
Training Recipe for N:M Structured Sparsity with Decaying Pruning Mask

Sheng-Chun Kao*, Amir Yazdanbakhsh*[†], Suvinay Subramanian[‡]

Shivani Agrawal[†], Utku Evci[†], Tushar Krishna

Georgia Tech [†]Google Research, Brain Team [‡]Google

felix@gatech.edu, ayazdan@google.com, suvinay@google.com

shivaniagrawal@google.com, evcu@google.com, tushar@ece.gatech.edu

(*Equal Contribution)

Abstract

Sparsity has become one of the promising methods to compress and accelerate Deep Neural Networks (DNNs). Among different categories of sparsity, structured sparsity has gained more attention due to its efficient execution on modern accelerators. Particularly, N:M sparsity is attractive because there are already hardware accelerator architectures that can leverage certain forms of N:M structured sparsity to yield higher compute-efficiency. In this work, we focus on N:M sparsity and extensively study and evaluate various training recipes for N:M sparsity in terms of the trade-off between model accuracy and compute cost (FLOPs). Building upon this study, we propose two new decay-based pruning methods, namely “pruning mask decay” and “sparse structure decay”. Our evaluations indicate that these proposed methods consistently deliver state-of-the-art (SOTA) model accuracy, comparable to unstructured sparsity, on a Transformer-based model for a translation task. The increase in the accuracy of the sparse model using the new training recipes comes at the cost of marginal increase in the total training compute (FLOPs).

1 Introduction

Deep Neural Networks (DNNs) have shown success in many domains such as computer vision, language modeling, machine translation, and so on. An trend of SOTA DNN models is that the model size increases quickly with time. For example T5 from Google [42], OPT from Meta [56] and GPT-3 from OpenAI [5] have over 100 billions parameters, making them hard to be deployed and inaccessible for many practitioners with limited compute resources. Another line of effort in the DNN community is to propose different methods to compress the models, such as quantization [45, 25, 55, 57, 53], sparsification [11, 18, 17, 21, 36, 52, 60, 15, 37, 38, 8, 39, 23, 9], and distillation [44, 22, 46, 49].

In this paper, we focus on sparsification (or pruning), which prunes a portion of the parameters in the model by setting their values to 0. It can reduce the amount of compute by skipping multiplications with 0, reduce memory usage by using compressed sparse representations such as COO, CSR, and so on [41], and save energy/power by reducing memory accesses and computations. It opens up the possibility of deploying a large model in resource-limited devices. However, sparsification is often about trading-off between model quality¹ and compression ratio. For example, many studies show promising results in sparsifying image classification models to around 90%-95% sparsity (5%-10% density) without quality loss [17, 19]. With the success of Transformers in natural language processing, there is rising interest in investigating sparsification in Transformer models, where around 80%-90% sparsity can be achieved. Sparsification in language models has huge potential benefits

¹In this paper, we refer to algorithmic-wise criteria such as accuracy, recall, and precision as *model quality*; we refer to model runtime/latency as *model performance*.

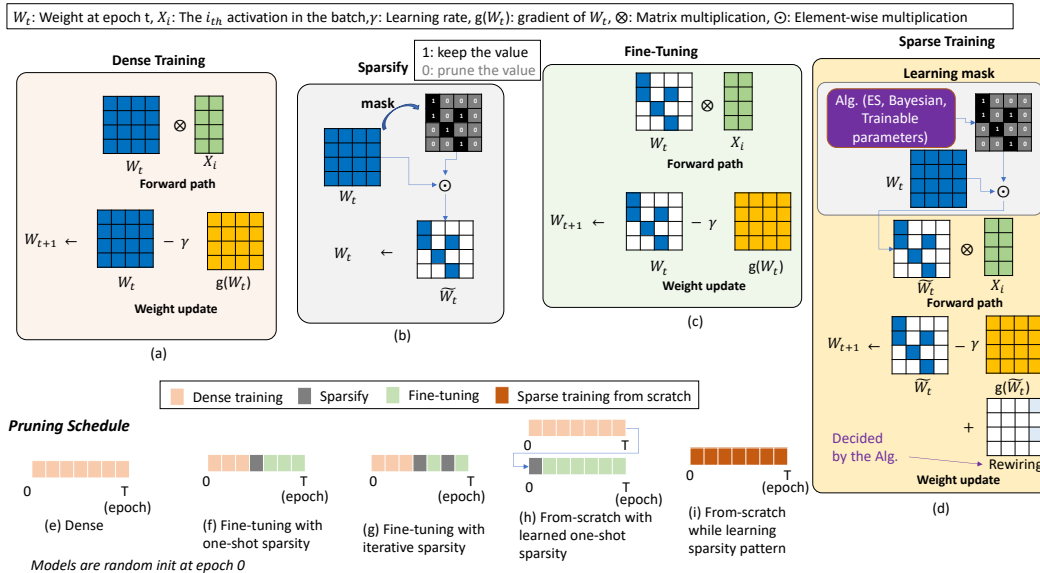


Fig. 1: The compute flows of (a) Dense training, (b) Sparsify, (c) Fine-tuning, and (d) Sparse training. The training schedule of (e) regular dense training, (f) fine-tuning with one-shot sparsifying, (g) fine-tuning with iterative sparsifying, (h) from-scratch with learned one-shot sparsity pattern, and (i) from-scratch while learning sparsity pattern. The sparsify algorithm in (d): ES such as [34], Bayesian Optimization such as [4], Trainable parameters such as [51, 6, 26].

especially in encoder-decoder tasks such as translation. Since decoder needs to be run iteratively for all N tokens in a sequence, even minor performance improvements in the decoder can improve performance significantly. In this paper, we demonstrate our method with an encoder-decoder Transformer-based translation model.

