# BAME: BLOCK-AWARE MASK EVOLUTION FOR EFFI CIENT N:M SPARSE TRAINING

Anonymous authors

Paper under double-blind review

## ABSTRACT

N:M sparsity stands as a progressively important tool for DNN compression, achieving practical speedups by stipulating at most N non-zero components within M sequential weights. Unfortunately, most existing works identify the N:M sparse mask through dense backward propagation to update all weights, which incurs exorbitant training costs. In this paper, we introduce BAME, a method that maintains consistent sparsity throughout the N:M sparse training process. BAME perpetually keeps both sparse forward and backward propagation, while iteratively performing weight pruning-and-regrowing within designated weight blocks to tailor the N:M mask. These blocks are selected through a joint assessment based on accumulated mask oscillation frequency and expected loss reduction of mask adaptation, thereby ensuring stable and efficient identification of the optimal N:M mask. Our empirical results substantiate the effectiveness of BAME, illustrating it performs comparably to or better than previous works that fully maintaining dense backward propagation during training. For instance, BAME attains a 72.0% top-1 accuracy while training a 1:16 sparse ResNet-50 on ImageNet, eclipsing SR-STE by 0.5%, despite achieving  $2.37 \times$  training FLOPs reduction. Code will be released.

026 027 028

025

004

010 011

012

013

014

015

016

017

018

019

021

### 1 INTRODUCTION

029 030

031 In recent years, the vision community has precipitously bolstered the performance of Deep Neural 032 Networks (DNNs) across various tasks, including image classification (He et al., 2016), object 033 detection (He et al., 2017a), and semantic segmentation (Girshick et al., 2014), etc. These progressions 034 are chiefly driven by an augmented parameter burden and an increasingly onerous computational cost. Regrettably, this tendency presents significant impediments for the deployment of DNNs on resource-constrained edge devices like smartphones and various Internet of Things (IoT) apparatuses. Consequently, there has been a proliferation of interest in model compression research (Hubara et al., 037 2016; Howard et al., 2017; Lin et al., 2020), with the explicit objective of reducing the model's computation and parameter complexity whilst preserving comparable performance to the original model, thereby alleviating the deployment tribulations experienced with DNNs. 040

Among these techniques, network sparsity has proven many successes (Han et al., 2015; LeCun 041 et al., 1989; Luo et al., 2017) by zeroizing weights to yield lightweight, sparse networks at different 042 granularity levels, from fine to coarse. Fine-grained sparsity (unstructured sparsity) (LeCun et al., 043 1989; Ding et al., 2019) removes individual weights and is demonstrated to well retain performance 044 even at high sparsity rates. Regrettably, the deployment of such fine-grained sparse networks onto mainstream hardware systems becomes exceptionally challenging, given the irregular matrix patterns 046 created by sparse weights. In contrast, coarse-grained sparsity, otherwise known as structured sparsity, 047 (He et al., 2017b; Lin et al., 2020) procures substantial acceleration, purging whole convolution filters 048 in the process (Liu et al., 2019; Lin et al., 2020). Nevertheless, structured sparsity can experience severe performance degradation, especially under high sparsity conditions. Recent developments indicate N:M sparsity as an auspicious avenue towards effectively balancing the dual requirements 051 of acceleration and performance retention (Zhou et al., 2021; Pool & Yu, 2021). By imposing a restriction of, at most, N non-zero elements within M sequential weights throughout the input channel 052 dimension, N:M sparsity can substantially enhance the performance of structured sparsity, while concurrently assuring swift inference, ably facilitated by the N:M sparse tensor core (Nvidia, 2020).

073

074

075 076



Figure 1: Framework of BAME. It iteratively performs weight pruning-and-regrowing through Loss-Aware Mask adaption (LMA) and Oscillation-aware Block Freezing (OBF), which leads to stable and efficient location for the optimal N:M mask.

077 The crux of maintaining the performance of N:M sparse networks lies in identifying the optimal N:M 078 sparsity mask. To achieve this, prevalent methodologies involve updating all weights during training 079 to determine the most effective N:M mask, adopting a straight-through estimator to approximate the 080 gradients of the pruned weights (Zhou et al., 2021; Zhang et al., 2023b) or learning the importance 081 criteria for all weights (Zhang et al., 2022). Despite their efficacy, the computation of dense gradients invariably imposes a substantial training overhead. Notably, the reduction of training costs has been a 083 focal research point within the sparsity comunity in recent years (Liu et al., 2021; Evci et al., 2020; Dettmers & Zettlemoyer, 2019). With the ever-growing size of cutting-edge models, the significant 084 computational demands and energy consumption of training sparse networks are escalating critical 085 environmental, ethical, and financial concerns. Consequently, the development of efficient and scalable N:M sparse training methods is paramount, potentially even more urgent, to support the 087 widespread accessibility and democratization of DNNs. 088

In this paper, we present BAME as a way of maintaining consistent sparsity in both forward and 089 backward propagation throughout the N:M sparse training process. As shown in Figure 1, BAME 090 escapes from dense weight's update through block-aware N:M mask evolution. It specifically executes 091 weight pruning-and-regrowing within each consecutive M weights in order to adapt the sparse mask. 092 Such mask evolution occurs solely when the detrimental effects on loss caused by pruning a certain weight is outweighed by the gain in loss from restoring another already pruned weight. Concurrently, 094 we selectively adapt the mask of N:M blocks, as some blocks are experimentally observed to exhibit 095 frequent oscillations on their masks during training, leading to unstables loss landscape. To this end, 096 we employ exponential moving averaging (EMA) to accumulate the incidence of mask fluctuations for each block, choosing those with fewer fluctuations for mask evolution to ensure stable optimization 098 for the N:M sparse network during training. In this manner, BAME can stably optimize the N:M 099 mask while conducting N:M sparse training in a dense-backward-free efficient manner.

