

SPARSIFICATION VIA COMPRESSED SENSING FOR AUTOMATIC SPEECH RECOGNITION

Kai Zhen^{1*}, Hieu Duy Nguyen², Feng-Ju Chang², Athanasios Mouchtaris², and Ariya Rastrow²

¹Indiana University Bloomington

²Alexa Machine Learning, Amazon, USA

zhenk@indiana.edu, {hieng, fengjc, mouchta, arastrow}@amazon.com

ABSTRACT

In order to achieve high accuracy for machine learning (ML) applications, it is essential to employ models with a large number of parameters. Certain applications, such as Automatic Speech Recognition (ASR), however, require real-time interactions with users, hence compelling the model to have as low latency as possible. Deploying large scale ML applications thus necessitates model quantization and compression, especially when running ML models on resource constrained devices. For example, by forcing some of the model weight values into zero, it is possible to apply zero-weight compression, which reduces both the model size and model reading time from the memory. In the literature, such methods are referred to as sparse pruning. The fundamental questions are when and which weights should be forced to zero, i.e. be pruned. In this work, we propose a compressed sensing based pruning (CSP) approach to effectively address those questions. By reformulating sparse pruning as a sparsity inducing and compression-error reduction dual problem, we introduce the classic compressed sensing process into the ML model training process. Using ASR task as an example, we show that CSP consistently outperforms existing approaches in the literature.

Index Terms— Model pruning, automatic speech recognition (ASR), sparsity, Recurrent Neural Network Transducer (RNN-T), compressed sensing.

1. INTRODUCTION

Automatic Speech Recognition (ASR) is an important component of a virtual assistant system. The main focus of ASR is to convert users’ voice command into transcription, based on which further processing will act upon. Recently, end-to-end (E2E) approaches have attracted much attention due to their ability of directly transducing audio frame features into sequence outputs [1]. Without explicitly imposing/injecting domain knowledge and manually tweaking intermediate components (such as lexicon model), building and maintaining E2E ASR system is much more efficient than a hybrid deep neural network (DNN)-Hidden Markov Model (HMM) model.

In order to provide the best user experience, an ASR system is required to achieve not only high accuracy but also small user-perceived latency. This motivates the trend of moving the processing from Cloud/remote servers to users’ device to reduce the latency further. More often than not, the hardware limitations impose strict constraints on the model complexity [2, 3]. Firstly, the hardware

may only support integer arithmetic operations for run-time inference, which necessitates model quantization. Secondly, a hardware often performs multiple tasks supported by different models in sequence. The hardware, with limited memory size, thus needs to move multiple models in and out of the processing units.

To compress the model for the hardware, two widely applied methods are (a) model selection/structured pruning, i.e. choosing a model structure with pruned layers/channels and small performance degradation [4, 5], and (b) zero-weight compression/sparse pruning, i.e. pruning small-value weights to zero [6, 7]. Model selection differs from sparse pruning in that it deletes entire channels or layers, showing a more efficient speedup during inference, yet with a more severe performance degradation [4, 5]. These two types of methods are usually complementary: after being structurally pruned, a model can also undergo further zero-weight compression to improve the inference speed. In this study, we focus on sparse pruning, targeting at lowering the memory storage and bandwidth requirement as it largely contributes to the latency for on-device ASR models.

A naïve approach for sparse pruning is to push the weight values smaller than a threshold to zero after training, which often leads to significant performance degradation [6]. To mitigate this problem, Tibshirani *et. al.* applied *LASSO* regularization to penalize large-value model weights [8]. The drawbacks are twofold: firstly, it does not exert an explicit specification of the target sparsity level; secondly, it is subject to the gradient vanishing issue when the model has more and more layers. Gradual pruning approach [6] resolves those concerns by defining a sparsity function that maps each training step to a corresponding intermediate sparsity level. During the model training, the pruning threshold is adjusted gradually according to the function to eventually reach the target sparsity level. However, gradual pruning assumes that pruning the values smaller than the threshold will lead to the least degradation, which is heuristic and sub-optimal. Consequently, gradual pruning techniques provide a guidance on “when to prune and by how much”, but a less satisfying answer for “which (weights) to prune”, thus leading to inefficiency.