While sparsification can effectively reduce the memory requirement, generally leveraging induced (unstructured) sparsity in the model for higher performance improvements is challenging. The irregularity of the sparsity pattern makes it challenging to be effectively leveraged by the dense accelerators such as GPU and TPU. The sparsified models often ends up with similar or worse performance (because of the extra complexity to compress and decompress the parameters) than their dense counterparts [2, 32, 43, 21, 30, 15, 59, 50, 10].

To this end, structured sparsity, which regularizes the sparsity pattern such as channel/filter sparsity [29, 50, 21], or block sparsity [32, 47], have become increasingly popular owing to their hardware-friendly nature. For example the dense accelerator can skip a full channel computation when it is sparsified without any low-level modification. The caveat is structured sparsity also introduces larger quality loss. Recent research [52, 24] found fine-grained N:M structured sparsity, which keeps N out of consecutive M elements in the the weight tensor, can ameliorate the quality loss. Moreover, with the launch of 2:4 structured-sparse tensor core in GPU Ampere architecture [2] developing sparse training recipes for N:M sparsity has acquired increased interest [40, 33, 3, 58].

In this paper, we demonstrate a training recipe for N:M structure sparsity in Transformer-based translation task and propose two techniques. We propose *Structure Decay* an iterative pruning approach tailored for N:M sparsity. We propose *Mask Decay*, which gradually decays the mask from 1, to 0.9, 0.8, ..., to 0, instead of the conventional 1/0 mask. We found these techniques can stabilize the training and achieve better quality and compression rate. We make following contributions:

- We compare *Structure Decay* and *Mask Decay* with the state-of-the-art N:M sparsity training recipes, SR-STE [58]. They achieve (geomean) 0.004 and (geomean) 0.006 accuracy improvement over SR-STE, respectively.
- *Mask Decay* enables “structured pruning” to achieve comparable quality and compression rate to “unstructured pruning”.

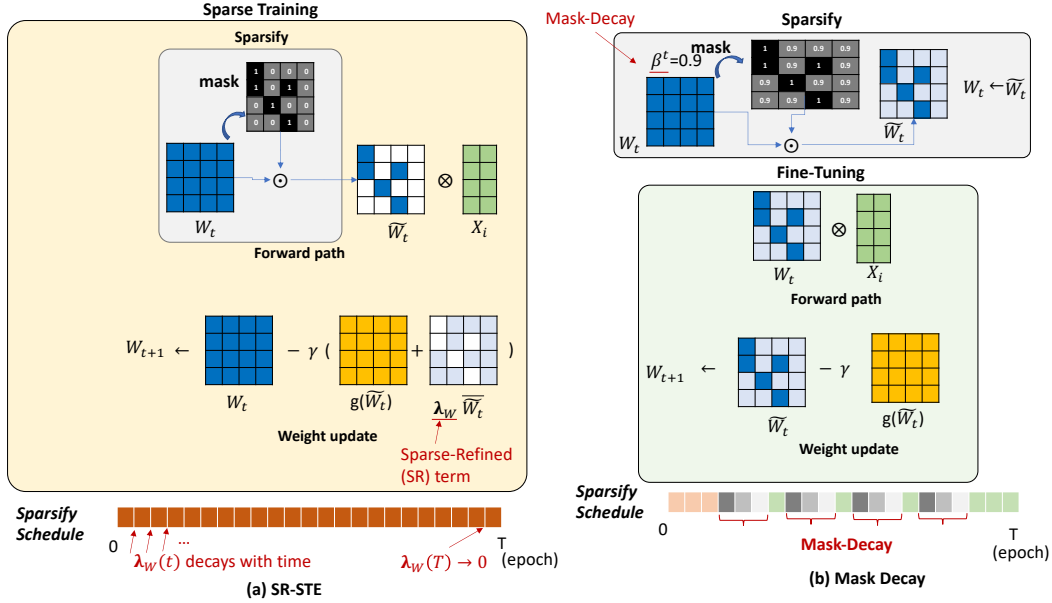


Fig. 2: The weight update scheme of (a) SR-STE [58] and (b) Mask-Decay.

2 Related Work

We primarily focus on weight sparsification in this work. A sparsification recipe includes: 1) pruning criteria, 2) pruning schedule, and 3) sparsity pattern.

Pruning criteria. Pruning criteria is the criteria to decide which elements to prunes in the weight tensor. Magnitude pruning, which selects the pruning elements by their absolute values, is the most widely used method [43, 17, 28, 12, 14, 60, 18, 31, 40, 33]. Some other metrics such as gradient-based [54, 9], Hessian based [27], connection sensitivity [28], and salient-based [35, 28] are also used. In this paper, we use magnitude pruning.

