b>	<b>97.49 ± 1.57</b>	61.46 ± 0.25	29.68 ± 0.1	51.83 ± 0.21	40.07 ± 0.49	19.43 ± 0.23	60.83 ± 0.05	12.47 ± 0.13	39.38 ± 0.02	71.43 ± 0.0
	1.0	88.49 ± 0.32	63.13 ± 0.36	100.53 ± 0.5	61.51 ± 0.27	<b>29.67 ± 0.1</b>	51.75 ± 0.09	<b>40.32 ± 0.42</b>	<b>19.88 ± 0.05</b>	<b>60.95 ± 0.12</b>	12.27 ± 0.13	<b>39.48 ± 0.06</b>	80.95 ± 4.76
0.9	112.33 ± 1.03	80.75 ± 1.03	120.89 ± 0.73	58.01 ± 1.1	28.92 ± 0.09	51.22 ± 0.47	37.56 ± 0.16	19.51 ± 0.21	59.32 ± 0.04	<b>13.4 ± 0.4</b>	38.28 ± 0.1	-	
0.6 (OWL)	0.0	76.14 ± 0.64	62.0 ± 0.52	104.62 ± 1.6	61.64 ± 0.17	30.36 ± 0.1	51.49 ± 0.54	40.32 ± 0.19	20.53 ± 0.3	60.23 ± 0.16	13.2 ± 0.46	39.68 ± 0.06	57.14 ± 8.25
	0.1	74.55 ± 0.73	60.11 ± 0.58	101.96 ± 0.17	61.72 ± 0.09	30.42 ± 0.11	52.25 ± 0.16	40.7 ± 0.23	20.34 ± 0.31	60.5 ± 0.03	13.6 ± 0.31	39.93 ± 0.13	61.9 ± 9.52
	0.25	<b>73.97 ± 0.19</b>	58.81 ± 0.28	100.07 ± 1.96	<b>61.81 ± 0.08</b>	30.47 ± 0.07	51.51 ± 0.37	40.92 ± 0.22	20.71 ± 0.15	60.36 ± 0.19	13.27 ± 0.35	39.86 ± 0.15	66.67 ± 4.76
	0.5	74.58 ± 0.49	58.75 ± 0.37	97.11 ± 0.77	61.72 ± 0.04	30.58 ± 0.11	52.01 ± 0.12	40.67 ± 0.28	20.31 ± 0.09	60.39 ± 0.08	13.8 ± 0.42	39.93 ± 0.13	52.38 ± 9.52
	0.75	74.74 ± 0.47	<b>57.56 ± 0.17</b>	<b>94.12 ± 1.37</b>	61.75 ± 0.16	30.58 ± 0.11	52.8 ± 0.28	40.81 ± 0.23	20.19 ± 0.17	60.68 ± 0.15	13.93 ± 0.48	40.11 ± 0.11	<b>71.43 ± 8.25</b>
	1.0	76.62 ± 0.57	58.53 ± 0.42	94.61 ± 0.92	61.73 ± 0.14	<b>30.66 ± 0.09</b>	<b>52.99 ± 0.39</b>	<b>41.18 ± 0.25</b>	20.59 ± 0.27	<b>61.08 ± 0.24</b>	13.87 ± 0.35	<b>40.3 ± 0.19</b>	<b>71.43 ± 0.0</b>
0.9	99.38 ± 1.37	73.0 ± 0.37	111.24 ± 1.83	61.28 ± 0.28	29.88 ± 0.11	52.07 ± 0.82	40.22 ± 0.09	<b>20.82 ± 0.05</b>	60.03 ± 0.2	<b>14.8 ± 0.12</b>	39.87 ± 0.17	-	
0.7	0.0	363.68 ± 3.74	393.51 ± 5.96	459.34 ± 12.4	38.81 ± 0.41	26.85 ± 0.04	49.14 ± 0.37	29.69 ± 0.24	18.86 ± 0.1	55.39 ± 0.25	12.4 ± 0.42	33.02 ± 0.07	52.38 ± 9.52
	0.1	354.91 ± 6.53	404.15 ± 1.97	438.68 ± 18.66	38.99 ± 0.61	26.9 ± 0.03	50.22 ± 0.62	29.87 ± 0.24	18.66 ± 0.25	55.6 ± 0.26	12.53 ± 0.07	33.25 ± 0.16	52.38 ± 17.17
	0.25	<b>342.24 ± 11.1</b>	370.67 ± 5.15	419.8 ± 37.86	39.01 ± 0.83	26.89 ± 0.09	50.38 ± 0.94	29.92 ± 0.13	18.57 ± 0.3	<b>55.79 ± 0.24</b>	12.2 ± 0.23	33.25 ± 0.29	57.14 ± 8.25
	0.5	356.93 ± 11.89	377.15 ± 1.18	415.21 ± 15.26	38.83 ± 0.59	26.95 ± 0.01	50.57 ± 0.39	<b>29.99 ± 0.17</b>	<b>19.14 ± 0.08</b>	55.64 ± 0.12	12.27 ± 0.18	33.34 ± 0.06	<b>61.9 ± 4.76</b>
	0.75	357.83 ± 11.36	<b>358.07 ± 2.4</b>	<b>384.48 ± 10.82</b>	38.84 ± 0.17	<b>26.98 ± 0.05</b>	51.62 ± 0.46	29.55 ± 0.13	18.91 ± 0.03	55.73 ± 0.21	12.53 ± 0.55	<b>33.45 ± 0.1</b>	57.14 ± 8.25
	1.0	383.24 ± 9.61	380.41 ± 5.38	420.72 ± 13.07	38.13 ± 0.06	26.93 ± 0.05	<b>51.85 ± 0.34</b>	29.81 ± 0.27	19.03 ± 0.21	55.46 ± 0.13	12.07 ± 0.52	33.33 ± 0.06	<b>61.9 ± 4.76</b>
0.9	730.58 ± 21.78	777.79 ± 20.03	1130.38 ± 57.74	<b>39.91 ± 0.83</b>	26.52 ± 0.05	50.86 ± 0.41	29.32 ± 0.17	18.86 ± 0.32	55.28 ± 0.28	<b>12.93 ± 0.44</b>	33.38 ± 0.15	-	
0.7 (Alpha-Pruning)	0.0	<b>372.76 ± 4.51</b>	399.6 ± 23.33	448.03 ± 16.78	40.22 ± 0.43	26.76 ± 0.04	48.93 ± 0.32	<b>29.81 ± 0.17</b>	18.6 ± 0.09	55.28 ± 0.19	11.93 ± 0.18	33.08 ± 0.07	47.62 ± 4.76
	0.1	380.49 ± 8.11	410.81 ± 13.71	467.38 ± 22.23	40.04 ± 1.0	26.78 ± 0.06	49.43 ± 0.63	29.66 ± 0.04	18.77 ± 0.05	55.08 ± 0.07	11.67 ± 0.18	33.06 ± 0.11	52.38 ± 9.52
	0.25	378.98 ± 10.46	395.64 ± 3.8	456.96 ± 32.25	39.91 ± 0.87	26.76 ± 0.08	48.91 ± 0.3	29.67 ± 0.21	<b>19.17 ± 0.25</b>	55.19 ± 0.27	12.2 ± 0.6	33.12 ± 0.19	<b>61.9 ± 4.76</b>
	0.5	388.09 ± 11.97	405.85 ± 6.13	426.49 ± 25.19	39.11 ± 0.19	26.8 ± 0.07							

Table 17: Test perplexity on C4, WikiText2 and PTB and zero-shot accuracy of Llama-3.2-3B using *MOON-SHOT-OSSCAR*. Zero-shot mean performance and win rate are computed over the 7 classification tasks. Results are averaged over 3 seeds with standard errors.

