

Transfer Learning for Structured Pruning under Limited Task Data

Lucio Dery
CMU

ldery@andrew.cmu.edu

David Grangier
Apple

lastname@apple.com

Awni Hannun
Apple

awni@apple.com

Abstract

Large, pre-trained models are problematic to use in resource constrained applications. Fortunately, task-aware structured pruning methods offer a solution. These approaches reduce model size by dropping structural units like layers and attention heads in a manner that takes into account the end-task. However, these pruning algorithms require more task-specific data than is typically available. We propose a framework which combines structured pruning with transfer learning to reduce the need for task-specific data. Our empirical results answer questions such as: How should the two tasks be coupled? What parameters should be transferred? And, when during training should transfer learning be introduced? Leveraging these insights, we demonstrate that our framework results in pruned models with improved generalization over strong baselines.

1 Introduction

Large pre-trained language models have been successfully applied to a wide variety of application scenarios (Bommasani et al., 2021; Anil et al., 2023). However, not all applications can justify the cost of running such large models. E.g. an interactive, offline spellchecker for a phone has strong memory limits compared to a server-side chat model (Dettmers et al., 2022). Even server-side, the benefit/cost of large models depends on the application. This situation motivates research into structured model pruning algorithms.

Structured pruning algorithms generate smaller, faster and yet reasonably accurate sub-models from large pre-trained ones by removing components (beyond individual parameters) like convolutional channels, attention heads and whole layers. Several works over the years (Wang et al., 2019; Sanh et al., 2020; Xia et al., 2022) have been proposed to perform task-specific structured pruning. Unfortunately, to the best of our knowledge, all existing algorithms have been

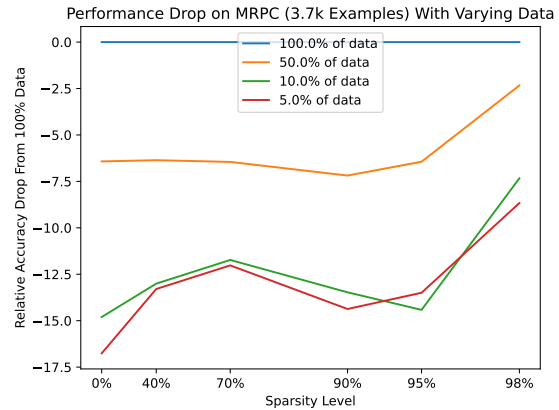


Figure 1: Accuracy degradation of CoFi (Xia et al., 2022) vs training data sizes. Sparsity level refers to the fraction of removed weights (excluding embeddings). Accuracy at 50% data is stable across sparsity levels (except for 98% sparsity) while more data-limited regimes (10%–5%) exhibit stronger sensitivity to the sparsity level.

developed without consideration for the amount of training data available for the target task. Thus, as Figure 1 shows that, even state-of-the-art methods like CoFi (Xia et al., 2022), do not gracefully handle scenarios with limited training data. We argue that the data-limited structured pruning setting is important since limited compute for inference and data scarcity for training tend to co-occur often in practice (Ahia et al., 2021). A popular remedy to the limited data problem at fixed model size, is to leverage transfer learning (Caruana, 1997; Erhan et al., 2010; Dery et al., 2022) by introducing external data or extra tasks. In this work, we investigate transfer learning based remedies for structured pruning under limited data. Structured pruning algorithms need to jointly learn both model weights and structural variables (which layers, attention heads, etc. to prune) for the final size-reduced model (Wang et al., 2019; Xia et al., 2022). This added complexity makes deploying transfer learning in the structured pruning setting

non-trivial and raises several questions. Do we only perform transfer learning for model weights or do we include structural variables too? How do we learn structural variables for the target task in a way that benefits from the presence of a transfer task? When is it best to introduce transfer learning so as to produce the most accurate pruned target model?

This work aims to provide answers to the questions above. We propose a simple modification to existing structured pruning algorithms to allow for effective transfer of both structural variables and model parameters. Overall, our analyses allow us to provide prescriptions to researchers about what, how and when to transfer during structured pruning. Our effort results in significant improvements in generalization performance even at compression ratios as high as $50\times$.