Pruning schedule. There are coarsely four different pruning schedules. 1) *Fine-tuning with one-shot pruning* (Figure 1f) [33, 40, 12, 28], which trains a dense models, prunes the weight with one-shot, and re-trains the model in order to recover the quality loss. 2) *Fine-tuning with iterative pruning* (Figure 1g) [11, 18, 17, 21, 36, 52, 60, 15, 37, 38, 8, 39, 23, 9], which trains a dense model and then iterates between pruning and re-training. This schemes are usually found to has higher ability to recover the quality loss. 3) *From-scratch with learned one-shot pruning pattern* (Figure 1h) [13, 11], which determines the sparsity pattern from the trained dense version and trains a sparse model from scratch. 4) *From-scratch while learning sparsity pattern* (Figure 1i) [51, 6, 14, 26, 9, 4, 34, 58, 10], which trains a sparse model from scratch while learning sparsity patterns simultaneously.

Sparsity pattern. *Unstructured Sparsity* prunes the model without any sparsity pattern constraint [43, 17, 28, 12, 14, 60, 18, 31, 51, 6, 14, 26, 9, 4, 34]. It is often found to be able to prune the model size to an order of magnitude smaller while keeping the model quality. However, it has the challenge of similar or worse (because of the additional complexity) runtime than the dense model owing to its irregular sparsity pattern. *Coarse-grained Structured Sparsity* constrains the pruning scheme to prune the model in a coarse-grained manner such as filter/channel pruning [29, 50, 21], block-wise pruning [50, 32, 38, 16], and so on. By skipping the full computation at a coarse-granularity of the computation, this scheme can often achieve speedup in dense computation accelerators such as GPUs and TPUs; however this often sacrifices some model quality. These studies often trade off between performance and quality for different application needs. *Fine-grained N:M Structure Sparsity*, which prunes (M-N) out of consecutive M elements. Some early works rely on special threading and grouping techniques [52] or specialized sparse accelerators [24] to leverage this fine-grained pattern. With the 2:4 structured-sparse GEMM support in tensor cores in GPU Ampere architecture [2], many recent works start to investigate in different training recipe for N:M sparsity pattern to leverage the existing hardware [40, 33, 3, 58].

Proposed recipe and the SOTA SR-STE [58]. In this paper, our sparse training recipe is (Pruning criteria: magnitude pruning, Pruning schedule: fine-tuning with iterative pruning, Sparsity pattern: fine-grained N:M structure sparsity). The recipe for SR-STE [58] is (Pruning criteria: magnitude pruning, Pruning schedule: from-scratch with learning iterative pruning, Sparsity pattern: fine-grained N:M structure sparsity). Our methods, Mask Decay and Structure Decay, are techniques to improve the training quality of “fine-tuning with iterative pruning”. SR-STE [58] is a “Sparse-Refined (SR)” technique to stabilize the training of “from-scratch with learning iterative pruning”. Both SR-STE [58] and us are proposing techniques to pursue high quality sparsification for fine-grained N:M structure sparsity pattern.

3 Methodology

Table 1: The compute and memory contributions of the three major layers in Transformers. Einsum: computation of attention scores and weighted sum of values by the attention scores. Projections: projecting inputs to key, query, and value, and projection weighted sum of values to outputs. Feed Forward: The multiple feed forward layers at the end of the attention layer. The other layers/operations such as ReLU, LayerNorm, Add, Softmax, embedding, and so on have little contributions to the FLOPS and parameters of the Transformers, hence not included in this estimation. The feed forward layers account for around 64% of overall FLOPs and 67% of parameters. These estimations are made under model configuration in Table 2.

	Einsum	Projections	Feed Forward	Einsum	Projections	Feed Forward
(T)FLOPS	1.6	13.2	26.4	4%	32%	64%
Params (MB)	0.0	50.3	100.7	0%	33%	67%

Workload and model. We evaluate different sparsification methods on the WMT translate task [1] that uses a Transformer-based model [48], and is a key benchmark in machine translation research. They hold several translation datasets across different languages. The encoder and decoder blocks in this model each have six attention layers with 16 heads. The embedding dimension for both input and query/key/value are 1024. The feed-forward blocks within each attention head has 4096 neurons. For all the experiments, we *only* induce sparsity in the feed-forward layers of both encoder and decoder blocks (Table 1 shows that feed-forward layers account for around 64% of FLOPS and 67% of parameters of the entire model. Therefore, we focus on feed-forward layers for sparsification). We follow the standard practice of fine-tuning using the final learning rate used during the original training phase [31].

Training details. Table 2 shows the details of training hyperparameters that we use for all the evaluations. For each experiment, we use a TPUv3 with 32 cores.