Sparsity	$\lambda$	C4 $\downarrow$	WikiText2 $\downarrow$	PTB $\downarrow$	BoolQ $\uparrow$	HellaSwag $\uparrow$	WinoGrande $\uparrow$	ARC-e $\uparrow$	ARC-c $\uparrow$	PIQA $\uparrow$	OBQA $\uparrow$	Mean $\uparrow$	Win Rate $\uparrow$
Dense	-	11.34	7.81	13.54	72.72	55.30	69.22	74.37	42.41	76.71	31.20	60.28	-
0.1	0.0	19.73 $\pm$ 0.91	15.94 $\pm$ 0.42	51.85 $\pm$ 5.02	60.2 $\pm$ 2.06	45.85 $\pm$ 0.78	60.59 $\pm$ 0.97	64.8 $\pm$ 1.09	31.85 $\pm$ 0.62	73.94 $\pm$ 0.67	25.0 $\pm$ 0.42	51.75 $\pm$ 0.7	9.52 $\pm$ 4.76
	0.1	18.86 $\pm$ 1.07	14.93 $\pm$ 0.6	50.43 $\pm$ 9.15	<b>65.4</b> $\pm$ 1.45	<b>51.07</b> $\pm$ 0.51	<b>62.4</b> $\pm$ 1.71	<b>67.86</b> $\pm$ 0.62	34.7 $\pm$ 0.53	74.86 $\pm$ 0.37	24.13 $\pm$ 1.43	<b>54.35</b> $\pm$ 0.84	57.14 $\pm$ 21.82
	0.25	18.5 $\pm$ 0.24	13.47 $\pm$ 0.73	27.07 $\pm$ 0.13	64.37 $\pm$ 0.98	50.83 $\pm$ 0.28	60.77 $\pm$ 1.78	66.89 $\pm$ 0.4	<b>35.13</b> $\pm$ 0.38	75.66 $\pm$ 0.52	22.33 $\pm$ 2.1	53.71 $\pm$ 0.74	61.9 $\pm$ 4.76
	0.5	<b>16.42</b> $\pm$ 0.22	<b>12.63</b> $\pm$ 0.33	23.22 $\pm$ 0.42	65.06 $\pm$ 1.18	50.36 $\pm$ 0.3	60.41 $\pm$ 1.79	66.23 $\pm$ 0.94	34.67 $\pm$ 0.78	75.55 $\pm$ 0.67	25.67 $\pm$ 0.74	53.99 $\pm$ 0.83	<b>66.67</b> $\pm$ 20.76
	0.75	16.56 $\pm$ 0.38	13.28 $\pm$ 0.55	<b>22.1</b> $\pm$ 0.28	64.91 $\pm$ 1.0	49.5 $\pm$ 0.62	61.43 $\pm$ 0.47	66.39 $\pm$ 0.83	34.41 $\pm$ 0.33	75.54 $\pm$ 0.45	26.8 $\pm$ 0.76	54.14 $\pm$ 0.2	<b>66.67</b> $\pm$ 4.76
	0.9	16.57 $\pm$ 0.41	13.58 $\pm$ 0.8	22.16 $\pm$ 0.24	64.36 $\pm$ 1.07	49.21 $\pm$ 0.9	61.09 $\pm$ 0.4	66.47 $\pm$ 0.93	34.36 $\pm$ 0.52	75.5 $\pm$ 0.56	27.07 $\pm$ 0.7	54.01 $\pm$ 0.23	61.9 $\pm$ 4.76
1.0	16.8 $\pm$ 0.46	13.76 $\pm$ 0.85	22.34 $\pm$ 0.43	64.38 $\pm$ 0.86	49.03 $\pm$ 0.86	61.3 $\pm$ 0.74	65.82 $\pm$ 0.64	34.1 $\pm$ 0.53	<b>75.7</b> $\pm$ 0.78	<b>27.33</b> $\pm$ 0.41	53.95 $\pm$ 0.08	0.0 $\pm$ 0.0	
0.15	0.0	25.74 $\pm$ 1.05	23.46 $\pm$ 0.84	48.04 $\pm$ 3.18	50.73 $\pm$ 1.87	38.44 $\pm$ 1.87	56.3 $\pm$ 0.71	61.53 $\pm$ 0.87	29.01 $\pm$ 0.67	72.25 $\pm$ 0.52	20.4 $\pm$ 1.55	46.95 $\pm$ 0.55	14.29 $\pm$ 8.25
	0.1	22.26 $\pm$ 1.04	18.81 $\pm$ 2.21	55.04 $\pm$ 19.76	<b>64.46</b> $\pm$ 1.75	<b>48.69</b> $\pm$ 0.24	<b>60.69</b> $\pm$ 1.21	<b>65.75</b> $\pm$ 0.59	32.65 $\pm$ 0.92	73.81 $\pm$ 0.57	21.87 $\pm$ 1.38	52.56 $\pm$ 0.85	76.19 $\pm$ 9.52
	0.25	<b>21.5</b> $\pm$ 0.22	<b>18.28</b> $\pm$ 0.41	<b>27.99</b> $\pm$ 1.03	64.34 $\pm$ 2.28	48.37 $\pm$ 0.21	58.77 $\pm$ 1.44	65.4 $\pm$ 0.71	<b>33.16</b> $\pm$ 0.67	<b>74.77</b> $\pm$ 0.55	<b>25.0</b> $\pm$ 1.15	<b>52.83</b> $\pm$ 0.84	<b>95.24</b> $\pm$ 4.76
	0.5	22.17 $\pm$ 0.48	20.19 $\pm$ 1.46	30.22 $\pm$ 0.2	63.65 $\pm$ 2.03	47.74 $\pm$ 0.53	57.98 $\pm$ 1.77	64.41 $\pm$ 1.03	32.54 $\pm$ 0.69	74.52 $\pm$ 0.68	23.93 $\pm$ 1.16	52.11 $\pm$ 1.02	90.48 $\pm$ 4.76
	0.75	22.26 $\pm$ 0.65	21.1 $\pm$ 1.9	30.73 $\pm$ 0.23	63.72 $\pm$ 2.37	47.42 $\pm$ 0.78	58.43 $\pm$ 1.76	64.28 $\pm$ 1.04	31.68 $\pm$ 0.85	73.76 $\pm$ 0.65	24.33 $\pm$ 0.94	51.95 $\pm$ 1.07	85.71 $\pm$ 0.0
	0.9	22.32 $\pm$ 0.65	21.32 $\pm$ 1.79	31.24 $\pm$ 0.3	63.48 $\pm$ 2.25	47.22 $\pm$ 0.84	58.14 $\pm$ 1.6	63.57 $\pm$ 1.28	31.48 $\pm$ 0.98	74.16 $\pm$ 0.56	24.13 $\pm$ 0.85	51.74 $\pm$ 1.09	80.95 $\pm$ 9.52
1.0	22.52 $\pm$ 0.58	21.44 $\pm$ 1.93	31.7 $\pm$ 0.68	63.39 $\pm$ 2.33	47.11 $\pm$ 0.78	57.64 $\pm$ 2.02	61.94 $\pm$ 1.86	30.94 $\pm$ 1.25	73.83 $\pm$ 0.68	23.27 $\pm$ 1.38	51.16 $\pm$ 1.31	0.0 $\pm$ 0.0	

Table 18: Test perplexity on C4, WikiText2 and PTB and zero-shot accuracy of Llama-3.2-3B using *MOON-SHOT*-SparseGPT. Zero-shot **mean performance and win rate** are computed over the 7 classification tasks. Results are averaged over 3 seeds with standard errors.