100 We conduct extensive experiments on validating the effectiveness and efficacy of BAME for N:M 101 sparse training. The results show that BAME is able to get state-of-the-art performance when training 102 N:M sparse networks across a wide range of sparse pattern, datasets, and prevailing DNNs, even 103 with much fewer training FLOPs compared with existing work. Illustratively, BAME attains a 72.0% 104 top-1 accuracy while training a 1:16 sparse ResNet-50 on ImageNet, eclipsing SR-STE (Zhou et al., 105 2021) by 0.5%, while using far less training FLOPs. Our work provides fresh insights in N:M sparse training without dense weight updates and we anticipate that BAME will not only equip practitioners 106 with a robust training tool but also lay the groundwork for subsequent explorations into the training 107 efficiency of N:M sparsity.

# 108 2 RELATED WORK

# 110

# 2.1 NETWORK SPARSITY

111 112

By removing redundant weights to eliminate the parameter and FLOPs burden, network sparsity 113 has emerged as a fervent area of research over the last decade (LeCun et al., 1989; Han et al., 114 2015; Louizos et al., 2017). Traditional approaches can broadly be classified into two categories 115 based on their pruning granularity: unstructured and structured sparsity. The former involves the 116 elimination of individual weights at any location within the network, achieving sparsity at a fine-117 grained level (Han et al., 2015; Lee et al., 2019; Ding et al., 2019). In essence, unstructured sparsity 118 can rival the performance of their dense counterparts even at exceedingly high sparsity ratios, such as 119 90% (Mostafa & Wang, 2019). Nonetheless, the generated sparse weight tensors generally precludes 120 acceleration on standard hardware platforms unless the sparsity ratio reaches or exceeds 95% (Wang). 121 Conversely, structured sparsity achieves notable acceleration by extensively removing entire weight rows or convolution filters (Luo & Wu, 2020; Lin et al., 2020). Regrettably, structured sparsity 122 often leads to substantial performance degradation at sparsity levels exceeding 50%, attributed to 123 the constraints imposed on sparsity flexibility. Diverging from conventional sparsity granularities, 124 this paper delves into N:M sparsity that removes weight in an mid-level granularity and has garnered 125 significant research interest in recent years (Zhou et al., 2021; Sun et al., 2021; Pool & Yu, 2021). 126

127 128

129

# 2.2 N:M SPARSITY

130 The recent development of N:M sparsity upholds the conservation of N-out-of-M consecutive weights 131 in DNNs (Nvidia, 2020; Pool & Yu, 2021; Sun et al., 2021; Zhou et al., 2021; Chmiel et al., 2021; 132 Hubara et al., 2016; Zhang et al., 2022). Supported by the NVIDIA Ampere Core (Ronny Krashinsky, 133 2020), N:M sparsity fosters superior storage and computational efficiency, establishing an immaculate 134 harmony between model efficiency and precision, outdoing both unstructured and structured sparsity. To illustrate, 2:4 sparsity can realize 2× speedups on an NVIDIA A100 GPU, while unstructured 135 sparsity might further decelerate the inference speed at identical levels of sparsity. As trailblazing 136 work, ASP (Nvidia, 2020) employs a traditional tri-phase workflow encompassing model pre-training, 137 high-magnitude weight extraction (Han et al., 2015), and network fine-tuning. Zhou et al. (2021) 138 subsequently proposed to train N:M sparse network from scratch by introducing the Sparse-refined 139 Straight-Through Estimator (SR-STE). More specifically, N-out-of-M weights of higher magnitudes 140 are selected in each forward pass, whileall weights are updated during the backward phase, utilizing 141 the STE estimator, paired with a uniquely designed sparse penalty term. LBC (Zhang et al., 2022) 142 further recasts N:M sparsity as a combinatorial problem, learning the optimal mask for each N:M 143 block. Despite their effectiveness in preserving the performance of sparse networks, most existing 144 works require dense backward propagation to update all weights to discover the optimal N:M mask, 145 leading to massive training burden and memory cost. Our proposed BAME in this paper diverges from existing N:M methods as it performs both sparse forward and backward propagation during the 146 entire training process, substantially alleviating the training cost. 147

148 149

# 2.3 SPARSE TRAINING

150 151

Sparse training, which dynamically adjusts the sparse masks throughout the training process, has re-152 cently emerged as a promising solution to enhance the training efficiency of network sparsity (Hoefler 153 et al., 2021; Evci et al., 2020; Han et al., 2015; Liu et al., 2021). The most representative method RigL 154 (Evci et al., 2020) prunes weights of smaller magnitudes during inference and subsequently regrows 155 the same quantity of weights based on their gradient values throughout backward propagation. Sparse 156 Momentum (Dettmers & Zettlemoyer, 2019) employs the mean momentum magnitude of each layer 157 as a benchmark for redistributing parameters. Kusupati et al. (2020) proffer layer-wise learnable 158 thresholds strategizing the reallocation of parameters across layers. Moreover, Liu et al. (2021) 159 proposed to gradually increase the sparsity level during training to further enhance the performance of sparse networks. While these approaches predominantly concentrate on boosting unstructured spar-160 sity, our endeavor in this paper differs by targeting the training of N:M sparse networks, innovatively 161 designing a block-aware selection mechanism for pruning and reviving N:M sparse weights.

# 162 3 METHODOLOGY

#### 3.1 BACKGROUND

We first recap basic preliminaries of N:M sparsity. For simplicity, we take the weights from a specific layer within DNNs for illustration. N:M sparsity forces at most N out of M consecutive weights in the weight row to have non-zero values. The weights can be therefore grouped into K blocks where each block contains M consecutive weights, denoted as  $\mathbf{W} \in \mathbb{R}^{K \times M}$ . And then, N:M sparsity can be formulated as multiplying W with a binary mask  $\mathbf{B} \in \mathbb{R}^{K \times M}$ , with the following objectives:

164

172

$$\min_{\mathbf{W},\mathbf{B}} \mathcal{L}(\mathbf{W} \odot \mathbf{B}; \mathcal{D}) \quad s.t. \quad \|\mathbf{B}_{k,:}\|_0 = \mathbf{N},$$
(1)

where k = 1, 2, ..., K,  $\odot$  is the point-wise element-wise multiplication,  $\mathcal{L}(\cdot)$  denotes training loss function and  $\mathcal{D}$  represents the observed training dataset, respectively. The zero elements in **B** indicate the removal of corresponding weights in the network, and vice versa.