In this work, we propose a compressed sensing (CS)-based pruning method, referred to as CSP subsequently, that is sparse-aware and addresses both “when to prune” and “which to prune”. CSP reformulates the feedforward operations in machine learning architectures, such as Long Short Term Memory (LSTM) or Fully-Connected (FC) cells, as a sensing procedure with the inputs and hidden states being random sensing matrices. Under that perspective, a sparsification process is to enhance the sparsity and reduce the compression error, due to pruning, simultaneously. Following [9], we adopt the ℓ_1 regularization to enforce sparsity and the ℓ_2 regularization to mitigate the compression loss, and reformulate the

*This work was conducted during Kai’s internship in Amazon Pittsburgh, PA, USA.

sparsification procedure as an optimization problem. We demonstrate the effectiveness of our method by compressing recurrent neural network transducer (RNN-T), one of the E2E ASR models. The RNN-T model is sparsified via a hybrid training mechanism in which CSP is conducted during the feedforward stage, along with the back propagation for the global optimization. Our proposed method constantly outperforms the state-of-the-art gradual pruning approach in terms of the word error rate (WER) under all settings. In particular, with a sparsity ratio of 50% where half of the weights are 0, CSP yields little to no performance degradation on LibriSpeech dataset.

The rest of the paper is structured as follows. In Sec. 2, we briefly review related pruning methods. Our CSP method is introduced in Sec. 3. Sec. 4 describes our experiment setup and results. Finally, we conclude our work with some remarks in Sec. 5.

2. RELATED WORK

One of the most straightforward approaches for sparse-aware training is applying ℓ_k regularization, where $k = 0, 1, 2$, etc. There has been a rich literature in comparing various forms of sparsity regularizers. Consider the model training: $\mathbb{W} \leftarrow \arg \min_{\mathbb{W}} \mathcal{L}_{\text{accuracy}}(\mathbb{W}) + \|\mathbb{W}\|_1$, where ℓ_1 norm is used on the regularization. The fundamental idea is to optimize the model prediction while penalizing large weight values. In DNN model compression, the regularization is usually implemented as an extra loss term for the training.

This training-aware sparsity regularization leads to promising pruning results especially for convolutional neural networks (CNN) [10] with residual learning techniques [11, 12], but may not apply well to models employing recurrent neural network (RNN) components such as LSTM. The error due to a global sparsity constraint ℓ_1 will be propagated to all time steps. Additionally, such drawback is much more severe for architectures, such as RNN-T, which contains feedback loop from one part of the model to the others.

Another well-known, state-of-the-art, pruning method for ML models is gradual-pruning [6]. This method does not resort to ℓ_1 regularization for sparsity, but dynamically updates the pruning threshold during model training, as is indicated by its name. To answer the question "when to prune", the authors defines a sparsity function parametrized by the target sparsity s_f at t_n step with an initial pruning step t_0 . Concretely, at training step t , the pruning threshold is adjusted to match the sparsity s_t calculated in Eq. 1. The main complication is to adjust the pruning procedure such that the model weights are relatively converged and the learning rate is still sufficiently large to reduce the pruning-induced loss.

$$s_t = s_f * \left(1 - \left(1 - \frac{t - t_0}{t_n - t_0} \right)^3 \right) \quad (1)$$

One concern is that finding an optimal setup for these hyperparameters can be hard without going through a rigorous ablation analysis. Furthermore, with gradual pruning, gradient-updating back-propagation is the only mechanism to limit the degradation. Most importantly, gradual pruning only addresses the question "when to prune" but not "which (weights) to prune". At each time step, the weights below the (gradually increased) threshold are to be pruned. This is based on the premise that the smaller the weight, the less important it is, which is heuristic and sub-optimal.

3. COMPRESSED SENSING BASED PRUNING

3.1. Adapting Compressed Sensing for Sparse Pruning

CS aims to compressing potentially redundant information into a sparse representation and reconstructing it efficiently [13, 14], which has facilitated a wide scope of engineering scenarios, such as medical resonance imaging (MRI) and radar imaging. For example, in MRI [15], high resolution scanned images are generated per millisecond or microsecond, leading to significant storage cost and transmission overhead. CS learns a sparse representation of each image with which, during the decoding time, CS can recover the reference image almost perfectly. Assume an image compression task with a reference image $\mathbf{x} \in \mathbb{R}^n$, where \mathbf{x} is usually decomposed into an orthogonal transformation basis ψ and the activation \mathbf{s} , as $\mathbf{x} = \psi * \mathbf{s}$. Given that \mathbf{s} satisfies \mathcal{K} -sparse property, CS is capable of locating those \mathcal{K} salient activation elements. Concretely, CS introduces a sensing matrix ϕ to project \mathbf{x} into \mathbf{y} , as $\mathbf{y} = \phi * \mathbf{x}$. In [9, 16], it has been proved that by optimizing the ℓ_2 loss in the sensing dimensionality while exerting the ℓ_1 norm regularizer to \mathbf{s} , a \mathcal{K} -sparse solution of \mathbf{s} , denoted as $\hat{\mathbf{s}}$, can be found in polynomial time. Consequently, $\psi * \hat{\mathbf{s}}$ can estimate the original image \mathbf{x} with high fidelity and relatively small latency.