Sparsity	$\lambda$	C4 $\downarrow$	WikiText2 $\downarrow$	PTB $\downarrow$	BoolQ $\uparrow$	HellaSwag $\uparrow$	WinoGrande $\uparrow$	ARC-e $\uparrow$	ARC-c $\uparrow$	PIQA $\uparrow$	OBQA $\uparrow$	Mean $\uparrow$	Win Rate $\uparrow$
Dense	-	11.34	7.81	13.54	72.72	55.30	69.22	74.37	42.41	76.71	31.20	60.28	-
0.5	0.0	18.12 $\pm$ 0.06	13.0 $\pm$ 0.11	21.83 $\pm$ 0.18	69.19 $\pm$ 1.07	45.96 $\pm$ 0.21	63.67 $\pm$ 0.8	61.52 $\pm$ 1.45	30.69 $\pm$ 0.93	73.38 $\pm$ 0.37	24.27 $\pm$ 0.66	52.67 $\pm$ 0.43	9.52 $\pm$ 4.76
	0.1	16.72 $\pm$ 0.04	11.93 $\pm$ 0.06	19.62 $\pm$ 0.17	70.62 $\pm$ 0.65	47.2 $\pm$ 0.05	64.69 $\pm$ 0.39	64.72 $\pm$ 0.69	32.51 $\pm$ 0.27	<b>73.72</b> $\pm$ 0.25	25.8 $\pm$ 0.31	54.18 $\pm$ 0.18	33.33 $\pm$ 12.6
	0.25	16.63 $\pm$ 0.02	11.86 $\pm$ 0.06	19.49 $\pm$ 0.19	70.59 $\pm$ 0.59	47.36 $\pm$ 0.08	64.09 $\pm$ 0.3	64.86 $\pm$ 0.89	32.51 $\pm$ 0.67	73.16 $\pm$ 0.19	26.0 $\pm$ 0.12	54.08 $\pm$ 0.21	33.33 $\pm$ 4.76
	0.5	16.54 $\pm$ 0.02	11.83 $\pm$ 0.06	19.34 $\pm$ 0.23	71.12 $\pm$ 0.63	47.47 $\pm$ 0.05	64.56 $\pm$ 0.16	65.22 $\pm$ 0.63	32.65 $\pm$ 0.28	73.27 $\pm$ 0.28	25.73 $\pm$ 0.37	54.29 $\pm$ 0.17	33.33 $\pm$ 4.76
	0.75	16.49 $\pm$ 0.04	11.78 $\pm$ 0.06	19.18 $\pm$ 0.19	71.0 $\pm$ 0.74	47.64 $\pm$ 0.05	64.33 $\pm$ 0.41	<b>65.4</b> $\pm$ 0.69	32.94 $\pm$ 0.39	73.3 $\pm$ 0.36	26.13 $\pm$ 0.24	54.39 $\pm$ 0.28	<b>47.62</b> $\pm$ 4.76
	1.0	16.45 $\pm$ 0.02	11.76 $\pm$ 0.05	19.05 $\pm$ 0.11	71.15 $\pm$ 0.76	<b>47.65</b> $\pm$ 0.08	65.09 $\pm$ 0.23	65.01 $\pm$ 0.88	33.3 $\pm$ 0.57	73.36 $\pm$ 0.18	26.27 $\pm$ 0.37	54.55 $\pm$ 0.28	<b>47.62</b> $\pm$ 4.76
		17.61 $\pm$ 0.08	12.58 $\pm$ 0.06	20.34 $\pm$ 0.13	<b>72.41</b> $\pm$ 0.77	47.32 $\pm$ 0.1	<b>65.93</b> $\pm$ 0.25	64.27 $\pm$ 0.53	<b>33.36</b> $\pm$ 0.37	72.98 $\pm$ 0.31	<b>26.8</b> $\pm$ 0.31	<b>54.72</b> $\pm$ 0.24	-
0.5 (Alpha-Pruning)	0.0	18.22 $\pm$ 0.02	13.07 $\pm$ 0.1	21.67 $\pm$ 0.21	68.6 $\pm$ 0.75	46.08 $\pm$ 0.12	63.98 $\pm$ 0.82	63.45 $\pm$ 0.22	30.94 $\pm$ 0.57	72.96 $\pm$ 0.23	23.87 $\pm$ 1.05	52.84 $\pm$ 0.28	4.76 $\pm$ 4.76
	0.1	16.8 $\pm$ 0.03	11.99 $\pm$ 0.08	19.59 $\pm$ 0.26	71.05 $\pm$ 0.34	47.37 $\pm$ 0.05	65.69 $\pm$ 0.15	65.7 $\pm$ 0.18	32.82 $\pm$ 0.1	73.16 $\pm$ 0.16	25.47 $\pm$ 0.64	54.47 $\pm$ 0.14	42.86 $\pm$ 8.25
	0.25	16.72 $\pm$ 0.02	11.91 $\pm$ 0.07	19.47 $\pm$ 0.26	71.35 $\pm$ 0.56	47.51 $\pm$ 0.05	65.25 $\pm$ 0.37	65.21 $\pm$ 0.22	32.99 $\pm$ 0.28	73.2 $\pm$ 0.42	25.27 $\pm$ 0.18	54.39 $\pm$ 0.11	42.86 $\pm$ 0.0
	0.5	16.66 $\pm$ 0.03	11.86 $\pm$ 0.06	19.36 $\pm$ 0.24	71.85 $\pm$ 0.4	47.69 $\pm$ 0.03	65.17 $\pm$ 0.27	65.47 $\pm$ 0.33	32.59 $\pm$ 0.81	73.49 $\pm$ 0.41	25.13 $\pm$ 0.29	54.48 $\pm$ 0.14	42.86 $\pm$ 0.0
	0.75	16.6 $\pm$ 0.03	11.83 $\pm$ 0.06	19.26 $\pm$ 0.24	71.64 $\pm$ 0.43	47.73 $\pm$ 0.06	65.22 $\pm$ 0.64	65.77 $\pm$ 0.16	33.02 $\pm$ 0.58	<b>73.7</b> $\pm$ 0.49	26.07 $\pm$ 0.24	54.74 $\pm$ 0.18	42.86 $\pm$ 0.0
	1.0	<b>16.56</b> $\pm$ 0.03	<b>11.79</b> $\pm$ 0.05	<b>19.01</b> $\pm$ 0.24	71.96 $\pm$ 0.28	<b>47.85</b> $\pm$ 0.03	65.19 $\pm$ 0.66	<b>66.12</b> $\pm$ 0.18	33.42 $\pm$ 0.86	73.45 $\pm$ 0.33	26.27 $\pm$ 0.07	54.89 $\pm$ 0.18	<b>52.38</b> $\pm$ 4.76