176 **Challenge of N:M sparse training.** The crux of optimizing Equation (1) falls into locating high-177 quality masks that correctly preserve important weights. As a pioneer work, ASP (Nvidia, 2020) 178 chooses to mask out weights that have lower magnitudes, intuitively reducing the output derivation 179 between dense pre-trained weights and N:M sparse weights. Nevertheless, the pre-training phase unavoidably carries huge training burden. In the literature, a more popular way to obtain the sparse 181 mask is performing training-time weight selection by updating all weights (Zhou et al., 2021; Zhang 182 et al., 2022; Fang et al., 2022; Zhang et al., 2023b). Particularly, the straight-through-estimator 183 (STE) (Bengio et al., 2013) is leveraged to calculate the gradient of all weights, since the currently removed weights always receive no gradient as their corresponding multiplied masks are 0s. Formally, the gradients of  $\mathbf{W}$  are derived as 185

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}} = \frac{\partial \mathcal{L}}{\partial (\mathbf{W} \odot \mathbf{B})} \odot \mathbf{B} \approx \frac{\partial \mathcal{L}}{\partial (\mathbf{W} \odot \mathbf{B})} \odot \mathbf{1}.$$
 (2)

In this vein, all weights can be updated during the training process. By dynamically selecting weights
with higher magnitude, such N:M sparse training can effectively boost the model performance, even
without reliance on pre-trained weights. Despite recent efforts to further enhance N:M sparse training
through additional norm constraints on pruned weights (Zhou et al., 2021) or gradual sparsity (Fang
et al., 2022), one significant concern remains that dense back-propagation and weight updates continue
to incur substantial resource consumption, posing challenges to scenarios with limited resources.

In this paper, we address the above hindrance of training inefficiency by proposing Block-Aware
Mask Evolution (BAME), a method that ensures consistent sparsity throughout the forward and
backward propagation phases of the N:M sparse training process. The unique contribution of BAME
encompasses loss-aware mask adaption (LMA) that prune-and-revive weights to effectively decrease
the training loss, and oscillation-aware block selection (OBS), limiting mask modifications within
blocks demonstrating high-frequency mask oscillations, thus stabilizing the N:M training process.
We meticulously present these two components as follows.

201 202

203

187 188

#### 3.2 LOSS-AWARE MASK ADAPTION

204 Owing to the great benefit of training cost reduction, adapting the sparse mask during training while 205 escaping from dense gradient calculation has been a hot topic within traditional unstructured sparsity 206 literature (Evci et al., 2020; Dettmers & Zettlemoyer, 2019; Liu et al., 2021; Jayakumar et al., 2020). The central philosophy of these methods involves performing a global pruning and revival based on 207 instantaneous gradient information every few training iterations. Specifically, several of the weights 208 with the highest gradients among all pruned weights are restored and the same number of weights 209 with the lowest magnitude among all retained weights are pruned, therefore reducing the loss to the 210 fastest extent. 211

Regrettably, prior methodologies for globally altering the sparse topology are unsuitable within the
context of N:M sparsity. Following a fixed sparsity budget for each N:M block, pruning-and-reviving
of weights can only be carried out in each independent N:M block. This presents substantial risks
for the mask adaptation: The gradients of the weights in the same block are likely to have minor
differences due to the continuous input received, as is the magnitude of the weights. Hence, directly

216 applying traditional sparse methods can have high possibility of resulting in the recovery of weights 217 yielding less loss benefit compared to the disruption caused by weight pruning, even if the pruned 218 weights have the smallest magnitude within the N:M block.

219 To address this challenge, we introduce loss-aware mask adaption (LMA) that ensures weight pruning-220 and-reviving always lead to loss decrease during N:M sparse training. LMA performs static sparse 221 training in both forward and backward propagation, while only calculating dense gradient to perform 222 mask adaption every  $\Delta T$  iteration. Here we use a specific N:M block  $\mathbf{W}_k \in \mathbb{R}^M$  for illustrating 223 the mask adaption procedure. Considering a currently preserved weight  $\mathbf{W}_{k,i}$  where  $\mathbf{B}_{k,i} = 1$ , the 224 loss change, denoted as  $\Delta \mathcal{L}(\mathbf{W}_{k,i})$ , upon its removal can be approximately derived using first-order 225 Taylor expansion (Molchanov et al., 2017) as:

$$\Delta \mathcal{L}(\mathbf{W}_{k,i}) = |\mathcal{L}(\mathbf{W} \odot \mathbf{B}; \mathcal{D}, \mathbf{B}_{k,i} = 0) - \mathcal{L}(\mathbf{W} \odot \mathbf{B}; \mathcal{D}, \mathbf{B}_{k,i} = 1)|$$

$$\approx |\mathcal{L}(\mathbf{W} \odot \mathbf{B}; \mathcal{D}, \mathbf{B}_{k,i} = 1) - \frac{\partial \mathcal{L}}{\partial (\mathbf{W} \odot \mathbf{B})_{k,i}} (\mathbf{W}_{k,i} - 0)$$

$$+ R_1(B_{k,i} = 0) - \mathcal{L}(\mathbf{W} \odot \mathbf{B}; \mathcal{D}, \mathbf{B}_{k,i} = 1)|.$$
(3)

If we ignore the first-order remainder  $R_1(\mathbf{B}_{k,i} = 0)$ , then:

$$\Delta \mathcal{L}(\mathbf{W}_{k,i}) \approx |\frac{\partial \mathcal{L}}{\partial (\mathbf{W} \odot \mathbf{B})_{k,i}} \mathbf{W}_{k,i}|.$$
(4)