In this work, we investigate the effectiveness of CS based pruning for ML models. We consider the ASR task, i.e. converting audio speech to transcriptions, using an RNN-T architecture. Due to the space constraint, we only describe the transformation of LSTM cell, which is a major building block in various E2E ASR models. It is straightforward to extend the transformations to other architectures/layers like the fully-connected (FC) network and CNN.

Consider a vanilla LSTM cell: the element-wise multiplication between the input at time step t , $\mathbf{x}^{(t)}$, and kernels is given in Eq. 2, while that of hidden states from the previous step $\mathbf{h}^{(t-1)}$ and recurrent kernels is in Eq. 3. Here, $\mathbf{W}_f, \mathbf{W}_c, \mathbf{W}_i$, and \mathbf{W}_o (correspondingly $\mathbf{U}_f, \mathbf{U}_c, \mathbf{U}_i$, and \mathbf{U}_o) denote the kernels (correspondingly recurrent kernels) weights of the cell (c), the input gate (i), output gate (o), and forget gate (f), respectively. All gating mechanisms to update the cell and hidden states are encapsulated in Eq. 4, where $\mathbf{C}^{(t)}$ is the cell state vector at time t and \mathcal{G} denotes the transformation of $\mathbf{z}_x, \mathbf{z}_h, \mathbf{C}^{(t-1)}$ into $\mathbf{C}^{(t)}, \mathbf{h}^{(t)}$. The bias terms are omitted for ease of presentation.

$$\mathbf{z}_x = [\mathbf{W}_f, \mathbf{W}_c, \mathbf{W}_i, \mathbf{W}_o] \odot [\mathbf{x}^{(t)}, \mathbf{x}^{(t)}, \mathbf{x}^{(t)}, \mathbf{x}^{(t)}] \quad (2)$$

$$\mathbf{z}_h = [\mathbf{U}_f, \mathbf{U}_c, \mathbf{U}_i, \mathbf{U}_o] \odot [\mathbf{h}^{(t-1)}, \mathbf{h}^{(t-1)}, \mathbf{h}^{(t-1)}, \mathbf{h}^{(t-1)}] \quad (3)$$

$$[\mathbf{C}^{(t)}, \mathbf{h}^{(t)}] = \mathcal{G}(\mathbf{z}_x, \mathbf{z}_h, \mathbf{C}^{(t-1)}), \quad (4)$$

To prune all kernel weights (denoted as $\mathbb{W} = [\mathbf{W}_f, \mathbf{W}_c, \mathbf{W}_i, \mathbf{W}_o]$) in LSTM cells, we adapt and reformulate a CS-like pruning procedure by adopting ℓ_1 regularization for sparsity-inducing and ℓ_2 regularization for compression-error reduction, as outlined in Fig. 1. The procedure starts by a midway inference to collect the activation inputs $[\mathbf{z}_x, \mathbf{z}_h]$ as in Eq. 2 and Eq. 3. When those kernels are sparsified, the activation inputs will be consequently updated as $[\mathbf{z}'_x, \mathbf{z}'_h]$ (see Fig. 1). The goal is to sparsify and prune the kernels while preserving the value of $[\mathbf{z}_x, \mathbf{z}_h]$ to minimize the pruning-induced loss. To that end, the ℓ_1 regularizer is applied to the input kernels while the ℓ_2 regularizer controls the reconstruction loss on \mathbf{z}_x . Hence, our CS solver is embedded in a local optimizer triggered periodically by feedforward steps in a stochastic manner. As illustrated in the restricted isometry property (RIP), the sensing matrix ϕ is expected to be random for an accurate signal reconstruction [16]. The proposed CS solver satisfies the RIP

Table 1: Model performance under various sparsity levels for far-field (left) and LibriSpeech (right) datasets.