		17.65 $\pm$ 0.06	12.59 $\pm$ 0.05	20.31 $\pm$ 0.17	<b>72.93</b> $\pm$ 0.22	47.31 $\pm$ 0.04	<b>66.09</b> $\pm$ 0.65	63.63 $\pm$ 0.28	<b>34.24</b> $\pm$ 0.55	72.98 $\pm$ 0.24	<b>26.8</b> $\pm$ 0.31	<b>55.0</b> $\pm$ 0.18	-
0.5 (OWL)	0.0	17.62 $\pm$ 0.02	12.92 $\pm$ 0.09	21.0 $\pm$ 0.15	69.13 $\pm$ 0.85	46.64 $\pm$ 0.12	64.63 $\pm$ 0.83	62.63 $\pm$ 0.49	31.11 $\pm$ 0.41	72.89 $\pm$ 0.33	24.0 $\pm$ 0.2	53.04 $\pm$ 0.1	9.52 $\pm$ 4.76
	0.1	16.54 $\pm$ 0.02	11.99 $\pm$ 0.06	19.46 $\pm$ 0.24	70.65 $\pm$ 0.61	47.82 $\pm$ 0.1	65.93 $\pm$ 0.42	66.13 $\pm$ 0.74	33.16 $\pm$ 0.37	72.98 $\pm$ 0.13	24.8 $\pm$ 0.31	54.5 $\pm$ 0.08	38.1 $\pm$ 12.6
	0.25	16.45 $\pm$ 0.01	11.95 $\pm$ 0.07	19.08 $\pm$ 0.21	70.83 $\pm$ 0.07	48.01 $\pm$ 0.04	66.04 $\pm$ 0.37	<b>66.61</b> $\pm$ 0.32	33.42 $\pm$ 0.28	73.07 $\pm$ 0.41	25.27 $\pm$ 0.24	54.75 $\pm$ 0.08	42.86 $\pm$ 8.25
	0.5	16.4 $\pm$ 0.0	11.89 $\pm$ 0.06	19.06 $\pm$ 0.3	71.04 $\pm$ 0.1	48.11 $\pm$ 0.05	65.77 $\pm$ 0.66	66.32 $\pm$ 0.46	32.94 $\pm$ 0.3	73.25 $\pm$ 0.2	24.93 $\pm$ 0.37	54.62 $\pm$ 0.15	42.86 $\pm$ 8.25
	0.75	16.36 $\pm$ 0.02	11.80 $\pm$ 0.06	18.96 $\pm$ 0.26	71.3 $\pm$ 0.18	48.15 $\pm$ 0.08	66.14 $\pm$ 0.36	66.41 $\pm$ 0.35	33.02 $\pm$ 0.6	73.25 $\pm$ 0.18	25.13 $\pm$ 0.18	54.77 $\pm$ 0.12	<b>52.38</b> $\pm$ 4.76
	1.0	<b>16.33</b> $\pm$ 0.02	<b>11.83</b> $\pm$ 0.07	<b>18.85</b> $\pm$ 0.29	71.49 $\pm$ 0.2	<b>48.21</b> $\pm$ 0.06	66.32 $\pm$ 0.07	66.59 $\pm$ 0.56	33.59 $\pm$ 0.51	<b>73.39</b> $\pm$ 0.19	25.27 $\pm$ 0.18	54.9 $\pm$ 0.11	42.86 $\pm$ 8.25
		17.14 $\pm$ 0.06	12.35 $\pm$ 0.03	19.81 $\pm$ 0.25	<b>72.73</b> $\pm$ 0.17	47.86 $\pm$ 0.06	<b>66.4</b> $\pm$ 0.17	65.25 $\pm$ 0.24	<b>34.36</b> $\pm$ 0.77	72.69 $\pm$ 0.09	<b>26.2</b> $\pm$ 0.31	<b>55.07</b> $\pm$ 0.16	-
0.6	0.0	40.78 $\pm$ 0.89	33.01 $\pm$ 1.0	58.58 $\pm$ 2.31	62.36 $\pm$ 1.25	35.6 $\pm$ 0.06	56.54 $\pm$ 0.47	49.93 $\pm$ 0.51	23.78 $\pm$ 0.58	66.87 $\pm$ 0.17	16.33 $\pm$ 0.66	44.49 $\pm$ 0.25	0.0 $\pm$ 0.0
	0.1	30.2 $\pm$ 0.09	23.74 $\pm$ 0.25	39.7 $\pm$ 1.07	66.47 $\pm$ 0.92	38.47 $\pm$ 0.06	60.43 $\pm$ 0.54	56.5 $\pm$ 0.39	26.96 $\pm$ 0.13	68.68 $\pm$ 0.13	19.93 $\pm$ 0.71	48.21 $\pm$ 0.13	80.95 $\pm$ 4.76
	0.25	29.21 $\pm$ 0.22	22.9 $\pm$ 0.23	37.44 $\pm$ 0.78	67.42 $\pm$ 0.06	38.88 $\pm$ 0.0	59.72 $\pm$ 0.4	56.66 $\pm$ 0.58	26.71 $\pm$ 0.21	69.24 $\pm$ 0.08	19.93 $\pm$ 0.77	48.37 $\pm$ 0.03	76.19 $\pm$ 4.76
	0.5	28.76 $\pm$ 0.13	22.81 $\pm$ 0.17	36.8 $\pm$ 0.62	67.67 $\pm$ 0.3	38.82 $\pm$ 0.1	60.38 $\pm$ 0.57	57.55 $\pm$ 0.86	27.1 $\pm$ 0.44	68.99 $\pm$ 0.25	19.47 $\pm$ 0.85	48.57 $\pm$ 0.1	85.71 $\pm$ 0.0
	0.75	28.57 $\pm$ 0.18	22.65 $\pm$ 0.21	36.17 $\pm$ 0.5	67.6 $\pm$ 0.13	38.96 $\pm$ 0.14	<b>61.27</b> $\pm$ 0.19	57.37 $\pm$ 0.83	27.39 $\pm$ 0.31	69.04 $\pm$ 0.57	19.73 $\pm$ 0.82	48.77 $\pm$ 0.12	90.48 $\pm$ 4.76
	1.0	<b>28.23</b> $\pm$ 0.11	<b>22.46</b> $\pm$ 0.17	<b>35.63</b> $\pm$ 0.68	<b>67.76</b> $\pm$ 0.35	<b>39.13</b> $\pm$ 0.07	61.01 $\pm$ 0.23	<b>57.59</b> $\pm$ 0.92	<b>27.79</b> $\pm$ 0.71	<b>69.44</b> $\pm$ 0.13	<b>20.0</b> $\pm$ 0.53	<b>48.96</b> $\pm$ 0.12	<b>95.24</b> $\pm$ 4.76