2 Background

Unstructured Pruning approaches sparsify models by zeroing out individual components of weight matrices (Frankle and Carbin, 2018; Sanh et al., 2020). The resulting sparse matrices reduce the memory overhead of the model but run-time gains cannot be realized unless on specialized hardware (Liu et al., 2018; Ma et al., 2021). Over the years, many criteria for choosing which parameters to remove have been explored. Some approaches like magnitude pruning (Han et al., 2015) and Wanda (Sun et al., 2023) prune parameters based on either their magnitudes or the magnitude of their product with previous layer activations respectively. Other approaches like (Frankle and Carbin, 2018; Sanh et al., 2020) use information about about much parameters have changed since initialization whilst others learn unstructured masks based using gradient descent (Ramanujan et al., 2020). Ahia et al. (2021) introduce the term *the low-resource double-bind* for the challenge of compressing models in data limited regimes. Unlike us, they study magnitude pruning, which as mentioned, does not ordinarily lead to run-time gains. They also do not propose a remedy for the limited-data problem, which we do in this paper.

Structured Pruning algorithms remove whole components from pre-trained models such as attention heads (Michel et al., 2019; Voita et al., 2019), whole layers (Fan et al., 2019) or intermediate dimensions of fully connected layers (Wang et al., 2019) in order to produce

faster, memory efficient sub-models without overly sacrificing downstream accuracy. Unlike unstructured pruning, there is no need for specialized hardware in order to realize the run-time speedups from compression. These approaches require optimizing over structural variables (to decide which model components to prune) and model weights (to adapt the final model to the disruption that results from removing whole components). Joint optimizations like these mean more variables to learn, resulting in the need for more end-task data points. To the best of our knowledge, we are the first consider the challenge of structured pruning under limited data.

Other model compression approaches

Quantization methods (Polino et al., 2018; Dettmers et al., 2022) reduce model size by reducing the number of bits required to represent each weight. These methods are generally complementary to pruning approaches but only achieve maximum size reductions on the order of $2-4\times$ before substantial model performance degradation. We are interested in achieving extreme compressions to the order of $50\times$ reduction without significant loss in performance. Distilling directly to a target task has been shown to be a data-hungry process (Jiao et al., 2019), often requiring a general distillation step (on abundant external data) to be able to achieve competitive performance with approaches modern structured pruning methods like CoFi (Xia et al., 2022).

Multitask Transfer Learning (Caruana, 1997)

is a common recipe for improving a models average performance on a desired end-task. When the end-task is data-limited, auxiliary tasks can be multi-tasked with the end-task (Dery et al., 2021a,b) to serve as proxy data. Previous work at the intersection of pruning and multitasking have only studied how to prune multi-task models (Garg et al., 2023; Yang et al., 2023). Unlike these, our starting point is not a multitask model but a generalist pre-trained model like BERT (Devlin et al., 2018). Our work is interested in using multitasking in as much as it improves generalization of the pruned model with respect to the data-starved end-task only.

3 Methodology

The goal of this paper is to improve the generalization of pruned models when the end-

task is data-limited without sacrificing memory and run-time gains. We assume that we are given a structured pruning algorithm that jointly learns structural/masking variables (which we denote as $\{\mathbf{z}_{\text{target}}^k\}$) and their corresponding parameters $\{\theta_{\text{target}}^k\}$ for the target end-task. k indexes the set of K structural variables being explored. We are primarily concerned with how to incorporate a transfer task by learning $[\{\mathbf{z}_{\text{transfer}}^k\}, \{\theta_{\text{transfer}}^k\}]$ such that we enjoy improved generalization with the target task’s final model. Since our focus is on the limited-data problem, we care less about a specific structured pruning algorithm and more about how to adapt any appropriate algorithm in the data starved setting. We therefore focus on building on top of a state-of-the-art structured pruning algorithm, CoFi (Xia et al., 2022) which we take as a representative algorithm. Whilst we describe CoFi below to provide sufficient background, for the rest of the paper, we will abstract away the details of the pruning algorithm and focus on the specifics of adapting transfer learning to this setting.

3.1 CoFi

CoFi (**C**oarse- and **F**ine-grained Pruning) is a mixed resolution structured pruning algorithm. Previous algorithms to prune transformer models (Vaswani et al., 2017) have focused on removing high level units like whole layers (Fan et al., 2019) or finer grained modules like attention heads (Voita et al., 2019) and dimensions of fully connected layers (Wang et al., 2019) but not both types. CoFi introduces variables that account for pruning at multiple levels of granularity.