		M-I, 38.7M	M-II, 60.0M	M-III, 34.0M			M-IV, 37.1M				
Sparsity (%)	Methods	Rel. Dgrd (%)		Sparsity (%)	Methods	WER	Abs. Dgrd	Rel. Dgrd	WER	Abs. Dgrd	Rel. Dgrd
0	–	–	–	0	–	7.27	–	–	9.58	–	–
25	A	1.83	0.80	50	A	17.45	10.18	140.03%	37.24	27.66	288.73%
	B	1.12	0.63		B	7.34	0.07	0.96%	10.35	0.77	8.04%
	Proposed	0.96	0.55		Proposed	7.26	-0.01	-0.14%	10.06	0.48	5.01%
50	A	23.40	20.48	75	A	99.76	92.49	1272.21%	95.06	85.48	892.28%
	B	8.77	7.14		B	8.13	0.86	11.83%	10.43	0.85	8.87%
	Proposed	6.32	5.05		Proposed	8.10	0.83	11.42%	10.14	0.56	5.84%

We compare our proposed CSP with two other pruning methods: naïve pruning and gradual pruning. In the naïve pruning approach, termed as method-A, the smallest weights are pruned post-training to reach the target sparsity level. Note that for a significantly over-parameterized model, achieving a certain sparsity level with little degradation may not be challenging even with method-A, since most of the weights are not effectively involved in the optimization. Method-A, although not being spare-aware during training, thus helps us probe the level of robustness for an RNN-T topology under various sparsity levels. The gradual pruning approach, denoted as method-B, is derived from [6]. As illustrated in Sec. 3, the pruning threshold is calibrated kernel-wise for a fair comparison.

The learning rate in all experiments is specified via a warm-hold-decay scheduler: the initial learning rate, $1e-7$, is raised up to $5e-4$ at 3K-th step, and is being held till 75K-th step, which then follows a decay stage with which it is reduced to $1e-5$ at 200K-th step. The pruning starts at 100K-th step and gradually increases the pruning threshold to reach the target sparsity level at 150K-th step. The intuition, similar to [6], is to neither apply pruning too early such that the weights are reasonably distributed, nor too late to allow recovering from the sparsification-induced degradation. All models are trained with 10% dropout. The sensing coefficient λ initialized as 0.1. Since λ is adjustable according to Eq. 6 during training, the results are not predominantly contingent on its initial value.

4.2. Experimental Results

Consider the performance of M-I and M-II, which do not have a joint network and are trained on the far-field dataset. It is observed in Table 1 that the models are relatively robust at the sparsity level of 25%. At 50%, the degradation becomes noticeable for all 3 methods. The hard pruning approach does not yield a desirable performance, while our proposed CSP method gives the lowest WER relative degradation in these experimental settings. As expected, the results also indicate a higher relative degradation when the model size decreases.

In Table 1, we also report absolute WERs from models trained on the LibriSpeech train dataset and decoded on the LibriSpeech test-clean dataset. Not surprisingly, the models trained with method-A also suffer significant degradation at the sparsity level 50%. Again, it is observed that the proposed CSP method consistently outperforms all other approaches. Comparing to M-IV, M-III indicates a higher robustness to pruning likely due to the J-N. At 75% sparsity, all methods experience substantial performance degradation, especially for models with J-N. One reason is that the additional layer actually exacerbates the error, thus leading to higher relative degradation. However, it is worth noting that RNN-T models with J-N still outperforms the counterparts without it by a large margin.

To understand the effect of CSP, we investigate how the model weights are being redistributed when CSP is applied. Fig. 3 shows the weight distribution when CSP is applied to M-III. It is observed

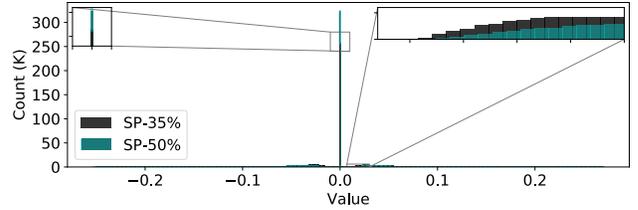


Fig. 3: Model weight histogram when the sparsity (SP) is up from 35% to 50%: the threshold is barely increased with newly pruned weights selected via the CSP method as shown in zoomed-in insets.

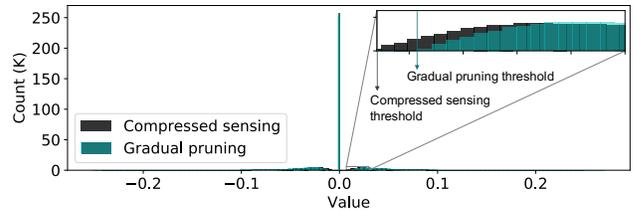


Fig. 4: Pruning threshold comparison at 50% sparsity: CSP can redistribute weights to approach the same sparsity level with a smaller pruning threshold than the gradual pruning method.

that to achieve a higher sparsity ratio from 35% to 50% for CSP, the pruning threshold does not move by a large margin. Instead, a set of the weights are driven towards 0, and consequently being pruned. In particular, from the top-right zoomed-in figure, we can see that most of the additional weights pruned at 50% in CSP are not those closest to the threshold at 35% sparsity level. In contrast, the gradual pruning approach will significantly increase the pruning threshold to accommodate the higher sparsity level, and then simply prune the weights with the smallest values (Fig. 4). Rather than just a hard reset of the threshold, CSP instead determines “which (weights) to prune” via a joint optimization on sparsity and reconstruction regularizers (see Eq.5), thus leading to much smaller WER degradation.