		33.63 $\pm$ 0.14	26.12 $\pm$ 0.23	42.69 $\pm$ 0.73	66.82 $\pm$ 0.6	38.14 $\pm$ 0.14	60.91 $\pm$ 0.65	53.89 $\pm$ 0.11	26.28 $\pm$ 0.18	67.75 $\pm$ 0.34	18.47 $\pm$ 0.44	47.47 $\pm$ 0.08	-
0.6 (Alpha-Pruning)	0.0	39.98 $\pm$ 0.8	32.14 $\pm$ 1.25	57.63 $\pm$ 0.55	64.43 $\pm$ 0.69	36.11 $\pm$ 0.05	57.27 $\pm$ 0.35	49.96 $\pm$ 0.18	23.83 $\pm$ 0.9	66.7 $\pm$ 0.39	17.47 $\pm$ 0.35	45.11 $\pm$ 0.14	0.0 $\pm$ 0.0
	0.1	30.39 $\pm$ 0.11	23.23 $\pm$ 0.16	37.83 $\pm$ 1.2	67.9 $\pm$ 0.49	38.53 $\pm$ 0.03	61.51 $\pm$ 0.27	55.19 $\pm$ 0.05	26.42 $\pm$ 0.29	68.39 $\pm$ 0.51	19.33 $\pm$ 0.18	48.18 $\pm$ 0.08	71.43 $\pm$ 8.25
	0.25	29.68 $\pm$ 0.05	23.01 $\pm$ 0.3	36.99 $\pm$ 0.82	68.92 $\pm$ 0.39	38.89 $\pm$ 0.09	61.72 $\pm$ 0.39	55.27 $\pm$ 0.27	26.02 $\pm$ 0.23	68.52 $\pm$ 0.35	19.07 $\pm$ 0.44	48.34 $\pm$ 0.08	80.95 $\pm$ 12.6
	0.5	29.14 $\pm$ 0.1	22.56 $\pm$ 0.07	36.56 $\pm$ 1.25	68.5 $\pm$ 0.05	39.19 $\pm$ 0.08	60.51 $\pm$ 0.37	56.41 $\pm$ 0.37	<b>26.99</b> $\pm$ 0.46	68.92 $\pm$ 0.34	19.87 $\pm$ 0.24	48.63 $\pm$ 0.07	71.43 $\pm$ 14.29
	0.75	28.96 $\pm$ 0.15	22.44 $\pm$ 0.04	35.89 $\pm$ 1.52	<b>68.93</b> $\pm$ 0.35	39.15 $\pm$ 0.14	61.3 $\pm$ 0.27	56.45 $\pm$ 0.46	26.82 $\pm$ 0.54	68.59 $\pm$ 0.14	20.0 $\pm$ 0.83	48.75 $\pm$ 0.05	80.95 $\pm$ 4.76
	1.0	<b>28.64</b> $\pm$ 0.25	<b>22.17</b> $\pm$ 0.19	<b>35.48</b> $\pm$ 1.52	68.73 $\pm$ 0.11	<b>39.34</b> $\pm$ 0.18	<b>62.27</b> $\pm$ 0.52	<b>56.52</b> $\pm$ 0.25	26.93 $\pm$ 0.37	<b>69.13</b> $\pm$ 0.32	<b>20.33</b> $\pm$ 0.24	<b>49.04</b> $\pm$ 0.06	<b>85.71</b> $\pm$ 8.25
		34.19 $\pm$ 0.4	26.06 $\pm$ 0.64	42.82 $\pm$ 0.13	68.17 $\pm$ 0.43	38.62 $\pm$ 0.16	61.98 $\pm$ 0.81	53.37 $\pm$ 0.7	26.0 $\pm$ 0.23	68.06 $\pm$ 0.38	19.8 $\pm$ 0.38	48.0 $\pm$ 0.38	-
0.6 (OWL)	0.0	32.72 $\pm$ 0.56	27.69 $\pm$ 0.6	45.87 $\pm$ 0.76	64.46 $\pm$ 0.51	38.09 $\pm$ 0.1	58.01 $\pm$ 0.25	53.38 $\pm$ 0.31	25.14 $\pm$ 0.33	67.94 $\pm$ 0.1	18.0 $\pm$ 0.4	46.43 $\pm$ 0.14	0.0 $\pm$ 0.0
	0.1	26.66 $\pm$ 0.16	21.77 $\pm$ 0.33	33.89 $\pm$ 0.66	66.67 $\pm$ 0.14	40.01 $\pm$ 0.11	61.93 $\pm$ 0.66	56.8 $\pm$ 0.29	27.79 $\pm$ 0.54	69.17 $\pm$ 0.4	20.67 $\pm$ 0.18	49.01 $\pm$ 0.08	61.9 $\pm$ 4.76
	0.25	26.17 $\pm$ 0.12	21.56 $\pm$ 0.35	33.74 $\pm$ 0.66	66.99 $\pm$ 0.37	40.38 $\pm$ 0.07	61.14 $\pm$ 0.49	56.93 $\pm$ 0.78	27.9 $\pm$ 0.69	69.48 $\pm$ 0.45	<b>21.4</b> $\pm$ 0.35	49.17 $\pm$ 0.16	<b>71.43</b> $\pm$ 14.29
	0.5	25.77 $\pm$ 0.16	21.16 $\pm$ 0.23	32.36 $\pm$ 0.55	<b>67.41</b> $\pm$ 0.49	40.51 $\pm$ 0.12	61.48 $\pm$ 0.2	57.15 $\pm$ 0.86	27.87 $\pm$ 0.72	69.7 $\pm$ 0.44	21.07 $\pm$ 0.07	49.31 $\pm$ 0.2	66.67 $\pm$ 12.6
	0.75	25.55 $\pm$ 0.12	21.0 $\pm$ 0.21	31.7 $\pm$ 0.41	67.35 $\pm$ 0.33	40.61 $\pm$ 0.13	61.48 $\pm$ 0.84	57.06 $\pm$ 1.27	27.7 $\pm$ 0.46	<b>69.99</b> $\pm$ 0.61			

Table 19: Test perplexity on C4, WikiText2 and PTB and zero-shot accuracy of Llama-3.2-3B using *MOON-SHOT*-Wanda. Zero-shot **mean performance and win rate** are computed over the 7 classification tasks. Results are averaged over 3 seeds with standard errors.