Similarly, if we consider reviving a currently removed weight  $\mathbf{W}_{k,i}$  back, the loss change  $\Delta \mathcal{L}(\mathbf{W}_{k,i}) = 0$  can be derived as:

$$\Delta \mathcal{L}(\mathbf{W}_{k,j}) = |\mathcal{L}(\mathbf{W} \odot \mathbf{B}; \mathcal{D}, \mathbf{B}_{k,j} = 1) - \mathcal{L}(\mathbf{W} \odot \mathbf{B}; \mathcal{D}, \mathbf{B}_{k,j} = 0)|$$

$$\approx |\mathcal{L}(\mathbf{W} \odot \mathbf{B}; \mathcal{D}, \mathbf{B}_{k,i} = 0) - \frac{\partial \mathcal{L}}{\partial (\mathbf{W} \odot \mathbf{B})_{k,j}} \left( 0 - (0 - \eta \frac{\partial \mathcal{L}}{\partial (\mathbf{W} \odot \mathbf{B})_{k,j}}) \right)$$

$$+ R_1(B_{k,i} = 1) - \mathcal{L}(\mathbf{W} \odot \mathbf{B}; \mathcal{D}, \mathbf{B}_{k,i} = 0)|$$

$$\approx \eta \left( \frac{\partial \mathcal{L}}{\partial (\mathbf{W} \odot \mathbf{B})_{k,i}} \right)^2,$$
(5)

242 243 244

255 256

236

$$pprox \eta\left(rac{\partial}{\partial(\mathbf{W})}
ight)$$

245 where  $\eta$  is the current learning rate. Based on the preceding derivation, we can articulate the following 246 conclusions. On one hand, Eq. (4) tells that for the preserved weights, pruning those with comparably 247 minor  $\Delta \mathcal{L}(\mathbf{W}_{k,i})$  ensures the loss does not undergo notable alterations. This perspective concurs with 248 traditional network sparsity knowledge (Molchanov et al., 2017; Zhang et al., 2023a). Conversely, 249 considering the presently pruned weights, their revival will invariably benefit the minimization of loss 250 as observed in the derivation of Eq. (5). Simultaneously, it bears mentioning that restoring weights 251 with significantly larger  $\Delta \mathcal{L}(\mathbf{W}_{k,i})$  will induce the most substantial degree of loss mitigation.

252 Therefore, at each mask adaption cycle, we first calculate the loss-aware metric  $\Delta \mathcal{L}(\mathbf{W}_{k,:})$  of all 253 weights in an N:M block using Eq. (4) and Eq. (5). Then, we adapt the mask of weights as follows: 254

$$\bar{\mathbf{B}}_{k,m} = \begin{cases} 0, \text{ if } \Delta \mathcal{L}(\mathbf{W}_{k,j}) < \text{Top}(\Delta \mathcal{L}(\mathbf{W}_{k,:}), \text{M-N}), \\ 1, \text{ otherwise,} \end{cases}$$
(6)

257 where m = 1, 2, ..., M and  $\mathbf{B}$  is the updated mask. Such mask adaptation perceptively prune-and-258 revive weights by looking at the effects imparted on the loss, conducting an inclusive ranking within 259 each N:M block. Paradoxically, preceding arts that mandates the pruning of lowest magnitude 260 weights while refurbishing those with highest gradients (Evci et al., 2020; Zhang et al., 2023a), 261 although justified when executed across the entire weight matrix, may potentially be harmful for N:M sparsity with limited amount of weights in each block. To explain, the increment to the loss prompted 262 by restored weights could indeed be considerably less than the disturbance to the loss distribution 263 induced by pruned weights. Hence, our proposed LMA effectively realizes loss-aware optimization 264 of sparse typologies. 265

266 It is also noteworthy that LMA necessitates the computation of dense gradients only intermittently, every  $\Delta T$  iterations, while primarily conducting truly sparse training and weight updates at other 267 times. This starkly contrasts previous N:M sparse training methods (Zhou et al., 2021; Fang et al., 268 2022; Zhang et al., 2023b), which obligate the calculation and storage of full gradients at each stage 269 of optimization, culminating in a significantly greater training impedance relative to LMA.



Figure 2: Sparse Architecture Divergence (SAD) of N:M blocks during training 2:4, 1:4, and 1:16 sparse ResNet-50 on ImageNet. The majority of blocks remain under minimal mask variation, yet a minority experience frequent mask oscillations.

#### 3.3 OSCILLATION-AWARE BLOCK FREEZING

Other than LMA which enables efficient mask adaption inner each N:M block, we further stabilize the N:M sparse training process by oscillation-aware block freezing (OBF). The impetus behind OBF stems from our observations of the high-frequency mask fluctuation for each block across different LMA cycles. In particular, we employ the Sparse Architecture Divergence (SAD) (Zhou et al., 2021) to calculate mask fluctuation at c-th and c + 1-th LMA cycle as follows:

$$SAD(\mathbf{B}_{k}^{c-1}, \mathbf{B}_{k}^{c}) = \sum_{m=1}^{M} |\mathbf{B}_{k,m}^{c-1} - \mathbf{B}_{k,m}^{c}|,$$
(7)