5. CONCLUSIONS

We propose a novel pruning approach for machine learning model compression based on compressed sensing, termed as CSP. Compared to existing sparsification methods which focus only on “when to prune”, CSP further addresses the question “which (weights) to prune” by considering both sparsity inducing and compression-error reduction mechanisms. We validate the effectiveness of CSP via the speech recognition task with RNN-T model. CSP achieves superior results compared to other sparsification approaches. The proposed method can be straightforwardly incorporated into other ML models and/or compression methods to further reduce model complexity.

6. REFERENCES

- [1] R. Prabhavalkar, K. Rao, T. N. Sainath, B. Li, L. Johnson, and N. Jaitly, "A comparison of sequence-to-sequence models for speech recognition," in *Proceedings of the Annual Conference of the International Speech Communication Association (INTERSPEECH)*, 2017, pp. 939–943.
- [2] S. Han, H. Z. Mao, and W. J. Dally, "Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding," in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2016.
- [3] S. Punjabi, H. Arisikere, Z. Raeesy, C. Chandak, N. Bhave, A. Bansal, M. Müller, S. Murillo, A. Rastrow, S. Garimella, et al., "Streaming end-to-end bilingual asr systems with joint language identification," *arXiv preprint arXiv:2007.03900*, 2020.
- [4] S. J. Cao, C. Zhang, Z. L. Yao, W. C. Xiao, L. S. Nie, D. C. Zhan, Y. X. Liu, M. Wu, and L. T. Zhang, "Efficient and effective sparse LSTM on FPGA with bank-balanced sparsity," in *Proceedings of the ACM/SIGDA International Symposium on Field-Programmable Gate Arrays (FPGA)*, 2019, pp. 63–72.
- [5] S. R. Wang, P. Lin, R. H. Hu, H. Wang, J. He, Q. J. Huang, and S. Chang, "Acceleration of LSTM with structured pruning method on FPGA," *IEEE Access*, vol. 7, pp. 62930–62937, 2019.
- [6] M. Zhu and S. Gupta, "To prune, or not to prune: exploring the efficacy of pruning for model compression," in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2018.
- [7] S. Narang, E. Elsen, G. Diamos, and S. Sengupta, "Exploring sparsity in recurrent neural networks," in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2017.
- [8] R. Tibshirani, M. Saunders, S. Rosset, J. Zhu, and K. Knight, "Sparsity and smoothness via the fused lasso," *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, vol. 67, no. 1, pp. 91–108, 2005.
- [9] D. L. Donoho, "Compressed sensing," *IEEE Transactions on Information Theory*, vol. 52, no. 4, pp. 1289–1306, 2006.
- [10] C. Louizos, M. Welling, and P. D. Kingma, "Learning sparse neural networks through l0 regularization," in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2018, pp. 1389–1397.
- [11] K. M. He, X. Y. Zhang, S. Q. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 770–778.
- [12] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 4700–4708.
- [13] X. Yuan and R. Haimi-Cohen, "Image compression based on compressive sensing: End-to-end comparison with JPEG," *IEEE Transactions on Multimedia*, 2020.
- [14] M. Qiao, Z. Y. Meng, J. W. Ma, and X. Yuan, "Deep learning for video compressive sensing," *APL Photonics*, vol. 5, no. 3, 2020.
- [15] M. Lustig, D. L. Donoho, J. M. Santos, and J. M. Pauly, "Compressed sensing MRI," *IEEE Signal Processing Magazine*, vol. 25, no. 2, pp. 72–82, 2008.
- [16] E. J. Candès, J. Romberg, and T. Tao, "Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information," *IEEE Transactions on Information Theory*, vol. 52, no. 2, pp. 489–509, 2006.
- [17] W. Chan, N. Jaitly, Q. Le, and O. Vinyals, "Listen, attend and spell: A neural network for large vocabulary conversational speech recognition," in *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2016, pp. 4960–4964.
- [18] A. Graves, S. Fernández, F. Gomez, and J. Schmidhuber, "Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks," in *Proceedings of the International Conference on Machine Learning (ICML)*, 2006, pp. 369–376.
- [19] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, "Librispeech: an asr corpus based on public domain audio books," in *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2015, pp. 5206–5210.