Sparsity	$\lambda$	C4 ↓	WikiText2 ↓	PTB ↓	BoolQ ↑	HellaSwag ↑	WinoGrande ↑	ARC-e ↑	ARC-c ↑	PIQA ↑	OBQA ↑	Mean ↑	Win Rate ↑
Dense	-	11.34	7.81	13.54	72.72	55.30	69.22	74.37	42.41	76.71	31.20	60.28	-
0.5	0.0	18.2 ± 0.04	12.52 ± 0.0	21.25 ± 0.03	64.89 ± 0.39	45.57 ± 0.04	63.04 ± 0.41	65.8 ± 0.33	32.71 ± 0.57	73.25 ± 0.15	25.53 ± 0.52	52.97 ± 0.19	42.86 ± 14.29
	0.1	18.16 ± 0.03	12.48 ± 0.0	21.17 ± 0.05	65.24 ± 0.16	45.6 ± 0.07	62.9 ± 0.16	65.98 ± 0.27	32.57 ± 0.5	73.47 ± 0.1	25.53 ± 0.41	53.04 ± 0.21	52.38 ± 12.6
	0.25	18.08 ± 0.03	12.42 ± 0.01	21.04 ± 0.07	64.22 ± 0.47	45.61 ± 0.07	62.75 ± 0.61	65.77 ± 0.35	32.48 ± 0.44	73.36 ± 0.35	25.6 ± 0.35	52.83 ± 0.09	52.38 ± 12.6
	0.5	18.0 ± 0.04	12.36 ± 0.01	20.87 ± 0.04	64.98 ± 0.31	45.73 ± 0.03	<b>63.11</b> ± 0.3	65.71 ± 0.59	32.11 ± 0.16	73.36 ± 0.15	25.67 ± 0.07	52.95 ± 0.1	52.38 ± 12.6
	0.75	<b>17.99</b> ± 0.04	<b>12.35</b> ± 0.02	20.76 ± 0.05	64.59 ± 0.63	<b>45.84</b> ± 0.04	62.56 ± 0.42	66.16 ± 0.35	32.65 ± 0.15	<b>73.61</b> ± 0.21	25.67 ± 0.66	53.01 ± 0.15	61.9 ± 12.6
	0.9	18.02 ± 0.02	12.37 ± 0.02	<b>20.74</b> ± 0.05	63.91 ± 0.83	45.81 ± 0.08	62.64 ± 0.41	<b>66.36</b> ± 0.34	<b>32.99</b> ± 0.34	73.32 ± 0.07	25.6 ± 0.42	52.95 ± 0.04	<b>66.67</b> ± 4.76
1.0	18.88 ± 0.03	13.0 ± 0.01	21.73 ± 0.07	<b>66.41</b> ± 0.39	45.64 ± 0.08	62.96 ± 0.14	65.84 ± 0.17	32.68 ± 0.27	72.67 ± 0.17	<b>25.8</b> ± 0.61	<b>53.14</b> ± 0.08	-	
0.5 (Alpha-Pruning)	0.0	18.19 ± 0.03	12.48 ± 0.02	21.31 ± 0.03	66.27 ± 0.68	45.69 ± 0.08	63.67 ± 0.25	<b>66.89</b> ± 0.3	33.42 ± 0.41	73.45 ± 0.2	25.33 ± 0.77	53.53 ± 0.2	52.38 ± 17.17
	0.1	18.14 ± 0.02	12.42 ± 0.02	21.17 ± 0.05	65.96 ± 0.73	45.7 ± 0.12	63.46 ± 0.33	66.72 ± 0.38	33.3 ± 0.1	73.12 ± 0.24	25.47 ± 0.68	53.39 ± 0.12	42.86 ± 16.5
	0.25	18.05 ± 0.04	12.37 ± 0.02	21.08 ± 0.08	66.55 ± 0.42	45.84 ± 0.03	63.64 ± 0.38	66.37 ± 0.27	32.79 ± 0.38	73.3 ± 0.17	25.27 ± 0.47	53.4 ± 0.06	42.86 ± 16.5
	0.5	17.97 ± 0.03	12.29 ± 0.01	20.87 ± 0.01	66.68 ± 0.58	45.9 ± 0.08	<b>63.83</b> ± 0.3	66.71 ± 0.14	33.13 ± 0.44	73.38 ± 0.08	25.53 ± 0.27	53.59 ± 0.12	52.38 ± 9.52
	0.75	<b>17.94</b> ± 0.03	<b>12.28</b> ± 0.02	20.75 ± 0.03	66.18 ± 0.8	46.05 ± 0.05	63.3 ± 0.19	66.62 ± 0.18	33.53 ± 0.18	<b>73.54</b> ± 0.16	<b>25.67</b> ± 0.44	53.57 ± 0.1	<b>57.14</b> ± 8.25
	0.9	17.97 ± 0.02	12.3 ± 0.04	<b>20.69</b> ± 0.01	66.25 ± 0.73	<b>46.08</b> ± 0.01	62.83 ± 0.41	<b>66.89</b> ± 0.11	33.5 ± 0.16	73.3 ± 0.24	25.6 ± 0.35	53.49 ± 0.1	52.38 ± 4.76
1.0	18.78 ± 0.01	12.83 ± 0.03	21.42 ± 0.04	<b>67.72</b> ± 0.78	45.9 ± 0.04	63.38 ± 0.4	66.69 ± 0.1	<b>34.04</b> ± 0.17	72.65 ± 0.13	25.07 ± 0.29	<b>53.64</b> ± 0.19	-	
0.5 (OWL)	0.0	18.24 ± 0.05	12.58 ± 0.02	20.7 ± 0.04	67.37 ± 0.34	46.05 ± 0.06	64.11 ± 0.23	65.88 ± 0.23	32.76 ± 0.32	73.09 ± 0.17	24.87 ± 0.35	53.45 ± 0.16	42.86 ± 8.25
	0.1	18.2 ± 0.04	12.55 ± 0.02	20.69 ± 0.06	67.12 ± 0.3	46.09 ± 0.01	63.48 ± 0.53	66.13 ± 0.23	32.88 ± 0.24	73.16 ± 0.13	24.67 ± 0.27	53.36 ± 0.17	33.33 ± 4.76
	0.25	18.16 ± 0.04	12.51 ± 0.01	20.6 ± 0.04	66.95 ± 0.22	46.13 ± 0.06	63.77 ± 0.24	66.25 ± 0.08	32.82 ± 0.23	73.25 ± 0.2	<b>25.2</b> ± 0.4	53.48 ± 0.1	42.86 ± 8.25