Table 2: Model configurations and hyper-parameters.

number of encoder layers	6
number of decoder layer	6
hidden dimension size	1024
feed forward dimension size	4096
number of head	16
max sequence length	256
training set	WMT-17
testing set	WMT-14
learning rate	0.0625
warmup steps	1000
decay factor	0.5
batch size	512
training steps	200K
Adam optimizer	beta1 = 0.9, beta2 = 0.92

3.1 Sparsification Method Baseline

Fine-grained N:M sparsity. We follow the proposed method in [58] to induce structured N:M sparsity from scratch, as shown in Figure 2(a). This method employs standard online magnitude-based pruning with an introduced sparse-refined regularization term. This regularization term applies refined gradients for pruned weights during backward pass. The authors use the refined gradient updates to increase the likelihood of pruning the same network weights at each training step, which purportedly leads to a more robust sparse training. While this work proposes to re-evaluate the pruning mask after each training iteration, we find this process time-consuming which significantly slows down the training process on TPU. Therefore, we moderately alter the frequency of updating the pruning mask to 1000 training steps. We ablate the importance of the frequency of updating pruning masks. Our results show that the model accuracy for WMT task is not sensitive to this parameter, as shown in Table 3.

Table 3: The effect of update frequency in SR-STE [58]. Raising the “Update Frequency” increases the training time significantly. Hence, we only include the results for “Update Frequency” 100 and 1000.

Accuracy	Update Frequency		
	every 1000 steps	every 100 steps	
Sparsity Target	1:16	0.709	0.710
	1:32	0.707	0.707
	1:64	0.706	0.706
	1:128	0.706	0.706

3.2 Proposed Sparsification Methods

In this section, we propose two sparsification methods that employ a decaying mechanism to gradually induce the target sparsity on the model. Note that, we do not alter the gradient update rule in either of these proposed methods. Instead, we simply employ various gradual update rules to the pruning mask itself.

Pruning mask decay. In the first approach, instead of using a binary pruning mask (e.g. “0” indicates pruning locations), we use a floating-point pruning mask with decaying, as shown in Figure 2(b). At the start of training, we employ an all-ones matrix as the pruning mask that simply indicates no pruning. At the beginning of sparse training phase, we use the same standard online magnitude-based pruning criteria to identify the locations of pruned weights. However, in contrast to prior work in which “0” is used to prune the weights, we use $0 < \beta < 1.0$ for the pruned weights in the pruning mask. We gradually decrease the value of β at different intervals following the formula β^d , where d indicates the decaying iteration index (e.g. $\beta^1, \beta^2, \beta^3, \dots$). After sufficient decaying intervals, we set β to zero to indicate the locations of pruned weights. We postulate that using a non-binary pruning mask enables the gradients of pruned weights to flow through the network leading to a more robust sparse training and better model performance.

Sparse structure decay. In the second proposed sparsification method, we apply a decaying mechanism on the structure of pruning mask, gradually increasing sparsification degree. At the beginning of sparse training phase, we start with M-1:M structured sparsity. As training progresses, we increase the sparsification degree by applying $\frac{M}{2^d} : M$ structured sparsity at different decaying intervals. Similar to previous method, d denotes the decaying iteration index. This method at its crux follows a similar hypothesis as the pruning mask decay. That is, enabling the gradient of pruned weights to flow through the network. However, because we still use a binary pruning mask, the contribution of the gradients of the pruned weights to the network reduces after each decaying interval.

4 Evaluation

4.1 Methodology

Task. We use translation as our target task. We use WMT dataset (En-De) [1].

Comparisons. In this study, we compare the effectiveness of different SOTA sparsification methods. All methods train for n steps.

- **Dense:** Dense training without sparsification for n steps (Figure 1(e)).
- **Dense-sparse:** Dense train for d steps, sparsify, and fine-tune for $(n-d)$ steps, as in [33] (Figure 1(f)).
- **Sparse:** SR-STE [58]-based sparse training for n steps (Figure 2(a)).
- **Structure Decay:** Dense train for d steps, structure decay the sparsity pattern for $(n-d-s)$ steps, and fine-tune for s steps (Figure 1(g)). The structure decay is set to decay by the power of 2 (§3.1). For example, when target sparsity pattern is 1:16, we divide $(n-d-s)$ steps to five equal time frame, and the sparsity pattern of each time frame is 15:16, 8:16, 4:16, 2:16, and 1:16, respectively.
- **Mask Decay:** Dense train for d steps, mask decay the sparsity pattern for $(n-d-s)$ steps, and fine-tune for s steps (Figure 2(b)). We use the mask decay rate (β) of 0.9 and mask update period of 1000 steps. In the above experiments, we use $n = 200\text{K}$, $d = 20\text{K}$, $s = 20\text{K}$.

4.2 Comparing with Baseline

Quality. We compare Structure Decay and Mask Decay, with two baseline Dense-Sparse [33] and Sparse [58] in Table 4. We evaluate the methods on different sparsity targets. Sparsity target is the final sparsity pattern we will achieve after model training. For example, sparsity target of 1:32 means the trained model will have only 1 non-zero parameters every 32 parameters. Table 4 shows that Dense-Sparse performs similarly to Sparse, and Mask Decay achieves the best accuracy across all sparsity targets. Structure Decay performs the second best. More interestingly, Mask Decay can help achieve similar or better accuracy than the “unstructured sparsity” ones.

Our results indicate that the “Mask Decay” pruning method on dense layers enables models to be pruned structurally while achieving comparable or even better accuracy to “unstructured pruning”.