294 where  $B_{\mu}^{c}$  denote the k-th block's mask at c-th LMA cycle. In Figure 2, we show the accumulated 295 SAD score of different N:M blocks during N:M sparse training. The observation reveals a significantly 296 higher frequency of mask alterations occurring in a certain number of blocks compared to others. 297 On reflecting upon the primary intent of LMA, the apex aim during the LMA process constitutes 298 the pruning of weights of lesser magnitude, making way for the revival of a more significant one, 299 thereby pinpointing an enhanced position while ensuring consistent training thereafter. Nevertheless, some blocks endure recurrent deviations, alternately zeroing the weights, unquestionably inducing 300 oscillations in loss, and thereby impeding network training. Consequently, we harness the capabilities 301 of Exponential Moving Average (EMA) to accumulate the episodes of mask perturbations across each 302 block, electing those exhibiting lesser fluctuations for mask evolution, thereby ensuring stabilized 303 optimization for the N:M sparse network during the course of training. Concretely, we devise a vector 304  $\mathbf{O} \in \mathbb{R}^{K}$ , equivalent in magnitude to the count of blocks, devised for logging the frequency of mask 305 alterations, as 306

$$\mathbf{O}_{k}^{c} = \gamma \ \mathbf{O}_{k}^{c-1} + (1-\gamma) \ SAD(\mathbf{B}_{k}^{c-1}, \mathbf{B}_{k}^{c}), \tag{8}$$

where  $\gamma$  is the momentum of EMA updating. Then, we restrict a  $\beta$  proportion of N:M blocks with the highest oscillation frequency from being updated by LMA as:

$$\mathbf{B}_{k,m}^{c} = \begin{cases} \mathbf{B}_{k,m}^{c-1}, \text{ if } \mathbf{O}_{k}^{c} > \operatorname{Top}(\mathbf{O}^{c}, \lfloor \beta \cdot k \rfloor), \\ \bar{\mathbf{B}}_{k,m}^{c}, \text{ otherwise.} \end{cases}$$
(9)

Furthermore, within the mask adaptation selection of LMA, the occurrence of gradients is sporadic, implying the prospect of a particularly extraordinary gradient for a specified weight at a given stage. This could conceivably prompt an incorrect pruning of a substantial weight, thereby precipitating a noteworthy effect on network performance. Consequently, we confine LMA to transpire solely within weight blocks of lesser magnitudes to circumvent such inadvertent erroneous mask adaptations.

$$\mathbf{B}_{k,m}^{c} = \begin{cases} \mathbf{B}_{k,m}^{c-1}, \text{ if } \mathbf{O}_{k}^{c} > \operatorname{Top}(\mathbf{O}^{c}, \lfloor \beta \cdot k \rfloor) \text{ and } ||\mathbf{W}_{k}^{c}||_{2} > \operatorname{Top}(\hat{\mathbf{W}}^{c}, \lfloor \alpha \cdot k \rfloor), \\ \bar{\mathbf{B}}_{k,m}^{c}, \text{ otherwise,} \end{cases}$$
(10)

319 320

318

307

308

280

281

282 283 284

285

291 292 293

where  $\mathbf{W}_{k}^{c} = ||\mathbf{W}_{k}^{c}||_{2}, k = 1, 2, ..., K$ . For the implementation of BAME, we follow Jayakumar et al. (2020) to perform a three-step sparse training pipline, with  $T_{i}$  and  $T_{f}$  evenly divides the training schedule. In particular, we first employ gradual pruning (Zhu & Gupta, 2017) to set the non-zeros parameters budget linearly decreased from M to the targeted N of each block in the early  $T_{i}$  iterations,

R	equire : Weights $W$ ; Loss function $\mathcal{L}$ ; Initial and final iterations for performing mask
	evolution $t_i$ and $t_f$ ; Update interval $\Delta T$ .
0	utput : Sparse weights $\overline{\mathbf{W}}$
1 <b>f</b> 0	$\mathbf{r} \ t \in [t_i, \ldots, t_f] \ \mathbf{do}$
2	if $t \% \Delta T == 0$ then
3	Calculate $\Delta \mathcal{L}(\mathbf{W})$ via Eq. (4) and Eq. (5)
4	Obtain the adapted mask $ar{ m B}$ via Eq. (6) // Loss-aware mask adaption
5	Get the restricted mask B via Eq. (9) // Oscillation-aware freezing
6	end
7	$ar{\mathbf{W}}$ = $\mathbf{W}\odot\mathbf{B}$ // Apply binary mask to the weights
8	Sparse Forward and backward propagation

Table 1: Results for sparsifying ResNet-32 and MobileNet-V2 on CIFAR-10.

339		1 2	0			
340	Model	Method	N:M	Top-1	Epochs	FLOPs
341			Pattern	Accuarcy (%)	(Train)	(Train)
342	ResNet-32	Baseline	-	94.52	300	1×(3.15e16)
343	ResNet-32	ASP	2:4	94.68	600	$1.5 \times$
244	ResNet-32	SR-STE	2:4	94.52	300	$0.83 \times$
344	ResNet-32	LBC	2:4	94.81	300	$0.72 \times$
345	ResNet-32	BAME(ours)	2:4	94.99	300	0.63×
346	PosNat 22	SD STE	1.4	04.52	200	0.74×
347	ResNet-32	Di Mask	1.4	94.32	200	$0.74 \times$
348	ResNet-32	DI-IVIASK DAME(ours)	1.4	94.43	200	0.49
040	Keshet-52	DAME(ours)	1.4	94./1	300	0.39 ×
349	ResNet-32	SR-STE	1:16	92.92	300	$0.67 \times$
350	ResNet-32	Bi-Mask	1:16	92.77	300	$0.37 \times$
351	ResNet-32	BAME(ours)	1:16	93.15	300	<b>0.29</b> ×
352	MobileNet-V2	Baseline	_	94.55	300	$1 \times (1.41e17)$
353	MobileNet-V2	SR-STE	1:16	93.14	300	0.67×
354	MobileNet-V2	Bi-Mask	1:16	92.48	300	$0.37 \times$
355	MobileNet-V2	BAME(ours)	1:16	93.32	300	<b>0.29</b> ×

which is shown to be effective for performance retention. Then, we perform BAME to find the best N:M mask from  $T_i$  to  $T_f$ . At last, we set the masks all freeze and conduct static training for the N:M sparse network within the remained iterations. The workflow for performing BAME for N:M sparse training is outlined in Alg. 1.

**EXPERIMENT** 

4.1 EXPERIMENTAL SETTINGS

Datasets and Networks. We validate the effectiveness of BAME by using it to train N:M sparse networks on image classification tasks on the CIFAR-10 (Krizhevsky et al., 2009) and ImageNet-1K datasets (Deng et al., 2009). For the networks, we sparsify ResNet-32 (He et al., 2016), MobileNet-V2 on CIFAR-10 dataset, and ResNet-18 (He et al., 2016), ResNet-50 (He et al., 2016), DeiT-small on ImageNet-1K dataset. 

