

506 A Supplementary Material

507 This **Supplementary Material** is structured as follows. We provide a formulation of our algorithm
 508 in Section B. To investigate the effectiveness of different components of our SUBP, we conduct
 509 ablation studies and provide additional experimental results in Section C. In Section D, we provide
 510 deployment results on the x86 platform of Intel(R) Xeon(R) Platinum 8260L CPU @ 2.30GHz to
 511 further explore the performance of $1 \times N$ sparse on different platforms. Finally, in Section E, we
 512 discuss the societal impact of our method.

513 B Algorithm Formulation

Algorithm 1: Overview of the SUBP method.

- 1 **Input:** An L -layer CNN model with weights $\mathbf{W} = \{W^i | 1 \leq i \leq L\}$; block binary mask matrices $\mathbf{M} = \{M_{j,k}^i \in \{0, 1\} | 1 \leq i \leq L, 0 \leq j \leq \frac{C_{i+1}}{N} - 1, 0 \leq k \leq C_i - 1\}$; indices of activate blocks with the top scores \mathbf{T} ; indices of regrow blocks based on importance sampling \mathbf{G} ; target prune rate p ; initial regrowing factor δ_0 ; importance balance coefficient λ ; sampling attention balance factor τ ; training epochs T_{total} ; start and end epoch in the pruning-regrowing stage t_s, t_e ; training set \mathcal{D} ;
 - 2 **Output:** A sub-model satisfying the target prune rate p , its optimal weight values \mathbf{W}^* and binary mask \mathbf{M}^* ;
 - 3 Randomly initialize the model weights \mathbf{W} ;
 - 4 Initialize $\{M_{j,k}^i | \forall i, \forall j, \forall k\}$ to 1 ;
 - 5 Reformat \mathbf{W} to \mathbf{B} according to Section 3 ;
 - 6 **for** each training epoch $t \in [T_{\text{total}}]$ **do**
 - 7 Sample a mini-batch from \mathcal{D} and update the model weights \mathbf{W} ;
 - 8 **if** $t_s < t \leq t_e$ **then**
 - 9 Reset $\{M_{j,k}^i | \forall i, \forall j, \forall k\}$ to 1 ;
 - 10 Compute the importance score S of block by Eq. 4 ;
 - 11 Get the indices of activate blocks with the top scores by Eq. 5 ;
 - 12 Prune the bottom-ranking block by set $\{M_{j,k}^i | k \notin \mathcal{T}_j^i\}$ to 0;
 - 13 Compute the importance sampling probabilities by

$$p_{j,k}^i = \exp\left(\frac{S_{j,k}^i}{\tau}\right) / \sum_{m \notin \mathcal{T}_j^i} \exp\left(\frac{S_{j,m}^i}{\tau}\right) ;$$
 - 14 Compute the regrowing factor by Eq. 6 ;
 - 15 Get the indices of regrow blocks based on importance sampling by

$$\mathcal{G}_j^i = \text{Multinomial}(\{p_{j,k}^i | k \notin \mathcal{T}_j^i\}, \delta_i C_i) \text{ without replacement ;}$$
 - 16 Regrow the blocks by resetting $\{M_{j,k}^i | k \in \mathcal{G}_j^i\}$ to 1 ;
-

514 C Ablation Analysis

Table 4: Compare different design choices in the regrowing stages of the SUBP method. All the experiments are based on the TinyImageNet with ResNet18(1×32). The random baseline is 57.0%.

Regrowing factor					
Design choices	$\delta_0 = 0.1$	$\delta_0 = 0.2$	$\delta_0 = 0.3$	$\delta_0 = 0.4$	Full
Top-1	57.6%	58.4%	57.9%	58.0%	58.5%
Decay scheduler for block regrowing					
Design choices	Default	Constant	Linear decay	Cosine decay	
Top-1	58.3%	57.5%	58.4%	58.3%	

515 In Table 4, we investigate the effectiveness of different design choices in our block regrowing stage.
 516 All the experiments are based on the TinyImageNet with ResNet18(1×32). Compared to the random

517 baseline with 57.0% top-1 accuracy, our SUBP achieves consistent improvement under the different
 518 settings.