Table 4: Comparisons between different sparsification strategies.

Accuracy		Dense	Structure Sparsity			Unstructure Sparsity		
Schedule	Dense	Structure Decay	Mask Decay	Dense-Sparse [33]	Sparse [58]	Dense-Sparse	Unstr Sparse	
Sparsity Target	1:16	0.747	0.717	0.717	0.714	0.709	0.714	0.714
	1:32	0.747	0.713	0.714	0.710	0.707	0.711	0.712
	1:64	0.747	0.710	0.711	0.708	0.707	0.711	0.711
	1:128	0.747	0.708	0.711	0.708	0.707	0.708	0.709

Performance. Note that there are no off-the-shelf accelerators that can support 1:16 or more aggressive sparsity patterns. To demonstrate the potential performance, we build a cost model to estimate the FLOPS and memory sizes. Table 5 shows that after 1:16 sparsification, the model sizes will reduce by 62% and inference FLOPS will reduce by 60%. More interestingly, since different methods have different sparsification schedules, the averaged training FLOPS across their training time will be different. “Dense-Sparse” is the most straight-forward and light-weighted training schedule in terms of FLOPS. “Sparse” relies on continuous mask update across the full training steps (Figure 2(a)), therefore performing worse in terms of FLOPS. Structure Decay becomes the most performant method to achieve the best quality at 1:16 sparsity. However, to sparsify more aggressively, Mask Decay is still the best method (Table 4).

Table 5: Comparisons of model performance (FLOPS) and quality (accuracy).

Sparsity Target (1:16)	Dense	Structure Decay	Mask Decay	Dense-Sparse	Sparse
Params (MB)	151.0	56.6	56.6	56.6	56.6
Inference TFLOPS	41.2	16.5	16.5	16.5	16.5
Training TFLOPS	123.7	108.0	121.2	101.4	123.7
Accuracy	0.75	0.72	0.72	0.714	0.709

4.3 Ablation Studies

Dense training v.s. training from scratch for SR-STE. SR-STE uses sparse training from scratch. All the other methods that we evaluated have a dense training phase at the first few steps (epochs).

This recipe has been proven to be effective as shown in the previous experiments and many prior works [33, 40, 12, 28, 11, 18, 17, 21, 36, 52, 60, 15, 37, 38, 8, 39, 23, 9]. Therefore, we experiment on adding a dense training phase at the beginning of SR-STE training, as shown in Table 6. We found that adding few steps of dense training (1.25% - 10% of the total training steps) can increase the accuracy by around 0.002 to 0.003. This tells that few steps of dense training does help achieve better performance even for SR-STE. Interestingly, the improved SR-STE becomes competitive to the proposed Structure Decay. However, Mask Decay is still consistently better.

Table 6: Ablation: SR-STE augmented with few epochs of dense training.

	Training Schedule		Accuracy			
	Dense steps	SR-STE -styled Sparse steps	1:16	1:32	1:64	1:128
SR-STE	0	200K	0.710	0.707	0.706	0.706
Dense + SR-STE	2.5K	197.5K	0.712	0.710	0.708	0.706
	5K	195K	0.712	0.709	0.707	0.708
	10K	190K	0.712	0.710	0.707	0.708
	20K	180K	0.713	0.710	0.708	0.707

Effect of dense training steps (d). Both our proposed methods include a dense training phase. We do an ablation study on different number of dense training steps in Table 7. We found that changing the dense step between 1.25% - 10% of the total training steps does not observably change the accuracy performance. However, empirically, we found that dense training phase is still essential. The model cannot achieve as competitive accuracy without few epochs of dense training.

Table 7: Ablation: The effect of number of dense training steps (d).

Sparsity Target	Accuracy	Mask Decay				Structure Sparsity			
	1:16	1:32	1:64	1:128	1:16	1:32	1:64	1:128	
Dense steps (d)	2.5 K	0.7155	0.7134	0.7106	0.7100	0.7157	0.7134	0.7108	0.7106
	5 K	0.7160	0.7127	0.7110	0.7093	0.7160	0.7136	0.7117	0.7100
	10 K	0.7157	0.7137	0.7103	0.7094	0.7164	0.7141	0.7107	0.7098
	20 K	0.7156	0.7126	0.7107	0.7104	0.7165	0.7128	0.7115	0.7107

Effects of fine-tuning steps (s). We also have a sets of study on number of fine-tuning steps in Table 8. We found that for both of our propose methods the fine-tuning steps between 10% - 20% of the total training steps does not observably change the accuracy performance. However, empirically, we also found few steps of fine-tuning at the end is essential to recover the accuracy.

Table 8: Ablation: The effect of number of fine-tuning steps (s).

Sparsity Target	Accuracy	Mask Decay				Structure Sparsity			
	1:16	1:32	1:64	1:128	1:16	1:32	1:64	1:128	
Fine-tuning steps (s)	20 K	0.7153	0.7130	0.7107	0.7098	0.7160	0.7125	0.7095	0.7072
	40 K	0.7161	0.7132	0.7106	0.7097	0.7121	0.7093	0.7081	0.7065

Effects of β in mask decay. Note that the extreme case of a mask decay rate ($\beta=0$) will turn the pruning mask back to the conventional 1/0 mask. As shown in Table 9, we found a mask decay rate of 0.9 is better than an aggressive one (0.001). It tells the Mask Decay technique does contribute and lead to better accuracy performance.