	0.5	18.13 ± 0.04	12.49 ± 0.01	20.5 ± 0.02	66.37 ± 0.44	46.26 ± 0.02	64.14 ± 0.23	65.91 ± 0.39	32.99 ± 0.27	<b>73.39</b> ± 0.32	25.0 ± 0.35	53.44 ± 0.13	47.62 ± 9.52
	0.75	<b>18.11</b> ± 0.03	<b>12.47</b> ± 0.01	20.35 ± 0.02	66.52 ± 0.63	46.25 ± 0.05	64.06 ± 0.07	66.37 ± 0.21	33.36 ± 0.27	73.14 ± 0.14	24.8 ± 0.31	53.5 ± 0.07	<b>52.38</b> ± 12.6
	0.9	<b>18.11</b> ± 0.03	<b>12.47</b> ± 0.03	<b>20.31</b> ± 0.02	66.3 ± 0.92	<b>46.39</b> ± 0.04	63.98 ± 0.35	<b>66.4</b> ± 0.22	33.59 ± 0.24	73.14 ± 0.1	24.93 ± 0.35	53.53 ± 0.17	<b>52.38</b> ± 4.76
1.0	18.82 ± 0.03	12.98 ± 0.01	21.06 ± 0.04	<b>68.98</b> ± 0.49	46.24 ± 0.05	<b>64.46</b> ± 0.27	66.05 ± 0.32	<b>34.22</b> ± 0.05	72.45 ± 0.04	24.73 ± 0.37	<b>53.88</b> ± 0.05	-	
0.6	0.0	40.28 ± 0.67	30.39 ± 0.33	52.83 ± 1.47	60.46 ± 0.73	35.1 ± 0.12	55.12 ± 0.29	51.57 ± 0.31	23.81 ± 0.44	66.78 ± 0.13	16.87 ± 0.35	44.24 ± 0.27	57.14 ± 8.25
	0.1	39.37 ± 0.65	29.45 ± 0.33	50.79 ± 1.44	60.4 ± 0.72	35.14 ± 0.13	54.96 ± 0.16	51.77 ± 0.23	23.95 ± 0.27	67.01 ± 0.09	17.0 ± 0.42	44.32 ± 0.21	57.14 ± 8.25
	0.25	38.87 ± 0.45	28.74 ± 0.3	49.29 ± 1.24	60.68 ± 0.55	35.25 ± 0.12	55.09 ± 0.66	52.3 ± 0.29	24.06 ± 0.15	67.14 ± 0.11	16.27 ± 0.29	44.4 ± 0.11	57.14 ± 0.0
	0.5	38.15 ± 0.24	28.04 ± 0.14	47.61 ± 0.7	60.45 ± 1.05	35.53 ± 0.07	55.17 ± 0.09	52.31 ± 0.2	23.98 ± 0.34	<b>67.23</b> ± 0.05	<b>17.07</b> ± 0.24	44.53 ± 0.22	<b>61.9</b> ± 4.76
	0.75	37.81 ± 0.05	27.75 ± 0.16	46.75 ± 0.16	61.86 ± 1.24	<b>35.54</b> ± 0.02	55.56 ± 0.6	<b>56.89</b> ± 0.22	24.49 ± 0.49	66.49 ± 0.14	16.53 ± 0.18	44.74 ± 0.26	<b>61.9</b> ± 4.76
	0.9	<b>37.73</b> ± 0.19	<b>27.71</b> ± 0.26	<b>46.47</b> ± 0.08	61.33 ± 0.91	35.53 ± 0.08	54.83 ± 0.14	52.53 ± 0.29	<b>24.69</b> ± 0.21	66.81 ± 0.03	16.6 ± 0.23	44.62 ± 0.12	<b>61.9</b> ± 4.76
1.0	41.98 ± 0.4	30.56 ± 0.32	51.0 ± 0.45	<b>64.82</b> ± 0.35	35.12 ± 0.07	<b>56.56</b> ± 0.46	50.58 ± 0.41	23.83 ± 0.12	65.58 ± 0.19	16.93 ± 0.03	<b>44.77</b> ± 0.1	-	
0.6 (Alpha-Pruning)	0.0	39.14 ± 0.81	29.06 ± 0.27	52.35 ± 1.36	60.69 ± 2.09	35.6 ± 0.16	56.49 ± 0.52	53.45 ± 0.25	24.43 ± 0.08	66.32 ± 0.22	17.0 ± 0.5	44.86 ± 0.4	47.62 ± 4.76
	0.1	38.33 ± 0.87	28.37 ± 0.37	50.78 ± 1.45	61.15 ± 1.57	35.66 ± 0.11	57.04 ± 0.34	53.48 ± 0.14	24.66 ± 0.3	66.59 ± 0.46	<b>17.2</b> ± 0.4	45.11 ± 0.26	61.9 ± 4.76
	0.25	37.79 ± 0.85	27.72 ± 0.35	49.1 ± 1.42	61.74 ± 2.05	35.76 ± 0.08	57.22 ± 0.39	53.82 ± 0.41	24.8 ± 0.1	66.67 ± 0.35	17.13 ± 0.57	45.31 ± 0.35	57.14 ± 0.0
	0.5	37.56 ± 0.43	27.26 ± 0.13	47.87 ± 0.88	63.15 ± 1.22	35.95 ± 0.08	56.7 ± 0.19	53.94 ± 0.33	25.11 ± 0.37	66.76 ± 0.3	16.13 ± 0.29	45.39 ± 0.17	47.62 ± 4.76
	0.75	<b>37.37</b> ± 0.21	<b>27.08</b> ± 0.16	<b>47.01</b> ± 0.29	63.38 ± 0.65	36.05 ± 0.1	<b>57.51</b> ± 0.56	53.15 ± 0.27	<b>25.43</b> ± 0.18	<b>66.96</b> ± 0.36	16.27 ± 0.18	<b>45.68</b> ± 0.1	61.9 ± 4.76
	0.9	37.57 ± 0.18	27.31 ± 0.1	47.45 ± 0.29	63.07 ± 0.61	<b>36.18</b> ± 0.11	57.2 ± 0.17	53.86 ± 0.3	25.23 ± 0.19	66.94 ± 0.31	16.73 ± 0.35	45.6 ± 0.07	<b>66.67</b> ± 4.76