519 We find that regrowing factor δ_0 significantly impacts the final quality of the model. Intuitively, a
 520 larger regrowing factor can provide a more extensive sampling space during training and retain the
 521 model’s capacity to a greater extent. However, a sizeable regrowing factor may also cause drastic
 522 sub-model structure changes, affecting stability during training. As shown, the accuracy is improved
 523 by 0.8% as the δ_0 increases from 0.1 to 0.2. When δ_0 increases again, the model’s accuracy drops
 524 until δ_0 is the full model size. This suggests that the relationship between the regrowing factor and
 525 the final quality of the model is varied, and selecting an appropriate regrowing factor in specific
 526 circumstances can improve the final quality.

527 We also investigate the decay scheduler for the block regrowing stage. We compare several decay
 528 schedulers, including default (Eq. 6), constant, linear, and cosine. The experiments show SUBP has
 529 good robustness to different decay schedulers, as default, linear, and cosine decay schedulers all show
 530 similar performance. With a decay scheduler, the sampling space can be gradually decreased, and the
 531 sub-model under training can converge stably.

532 D Deployment on x86 Platform

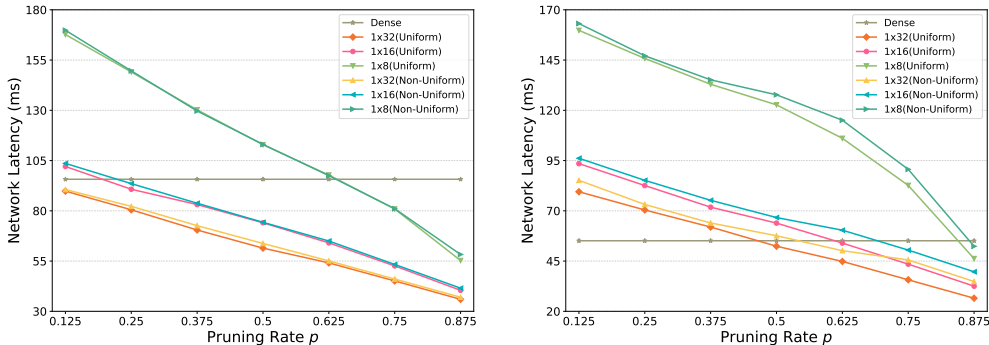


Figure 5: Network latency comparison between uniform $1 \times N$ sparse against non-uniform and dense model with varying N and prune rates. The experiment is conducted using ResNet18 and set the input shape as (4, 3, 224, 224) on the x86 platform of Intel(R) Xeon(R) Platinum 8260L CPU @ 2.30GHz with single thread (left) and two threads (right). Best view in colors.

533 In order to further explore the performance of $1 \times N$ sparse DNNs on different platforms, as shown
 534 in Fig. 5, we also conducted corresponding experiments on the x86 platform of Intel(R) Xeon(R)
 535 Platinum 8260L CPU @ 2.30GHz and obtain the similar results in general: 1) The gain of vanilla
 536 convolution in multithreading scenarios is much greater than that of $1 \times N$ sparse convolution. 2) The
 537 inference speed of uniform $1 \times N$ is slightly faster than that of non-uniform in the case of multithread-
 538 ing, indicating the importance of workload balance again. However, unlike the performance on the
 539 arm platform of Apple M1 Pro CPU @ 3.20GHz, the $1 \times N$ sparse DNNs are significantly accelerated
 540 when N is set to 16 and 32 on the Platinum 8260L CPU @ 2.30GHz. We can also notice that in most
 541 cases, $N=32$ achieves a fast inference speed.

542 E Societal Impact

543 Our method can reduce the computational overhead of training and inferencing stages while achieving
 544 satisfactory accuracy on modern CNN models. This can facilitate the application of CNN models on
 545 edge devices and is of high value for the community and society to realize Green AI. At the same
 546 time, our $1 \times N$ sparse DNNs are based on new sparse operators, which can promote the progress of
 547 related hardware and algorithms to a certain extent.