Table 9: Ablation: The effect of mask decay rate (β).

Sparsity Target	Accuracy	Mask Decay			
	1:16	1:32	1:64	1:128	
Mask decay rate (β)	0.9	0.715	0.713	0.711	0.710
	0.001	0.712	0.709	0.708	0.707

5 Limitations

This paper studies only translation task with one sparsification recipe (Criteria: magnitude pruning, Schedule: Structure Decay or Mask Decay, Pattern N:M structured pruning) and only prunes feed-forward layers. 1) We show we can effectively prune the most compute- and parameter-heavy layers in our model, feed forward layers. An interesting next-step is to prune other layers such as projections layers as well to further compress the model. 2) Studies on other language tasks or visions tasks and on different models such as on Resnet [20] or ViTs [7]) would be an interesting follow-up. 3) In addition, we only study one combination of sparsification recipe. We might discover better recipe by exploring other combinations such as salient-based pruning + Mask Decay + N:M structured pruning, magnitude pruning + Structure Decay + unstructured sparsity, or many others. 4) Lastly, in the evaluations, most of the hyper-parameters are set manually. We did a limited scope of hyper-parameters sweep in §4.3. A full-fledged hyper-parameter search might discover more performance improvement in Mask Decay and Structure Decay.

6 Conclusion

In this work, we study and evaluate various training recipes for N:M structured sparsity. Building on this study, we propose and compare two new training recipes for N:M structured sparsity based on decaying mechanisms. We study the trade-off between model accuracy and training compute cost (FLOPs) across these training recipes. We show that gradual decay of pruning mask values consistently yield better model accuracy, on-par with unstructured sparsity, on translate task at the cost of modest increase in the training compute cost. While structured sparsity seems to be better positioned for hardware acceleration, its associated training cost should not be overlooked. This work represents a first step in evaluating training recipes for structured sparsity from the perspective of trade-off between model accuracy and compute cost. As future work, we plan to expand the pool of models to other tasks and models and develop a platform to systematically evaluate and compare various sparse training recipes both for model accuracy and training cost.

Acknowledgements

We would like to extend our gratitude towards Jeremiah Willcock, Penporn Koanantakool, Chandu Thekkath, and our extended team at Google Research, Brain Team.

References

- [1] EMNLP 2017 Second Conference on Machine Translation (WMT17). <https://www.statmt.org/wmt17/>, 2017.
- [2] NVIDIA Ampere Architecture Whitepaper. <https://images.nvidia.com/aem-dam/en-zz/Solutions/data-center/nvidia-ampere-architecture-whitepaper.pdf>, 2021.
- [3] NVIDIA ASP (Automatic Sparsity). <https://github.com/NVIDIA/apex/tree/master/apex/contrib/sparsity>, 2021.
- [4] Guillaume Bellec, David Kappel, Wolfgang Maass, and Robert Legenstein. *Deep Rewiring: Training Very Sparse Deep Networks*. *arXiv preprint arXiv:1711.05136*, 2017.
- [5] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. *Language Models are Few-Shot Learners*. *NeurIPS*, 2020.
- [6] Tim Dettmers and Luke Zettlemoyer. *Sparse Networks from Scratch: Faster Training without Losing Performance*. *arXiv preprint arXiv:1907.04840*, 2019.