1.0	40.03 ± 0.04	29.19 ± 0.2	50.24 ± 0.27	<b>65.93</b> ± 0.24	35.95 ± 0.01	<b>57.51</b> ± 0.09	51.09 ± 0.16	24.32 ± 0.36	66.03 ± 0.15	16.87 ± 0.24	45.39 ± 0.03	-	
0.6 (OWL)	0.0	35.96 ± 0.47	27.31 ± 0.4	45.24 ± 0.7	61.81 ± 0.96	37.06 ± 0.16	57.51 ± 0.39	53.72 ± 0.64	26.11 ± 0.13	67.21 ± 0.13	<b>17.6</b> ± 0.5	45.86 ± 0.14	57.14 ± 0.0
	0.1	35.4 ± 0.33	26.77 ± 0.33	43.81 ± 0.77	61.4 ± 1.55	37.16 ± 0.2	57.35 ± 0.34	54.0 ± 0.67	25.8 ± 0.22	67.27 ± 0.15	17.27 ± 0.41	45.75 ± 0.25	47.62 ± 4.76
	0.25	34.99 ± 0.17	26.44 ± 0.23	42.75 ± 0.89	61.95 ± 1.53	37.22 ± 0.12	57.56 ± 0.16	54.22 ± 0.44	26.08 ± 0.27	<b>67.3</b> ± 0.23	17.13 ± 0.18	45.92 ± 0.24	61.9 ± 4.76
	0.5	34.71 ± 0.22	25.94 ± 0.15	40.98 ± 0.64	63.11 ± 1.86	37.28 ± 0.14	57.91 ± 0.43	53.82 ± 0.44	<b>26.31</b> ± 0.03	<b>67.3</b> ± 0.06	17.2 ± 0.23	46.2 ± 0.29	<b>66.67</b> ± 4.76
	0.75	<b>34.5</b> ± 0.08	25.64 ± 0.08	<b>40.35</b> ± 0.63	64.06 ± 1.0	<b>37.41</b> ± 0.12	58.33 ± 0.25	<b>54.41</b> ± 0.39	25.85 ± 0.18	67.03 ± 0.17	17.4 ± 0.12	46.35 ± 0.17	61.9 ± 4.76
	0.9	34.56 ± 0.11	<b>25.59</b> ± 0.06	40.41 ± 0.78	63.73 ± 0.92	37.4 ± 0.12	58.93 ± 0.35	54.31 ± 0.07	25.68 ± 0.13	67.05 ± 0.25	17.0 ± 0.23	46.3 ± 0.13	<b>66.67</b> ± 9.52
1.0	37.35 ± 0.14	27.93 ± 0.23	44.25 ± 0.74	<b>67.26</b> ± 0.07	37.18 ± 0.13	<b>59.06</b> ± 0.73	51.89 ± 0.34	25.54 ± 0.23	66.47 ± 0.1	17.2 ± 0.12	<b>46.37</b> ± 0.18	-	
0.7	0.0	246.47 ± 5.88	232.57 ± 8.94	323.16 ± 2.66	40.2 ± 1.09	27.32 ± 0.04	<b>49.33</b> ± 0.62	33.52 ± 0.2	17.83 ± 0.17	56.67 ± 0.14	12.0 ± 0.42	33.84 ± 0.28	<b>52.38</b> ± 4.76
	0.1	240.18 ± 3.7	224.64 ± 6.04	318.96 ± 7.59	40.63 ± 0.96	27.39 ± 0.05	48.38 ± 0.16	33.33 ± 0.21	17.26 ± 0.15	57.02 ± 0.19	12.0 ± 0.23	33.72 ± 0.21	42.86 ± 0.0
	0.25	240.29 ± 4.53	225.62 ± 7.31	307.06 ± 10.91	40.39 ± 1.24	27.43 ± 0.03	48.17 ± 0.7	33.49 ± 0.16	17.12 ± 0.17	57.05 ± 0.25	11.8 ± 0.0	33.64 ± 0.27	47.62 ± 4.76
	0.5	247.47 ± 3.63	222.78 ± 6.1	310.9 ± 14.25	40.31 ± 1.26	27.55 ± 0.02	48.51 ± 0.88	33.54 ± 0.19	17.12 ± 0.12	57.16 ± 0.17	<b>12.13</b> ± 0.13	33.76 ± 0.34	47.62 ± 4.76
	0.75	239.91 ± 9.44	206.01 ± 6.04	298.6 ± 4.91	38.66 ± 0.18	27.67 ± 0.04	48.46 ± 0.6	<b>33.87</b> ± 0.47	17.21 ± 0.42	57.25 ± 0.28	11.87 ± 0.07	33.57 ± 0.15	47.62 ± 4.76
	0.9	<b>235.73</b> ± 6.87	<b>195.34</b> ± 4.93	<b>297.54</b> ± 7.2	38.85 ± 0.35	27.73 ± 0.02	47.91 ± 0.32	33.54 ± 0.25	17.52 ± 0.32	<b>57.56</b> ± 0.08	12.07 ± 0.41	33.6 ± 0.16	47.62 ± 4.76
1.0	250.62 ± 7.51	230.39 ± 4.31	332.57 ± 5.15	<b>50.06</b> ± 0.28	<b>27.79</b> ± 0.06	48.57 ± 0.18	32.74 ± 0.28	<b>17.89</b> ± 0.2	56.6 ± 0.24	10.73 ± 0.24	<b>34.91</b> ± 0.13	-	
0.7 (Alpha-Pruning)	0.0	229.42 ± 11.77	202.94 ± 9.43	282.4 ± 14.97	49.63 ± 2.65	27.27 ± 0.03	48.46 ± 0.48	32.49 ± 0.47	17.72 ± 0.3	57.04 ± 0.02	12.27 ± 0.13	34.98 ± 0.42	38.1 ± 12.6
	0.1	222.38 ± 9.21	190.27 ± 6.69	261.37 ± 14.19	50.94 ± 2.29	27.29 ± 0.04	<b>49.57</b> ± 0.55	32.72 ± 0.42	<b>17.78</b> ± 0.23	57.29 ± 0.34	12.27 ± 0.18	35.41 ± 0.34	42.86 ± 14.29
	0.25	216.49 ± 6.86	183.02 ± 8.68	250.18 ± 20.32	50.7 ± 1.28	27.31 ± 0.1	48.43 ± 0.25	<b>33.07</b> ± 0.36	17.61 ± 0.06	57.09 ± 0.29	11.87 ± 0.33	35.15 ± 0.13	38.1 ± 4.76
	0.5	215.37 ± 6.76	1										