Implementation Details. We train N:M sparse networks from scratch via the Stochastic Gradient Descent (SGD) optimizer, paired with a momentum of 0.9 and a batch size of 256. The initial learning rate is set to 0.1 and gradually decayed based on the cosine annealing scheduler. Following previous works, we train all networks for 300 epochs on CIFAR-10, with a weight decay of 0.005. On ImageNet, 120 epochs are given for ResNet and 300 epochs for DeiT-small. For the implementation of BAME, we set the LMA update interval  $\Delta T = 100$  and 0.5 for both  $\alpha$  and  $\beta$  in OBF. ALL experiments are implemented based on PyTorch and executed on NVIDIA Tesla A100 GPUs.

379	·					
380	Model	Method	N:M	Top-1	Epochs	FLOPs
381			Pattern	Accuarcy (%)	(Train)	(Train)
382	ResNet-50	Baseline	-	77.1	120	$1 \times (3.2e18)$
202	ResNet-50	ASP	2:4	76.8	200	$1.24 \times$
505	ResNet-50	SR-STE	2:4	77.0	120	$0.83 \times$
384	ResNet-50	LBC	2:4	77.2	120	$0.72 \times$
385	ResNet-50	BAME(ours)	2:4	77.4	120	<b>0.63</b> ×
386	ResNet-50	SR-STE	1.4	75.3	120	0.74×
387	ResNet-50	Bi-Mask	1:4	75.6	120	$0.49\times$
388	ResNet-50	BAME(ours)	1:4	76.1	120	0.39×
389	ResNet-50	SR-STE	1:16	71.5	120	0.69×
390	ResNet-50	Bi-Mask	1:16	71.5	120	$0.37 \times$
391	ResNet-50	BAME(ours)	1:16	72.0	120	<b>0.29</b> ×
392	DeiT-small	Baseline	-	79.8	300	1x(8.9e18)
393	DeiT-small	SR-STE	2:4	79.6	300	0.83×
394	DeiT-small	Bi-Mask	2:4	79.4	300	0.72  imes
395	DeiT-small	BAME(ours)	2:4	79.7	300	<b>0.63</b> ×

Table 2: Results for sparsifying ResNet-50 and DeiT-small on ImageNet.

**Performance Metrics and Baselines.** We juxtapose BAME with several state-of-the-art N:M sparsity methods, including ASP (Nvidia, 2020), SR-STE (Zhou et al., 2021), LBC (Zhang et al., 2022), Bi-Mask (Zhang et al., 2023b). We experiment with a wide range of N:M patterns for comparison, including 2:4, 1:4, and 1:16. We report the Top-1 accuracy, the training/inference float-point operations (FLOPs) and parameter burden of N:M sparse networks.

#### 401 402 403

404

396 397

398

399

400

378

# 4.2 IMAGE CLASSIFICATION

405 **CIFAR-10.** We first evaluate the efficacy of BAME for training sparse ResNet-32 and MobileNet-V2 406 on the CIFAR-10 dataset, which includes 50,000 training images and 10,000 validation images 407 within 10 classes. Tab 1 showcases the performance comparison under different N:M patterns. 408 BAME achieves state-of-the-art accuracy at all scenarios, even utilizing far fewer training FLOPs 409 and parameters compared with other methods. For instance, BAME achieves 94.71% top-1 accuracy 410 when training 1:4 sparse ResNet-32, surpassing the recent baseline Bi-Mask that also pursues efficient 411 backward propagation by 0.28%. Moreover, even compared with SR-STE which conducts dense 412 gradient calculation training, BAME still achieves better performance retention for all N:M patterns even with sparse backward propagation. For example, when training 1:16 sparse MobileNet-V2, 413 BAME yields 93.32 top-1 accuracy, surpassing SR-STE by 0.18% while only using 0.29% training 414 FLOPs (0.67% for SR-STE). 415

416 ImageNet. For the large-scale ImageNet-1K dataset that contains over 1.2 million images for training and 50,000 images for validation in 1,000 categories, we first present the quantitative results for 417 training sparse ResNet with depths of 18 and 50, along with DeiT-small in Table 2. Again, BAME 418 substantially enlarges the performance of existing methods at all testing scenarios, with the minimum 419 training FLOPs by efficient weight pruning and growing. For instance, it surpasses SR-STE by 0.5% 420 Top-1 accuracy when training 1:4 sparse ResNet-50 (76.1% for BAME and 75.3% for SR-STE), while 421 consumes far fewer training FLOPs ( $0.39 \times$  for BAME and  $0.74 \times$  for SR-STE). When juxtaposed 422 with Bi-Mask which also focuses on training efficiency, BAME still lead to better performance at 423 all N:M patterns. It is also worth mentioning that BAME holds its advantages when training sparse 424 DeiT-small compared with other methods, demonstrating its scalability for other types of model 425 structures beyond convolution neural networks. These results demonstrated the efficacy of BAME for 426 performing loss-aware mask adaption to efficiently locate the best N:M mask during training.

427 428 429

430

#### 4.3 PERFORMANCE ANALYSIS

431 In this section, we conduct performance analysis of BAME, including its main components and the hyper-parameters setting. All experiments are based on training 1:16 ResNet-32 on CIFAR-10.



Table 3: Results for applying different sparse training schedule to BAME.

$t_i$	$t_f$	Top-1 Accuracy (%)	FLOPs (Train)
0	100	91.87	$0.07 \times$
0	200	92.57	0.07  imes
0	300	92.22	0.07  imes
100	300	92.81	$0.29 \times$
200	300	93.05	$0.41 \times$
100	200	93.15	$0.29 \times$

Figure 3: Results for applying different hyperparameters to BAME.