- [7] 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](#). *ICLR*, 2020.
- [8] Erich Elsen, Marat Dukhan, Trevor Gale, and Karen Simonyan. [Fast Sparse Convnets](#). In *CVPR*, 2020.
- [9] Utku Evci, Trevor Gale, Jacob Menick, Pablo Samuel Castro, and Erich Elsen. [Rigging the Lottery: Making all Tickets Winners](#). In *ICML*, 2020.
- [10] Utku Evci, Yani A Ioannou, Cem Keskin, and Yann Dauphin. [Gradient flow in Sparse Neural Networks and how Lottery Tickets Win](#). *AAAI*, 2020.
- [11] Utku Evci, Fabian Pedregosa, Aidan Gomez, and Erich Elsen. [The Difficulty of Training Sparse Neural Networks](#). *arXiv preprint arXiv:1906.10732*, 2019.
- [12] Jonathan Frankle and Michael Carbin. [The Lottery Ticket Hypothesis: Finding Sparse, Trainable Neural Networks](#). *arXiv preprint arXiv:1803.03635*, 2018.
- [13] Jonathan Frankle, Gintare Karolina Dziugaite, Daniel M Roy, and Michael Carbin. [Pruning Neural Networks at Initialization: Why are we Missing the Mark?](#) *arXiv preprint arXiv:2009.08576*, 2020.
- [14] Trevor Gale, Erich Elsen, and Sara Hooker. [The State of Sparsity in Deep Neural Networks](#). *arXiv preprint arXiv:1902.09574*, 2019.
- [15] Noah Gamboa, Kais Kudrolli, Anand Dhoot, and Ardavan Pedram. [Campfire: Compressible, Regularization-free, Structured Sparse Training for Hardware Accelerators](#). *arXiv preprint arXiv:2001.03253*, 2020.
- [16] Scott Gray, Alec Radford, and Diederik P Kingma. [GPU Kernels for Block-sparse Weights](#). *arXiv preprint arXiv:1711.09224*, 2017.
- [17] Yiwen Guo, Anbang Yao, and Yurong Chen. [Dynamic Network Surgery for Efficient DNNs](#). *NeurIPS*, 2016.
- [18] Song Han, Huizi Mao, and William J Dally. [Deep Compression: Compressing Deep Neural Networks with Pruning, Trained Quantization and Huffman Coding](#). *arXiv preprint arXiv:1510.00149*, 2015.
- [19] Song Han, Jeff Pool, John Tran, and William Dally. [Learning both Weights and Connections for Efficient Neural Network](#). *NeurIPS*, 2015.
- [20] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. [Deep Residual Learning for Image Recognition](#). In *CVPR*, 2016.
- [21] Yihui He, Xiangyu Zhang, and Jian Sun. [Channel Pruning for Accelerating Very Deep Neural Networks](#). In *ICCV*, 2017.
- [22] Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and Qun Liu. [TinyBERT: Distilling BERT for Natural Language Understanding](#). *arXiv preprint arXiv:1909.10351*, 2019.
- [23] Nal Kalchbrenner, Erich Elsen, Karen Simonyan, Seb Noury, Norman Casagrande, Edward Lockhart, Florian Stimberg, Aaron Oord, Sander Dieleman, and Koray Kavukcuoglu. [Efficient Neural Audio Synthesis](#). In *ICML*, 2018.
- [24] Hyeong-Ju Kang. [Accelerator-aware Pruning for Convolutional Neural Networks](#). *IEEE Transactions on Circuits and Systems for Video Technology*, 2019.
- [25] Sehoon Kim, Amir Gholami, Zhewei Yao, Michael W Mahoney, and Kurt Keutzer. [I-BERT: Integer-only BERT Quantization](#). *ICML*, 2021.

- [26] Aditya Kusupati, Vivek Ramanujan, Raghav Somani, Mitchell Wortsman, Prateek Jain, Sham Kakade, and Ali Farhadi. [Soft Threshold Weight Reparameterization for Learnable Sparsity](#). In *ICML*, 2020.
- [27] Yann LeCun, John Denker, and Sara Solla. [Optimal Brain Damage](#). *NeurIPS*, 1989.
- [28] Namhoon Lee, Thalaiyasingam Ajanthan, and Philip HS Torr. [SNIP: Single-shot Network Pruning based on Connection Sensitivity](#). *arXiv preprint arXiv:1810.02340*, 2018.
- [29] Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. [Pruning Filters for Efficient Convnets](#). *arXiv preprint arXiv:1608.08710*, 2016.
- [30] Mingbao Lin, Liujuan Cao, Shaojie Li, Qixiang Ye, Yonghong Tian, Jianzhuang Liu, Qi Tian, and Rongrong Ji. [Filter Sketch for Network Pruning](#). *IEEE Transactions on Neural Networks and Learning Systems*, 2021.
- [31] Zhuang Liu, Mingjie Sun, Tinghui Zhou, Gao Huang, and Trevor Darrell. [Rethinking the Value of Network Pruning](#). *arXiv preprint arXiv:1810.05270*, 2018.
- [32] Xiaolong Ma, Sheng Lin, Shaokai Ye, Zhezhi He, Linfeng Zhang, Geng Yuan, Sia Huat Tan, Zhengang Li, Deliang Fan, Xuehai Qian, Xue Lin, Kaisheng Ma, and Yanzhi Wang. [Non-Structured DNN Weight Pruning – Is It Beneficial in Any Platform?](#) *IEEE Transactions on Neural Networks and Learning Systems*, 2022.
- [33] Asit Mishra, Jorge Albericio Latorre, Jeff Pool, Darko Stosic, Dusan Stosic, Ganesh Venkatesh, Chong Yu, and Paulius Micikevicius. [Accelerating Sparse Deep Neural Networks](#). *arXiv preprint arXiv:2104.08378*, 2021.
- [34] Decebal Constantin Mocanu, Elena Mocanu, Peter Stone, Phuong H Nguyen, Madeleine Gibescu, and Antonio Liotta. [Scalable Training of Artificial Neural Networks with Adaptive Sparse Connectivity Inspired by Network Science](#). *Nature communications*, 2018.
- [35] Pavlo Molchanov, Arun Mallya, Stephen Tyree, Iuri Frosio, and Jan Kautz. [Importance Estimation for Neural Network Pruning](#). In *CVPR*, 2019.
- [36] Pavlo Molchanov, Stephen Tyree, Tero Karras, Timo Aila, and Jan Kautz. [Pruning Convolutional Neural Networks for Resource Efficient Inference](#). *arXiv preprint arXiv:1611.06440*, 2016.
- [37] Sharan Narang, Erich Elsen, Gregory Diamos, and Shubho Sengupta. [Exploring Sparsity in Recurrent Neural Networks](#). *arXiv preprint arXiv:1704.05119*, 2017.
- [38] Sharan Narang, Eric Undersander, and Gregory Diamos. [Block-sparse Recurrent Neural Networks](#). *arXiv preprint arXiv:1711.02782*, 2017.
- [39] Mi Sun Park, Xiaofan Xu, and Cormac Brick. [SQuantizer: Simultaneous Learning for both Sparse and Low-precision Neural Networks](#). *arXiv preprint arXiv:1812.08301*, 2018.
- [40] Jeff Pool and Chong Yu. [Channel Permutations for N:M Sparsity](#). *NeurIPS*, 2021.
- [41] Eric Qin, Geonhwa Jeong, William Won, Sheng-Chun Kao, Hyoukjun Kwon, Sudarshan Srinivasan, Dipankar Das, Gordon E Moon, Sivasankaran Rajamanickam, and Tushar Krishna. [Extending Sparse Tensor Accelerators to Support Multiple Compression Formats](#). In *IPDPS*, 2021.
- [42] 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](#). *arXiv preprint arXiv:1910.10683*, 2019.
- [43] Alex Renda, Jonathan Frankle, and Michael Carbin. [Comparing Rewinding and Fine-tuning in Neural Network Pruning](#). *arXiv preprint arXiv:2003.02389*, 2020.
- [44] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. [DistilBERT, a Distilled version of BERT: Smaller, Faster, Cheaper and Lighter](#). *arXiv preprint arXiv:1910.01108*, 2019.

- [45] Sheng Shen, Zhen Dong, Jiayu Ye, Linjian Ma, Zhewei Yao, Amir Gholami, Michael W Mahoney, and Kurt Keutzer. [Q-BERT: Hessian based Ultra Low Precision Quantization of BERT](#). In *AAAI*, 2020.
- [46] Zhiqing Sun, Hongkun Yu, Xiaodan Song, Renjie Liu, Yiming Yang, and Denny Zhou. [MobileBERT: A Compact Task-agnostic BERT for Resource-Limited Devices](#). *arXiv preprint arXiv:2004.02984*, 2020.
- [47] Zhanhong Tan, Jiebo Song, Xiaolong Ma, Sia-Huat Tan, Hongyang Chen, Yuanqing Miao, Yifu Wu, Shaokai Ye, Yanzhi Wang, Dehui Li, and Kaisheng Ma. [PCNN: Pattern-based Fine-grained Regular Pruning Towards Optimizing CNN Accelerators](#). In *DAC*, 2020.
- [48] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. [Attention is All you Need](#). In *NeurIPS*, 2017.
- [49] Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. [MiniLM: Deep Self-Attention Distillation for Task-Agnostic Compression of Pre-Trained Transformers](#). *NeurIPS*, 2020.
- [50] Wei Wen, Chunpeng Wu, Yandan Wang, Yiran Chen, and Hai Li. [Learning Structured Sparsity in Deep Neural Networks](#). *NeurIPS*, 2016.
- [51] Mitchell Wortsman, Ali Farhadi, and Mohammad Rastegari. [Discovering Neural Wirings](#). *NeurIPS*, 2019.
- [52] Zhuliang Yao, Shijie Cao, Wencong Xiao, Chen Zhang, and Lanshun Nie. [Balanced Sparsity for Efficient DNN Inference on GPU](#). In *AAAI*, 2019.
- [53] Amir Yazdanbakhsh, Ahmed T Elthakeb, Prannoy Pilligundla, F Mireshghallah, and Hadi Esmaeilzadeh. [ReleQ: An Automatic Reinforcement Learning Approach for Deep Quantization of Neural Networks](#). *arXiv preprint arXiv:1811.01704*, 2018.
- [54] Seul-Ki Yeom, Philipp Seegerer, Sebastian Lapuschkin, Alexander Binder, Simon Wiedemann, Klaus-Robert Müller, and Wojciech Samek. [Pruning by Explaining: A Novel Criterion for Deep Neural Network Pruning](#). *Pattern Recognition*, 2021.
- [55] Ofir Zafrir, Guy Boudoukh, Peter Izsak, and Moshe Wasserblat. [Q8BERT: Quantized 8bit BERT](#). *arXiv preprint arXiv:1910.06188*, 2019.
- [56] 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](#). *arXiv preprint arXiv:2205.01068*, 2022.
- [57] Wei Zhang, Lu Hou, Yichun Yin, Lifeng Shang, Xiao Chen, Xin Jiang, and Qun Liu. [Ternary-BERT: Distillation-aware Ultra-low Bit BERT](#). *arXiv preprint arXiv:2009.12812*, 2020.
- [58] 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](#). *arXiv preprint arXiv:2102.04010*, 2021.
- [59] Maohua Zhu, Tao Zhang, Zhenyu Gu, and Yuan Xie. [Sparse Tensor Core: Algorithm and Hardware Co-design for Vector-wise Sparse Neural Networks on Modern GPUs](#). In *MICRO*, 2019.
- [60] Michael Zhu and Suyog Gupta. [To Prune, or not to Prune: Exploring the Efficacy of Pruning for Model Compression](#). *arXiv preprint arXiv:1710.01878*, 2017.