445 **Hyper-parameters.** We first investigate the influence of hyper-parameters within BAME, including 446 the two restriction factors  $\alpha$  and  $\beta$ , and the updating interval  $\Delta T$ . As shown in Fig.3, the best 447 performance is obtained with  $\Delta T = 100, \alpha = 0.5, \beta = 0.5$ . To analyze, smaller  $\alpha$  and  $\beta$ , larger  $\epsilon$ 448 all lead to an insufficient procedure for mask exploration during the training schedule. Setting these 449 hyper-parameters in the contrast direction, also resulted in poor performance, which fits into our claim that high frequency of mask oscillations can unavoidably harm the training stability and lead to 450 sub-optimal results. Nevertheless, it serves as a promising directions to automatically perform N:M 451 sparse training without hyper-parameter choosen. 452

453 **Training Schedule.** Further, we analyze the training schedule of BAME, *i.e.*  $t_i$  and  $t_f$  for stooping the 454 gradual pruning and performing mask adaption. Tab. 3 delineates the quantitative results. Intuitively, 455 establishing a larger  $t_i$  indicates an increase in training iterations for pre-training with a gradual attainment of the desired sparsity level. Though this consequently induces a significant training 456 cost, neglecting gradual pruning simultaneously results in a considerable performance reduction. To 457 explain, the randomly-initialized weights require a certain degree of pre-training to initiate an effective 458 importance selection, which is validated in traditional sparsity work (Liu et al., 2021; Jayakumar et al., 459 2020). Regarding the mask adaptation schedule, prematurely halting BAME leads to a performance 460 downturn due to inadequate identification of the optimal mask. In stark contrast, prolonging BAME 461 until the termination of training, that is, designating  $t_f$  to the final iteration, results in an even more 462 precipitous performance degradation. To explain, the freshly grown weights, initialized to zeros as 463 per Alg. 1, mandate substantial training following restoration to enhance the performance of the 464 sparse network.

Mask Adaption. At last, we investigate the effectiveness of 466 our mask adaption methods including LMA and OBF. We set 467 static training as the baseline, which means the binary masks 468 are randomly initialized and kept frozen during the entire sparse 469 training procedure. In addition, we run RigL (Evci et al., 2020), 470 a representative method for tradition network sparse training 471 that take weight magnitude and gradient for pruning and re-472 viving, respectively. As shown in Tab. 4, both LMA and OBF contributes to the overall or sparse training performance. 473

Table 4:	Ablation	study	of BAME
----------	----------	-------	---------

Method	Top-1
	Accuracy (%)
Stastic	90.03
RigL	92.01
LMA	92.98
LMA+OBF	93.15
LMA+OBF	93.15

# 5 CONCLUSION

475 476

474

465

477 N:M sparsity has become an increasingly crucial DNN compression tool, delivering functional speed 478 ups by imposing a maximum of N non-zero constituents within M consecutive weights. We introduce 479 BAME, a method that enhances the efficiency of the contemporary N:M sparsity methods while 480 preserving the model's performance. BAME's fundamental principle involves carrying out loss-aware 481 mask adaptation to prune and revitalize weights within specific N:M blocks, whilst maintaining 482 the stability of frequently-oscillating blocks. BAME surpasses existing methods in sparsifying 483 mainstream networks across various vision tasks, all while greatly reducing the training FLOPs and the parameter strain by keeping both sparse forward and backward propagation through training. 484 Hopefully, BAME will not only provide practitioners with a robust N:M sparse training instrument, 485 but also set the groundwork for further investigations into efficient N:M sparsity.

# 486 REFERENCES

510

517

526

527

528

529

539

Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients through
 stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*, 2013.

- Brian Chmiel, Itay Hubara, Ron Banner, and Daniel Soudry. Optimal fine-grained N:M sparsity for activations and neural gradients. In *International Conference on Learning Representations (ICLR)*, 2021.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale
   hierarchical image database. In *IEEE Conference on Computer Vision and Pattern Recognition* (CVPR), pp. 248–255, 2009.
- Tim Dettmers and Luke Zettlemoyer. Sparse networks from scratch: Faster training without losing performance. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- Xiaohan Ding, Xiangxin Zhou, Yuchen Guo, Jungong Han, Ji Liu, et al. Global sparse momentum sgd for pruning very deep neural networks. In *Advances in Neural Information Processing Systems* (*NeurIPS*), pp. 6382–6394, 2019.
- Utku Evci, Trevor Gale, Jacob Menick, Pablo Samuel Castro, and Erich Elsen. Rigging the lottery: Making all tickets winners. In *International Conference on Machine Learning (ICML)*, pp. 2943–2952, 2020.
- 507 Chao Fang, Aojun Zhou, and Zhongfeng Wang. An algorithm–hardware co-optimized framework for
   508 accelerating N:M sparse transformers. *IEEE Transactions on Very Large Scale Integration (VLSI)* 509 Systems, 30(11):1573–1586, 2022.
- Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *IEEE International Conference on Computer Vision (ICCV)*, pp. 580–587, 2014.
- Song Han, Jeff Pool, John Tran, and William Dally. Learning both weights and connections for
   efficient neural network. In *Advances in Neural Information Processing Systems (NeurIPS)*, pp. 1135–1143, 2015.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 770–778, 2016.
- Kaiming He, Georgia Gkioxari, Piotr Dollar, and Ross Girshick. Mask r-cnn. In *IEEE International Conference on Computer Vision (ICCV)*, 2017a.
- Yihui He, Xiangyu Zhang, and Jian Sun. Channel pruning for accelerating very deep neural networks.
   In *IEEE International Conference on Computer Vision (ICCV)*, pp. 1389–1397, 2017b.
  - Torsten Hoefler, Dan Alistarh, Tal Ben-Nun, Nikoli Dryden, and Alexandra Peste. Sparsity in deep learning: Pruning and growth for efficient inference and training in neural networks. *Journal of Machine Learning Research*, 22:1–124, 2021.
- Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand,
   Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for
   mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.
- Itay Hubara, Matthieu Courbariaux, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. Binarized
   neural networks. *Advances in Neural Information Processing Systems (NeurIPS)*, 29, 2016.
- Siddhant Jayakumar, Razvan Pascanu, Jack Rae, Simon Osindero, and Erich Elsen. Top-kast: Top-k
   always sparse training. In *Advances in Neural Information Processing Systems (NeurIPS)*, pp. 20744–20754, 2020.
  - Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.

540 541 542	Aditya Kusupati, Vivek Ramanujan, Raghav Somani, Mitchell Wortsman, Prateek Jain, Sham Kakade, and Ali Farhadi. Soft threshold weight reparameterization for learnable sparsity. In <i>International Conference on Machine Learning (ICML)</i> , pp. 5544–5555, 2020.
543 544 545	Yann LeCun, John Denker, and Sara Solla. Optimal brain damage. In Advances in Neural Information Processing Systems (NeurIPS), pp. 598–605, 1989.
546 547 548	Namhoon Lee, Thalaiyasingam Ajanthan, and Philip Torr. Snip: Single-shot network pruning based on connection sensitivity. In <i>International Conference on Learning Representations (ICLR)</i> , 2019.
549 550 551	Mingbao Lin, Rongrong Ji, Yan Wang, Yichen Zhang, Baochang Zhang, Yonghong Tian, and Ling Shao. Hrank: Filter pruning using high-rank feature map. In <i>IEEE Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pp. 1529–1538, 2020.
552 553 554 555 556	Shiwei Liu, Tianlong Chen, Xiaohan Chen, Zahra Atashgahi, Lu Yin, Huanyu Kou, Li Shen, Mykola Pechenizkiy, Zhangyang Wang, and Decebal Constantin Mocanu. Sparse training via boosting pruning plasticity with neuroregeneration. In <i>Advances in Neural Information Processing Systems (NeurIPS)</i> , 2021.
557 558 559	Zechun Liu, Haoyuan Mu, Xiangyu Zhang, Zichao Guo, Xin Yang, Tim Kwang-Ting Cheng, and Jian Sun. Metapruning: Meta learning for automatic neural network channel pruning. In <i>IEEE International Conference on Computer Vision (ICCV)</i> , pp. 3296–3305, 2019.
560 561 562	Christos Louizos, Max Welling, and Diederik P Kingma. Learning sparse neural networks through <i>l</i> <sub>-</sub> 0 regularization. In <i>International Conference on Learning Representations (ICLR)</i> , 2017.
563 564	Jian-Hao Luo and Jianxin Wu. Autopruner: An end-to-end trainable filter pruning method for efficient deep model inference. <i>Pattern Recognition (PR)</i> , pp. 107461, 2020.
565 566 567 568	Jian-Hao Luo, Jianxin Wu, and Weiyao Lin. Thinet: A filter level pruning method for deep neural network compression. In <i>IEEE International Conference on Computer Vision (ICCV)</i> , pp. 5058–5066, 2017.
569 570 571	Pavlo Molchanov, Stephen Tyree, Tero Karras, Timo Aila, and Jan Kautz. Pruning convolutional neural networks for resource efficient inference. In <i>International Conference on Learning Representations (ICLR)</i> , 2017.
572 573 574 575	Hesham Mostafa and Xin Wang. Parameter efficient training of deep convolutional neural networks by dynamic sparse reparameterization. In <i>International Conference on Machine Learning (ICML)</i> , pp. 4646–4655, 2019.
576 577 578	Nvidia. Nvidia a100 tensor core gpu architecture. https://www. nvidia.com/content/dam/en-zz/Solutions/Data-Center/ nvidia-ampere-architecture-whitepaper.pdf, 2020.
579 580 581	Jeff Pool and Chong Yu. Channel permutations for N:M sparsity. In Advances in Neural Information Processing Systems (NeurIPS), 2021.
582 583 584	Olivier Giroux et al. Ronny Krashinsky. Nvidia ampere sparse tensor core. https://developer. nvidia.com/blog/nvidia-ampere-architecture-in-depth/, 2020.
585 586 587	Wei Sun, Aojun Zhou, Sander Stuijk, Rob Wijnhoven, Andrew O Nelson, Henk Corporaal, et al. Dominosearch: Find layer-wise fine-grained N:M sparse schemes from dense neural networks. In <i>Advances in Neural Information Processing Systems (NeurIPS)</i> , 2021.
588 589 590 591	Ziheng Wang. Sparsert: Accelerating unstructured sparsity on gpus for deep learning inference. In <i>Proceedings of the ACM International Conference on Parallel Architectures and Compilation</i> <i>Techniques (ICPACT)</i> , pp. 31–42.
592 593	Yuxin Zhang, Mingbao Lin, Zhihang Lin, Yiting Luo, Ke Li, Fei Chao, Yongjian Wu, and Rongrong Ji. Learning best combination for efficient N:M sparsity. In <i>Advances in Neural Information Processing Systems (NeurIPS)</i> , 2022.

- Yuxin Zhang, Mingbao Lin, Yunshan Zhong, Fei Chao, and Rongrong Ji. Lottery jackpots exist in
   pre-trained models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023a.
  - Yuxin Zhang, Yiting Luo, Mingbao Lin, Yunshan Zhong, Jingjing Xie, Fei Chao, and Rongrong Ji. Bi-directional masks for efficient n: M sparse training. In *International Conference on Machine Learning*, pp. 41488–41497. PMLR, 2023b.
- Aojun Zhou, Yukun Ma, Junnan Zhu, Jianbo Liu, Zhijie Zhang, Kun Yuan, Wenxiu Sun, and
   Hongsheng Li. Learning N:M fine-grained structured sparse neural networks from scratch. In
   *International Conference on Learning Representations (ICLR)*, 2021.
  - Michael Zhu and Suyog Gupta. To prune, or not to prune: exploring the efficacy of pruning for model compression. In *International Conference on Learning Representations Workshop (ICLRW)*, 2017.