Coarse Grain: Each transformer layer consists of a multi-headed attention component that feeds into a fully connected two-layer non-linear perceptron (Vaswani et al., 2017). CoFi introduces variables sets $\{z_{\text{MHA}}^i\}_{i \in [N]}$ and $\{z_{\text{FFN}}^i\}_{i \in [N]}$ for each of the model N layers. z_{MHA}^i denotes the probability that the **whole** attention component of the i th layer is removed whilst z_{FFN}^i is similarly defined for the fully connected component of the specified layer. CoFi also removes whole columns of the residual stream: $z^\ell \in \mathbb{R}^d \rightarrow \hat{z}^\ell \in \mathbb{R}^{\hat{d}} \quad \forall \ell \in [N]$. For a BERT model, $d = 768$ is typically reduced to $d \approx 750$. (Xia et al., 2022) find that though relatively few columns are dropped, including columns as structural variables is important for producing performant compressed models.

Fine Grain: Given a particular layer i , CoFi prunes subsets of the attention heads available. The variables $\{z_{j,\text{head}}^i\}_{j \in [n_h]}$ represent the j th attention head in the i layer which has n_h total attention heads. A similar set of variables is defined for the fully connected units within a layer: $\{z_{j,\text{fc}}^i\}_{j \in [n_f]}$ where the i th fully connected layer has n_f units.

For the j th attention head of the i th layer, the likelihood that this head is left unpruned is proportional to $z_{\text{MHA}}^i \cdot z_{j,\text{head}}^i$. This allows the algorithm to make coupled fine and coarse grained decisions that lead to improved results. We collectively represent $\{\mathbf{z}\}$ as the set of all structural variables that are learned by CoFi. For a model with parameters θ , $\{\mathbf{z}\}$ are learned by applying the reparameterisation trick on the hard concrete distribution (Louizos et al., 2018) and minimizing a joint loss wrt $\{\mathbf{z}, \theta\}$ that includes

1. distance from target size. CoFi follows Wang et al. (2019) and adds a lagrangian term that penalizes deviations from the target sparsity.
2. target task loss. Practitioners ultimately want a pruned model that generalizes well on their end-task. CoFi jointly optimizes the target task loss along with the pruning objective in order to produce performant pruned models.
3. a distillation objective on the original large model. Following Sanh et al. (2020), CoFi jointly performs distillation and structured pruning by introducing a layer-wise distillation objective.

With these high level details in mind, we proceed to present our simple, transfer learning based modification to CoFi that leads to improved results in data-limited settings.

3.2 Transfer Learning for Structured Pruning under limited data

Given a target task \mathbf{T} with limited training data, we want to improve the final model generated by CoFi through leveraging additional training data from an auxiliary task \mathbf{A} . Let $\{\mathbf{z}_{\mathbf{T}}, \theta\}$ be the initial set of all structural variables and model parameters for the target task and $\{\hat{\mathbf{z}}_{\mathbf{T}}, \hat{\theta}\}^\gamma$ be final output of CoFi at a chosen sparsity level γ . $\hat{\mathbf{z}}$ are binary variables $\hat{z}_i \in \{0, 1\}$ which indicate whether component i is dropped/masked out (0) or is retained (1). We would like a procedure that leverages the auxiliary task (with its own set of variables $\{\mathbf{z}_{\mathbf{A}}\}$ such

that the generalization performance of the pruned model when using data from **A** and **T** jointly improves upon using only data from **T**.

There are several design questions that arise in this setting when thinking about how to effectively utilize **A**. In the following sections, we discuss some of these pertinent questions and propose some reasonable choices which we will later experimentally validate.

3.2.1 What criteria do we use to select the auxiliary task **A** ?

The choice of auxiliary task, **A**, is an important design decision that must be considered carefully. A poor choice could result in poor generalization performance (with respect to **T**) of the pruned model instead of being helpful. To this end, inspired by existing literature, we propose two criteria for evaluating what auxiliary task to leverage:

(1) resourcedness: Previous work on transfer learning for learning model parameters has demonstrated the benefits of leveraging large pools of data (which may possibly be unrelated to the eventual end-task) for pre-training (Anil et al., 2023) or multi-tasking (Dery et al., 2021b). We therefore have a strong prior that using data-rich auxiliary tasks might be helpful for also learning structural parameters even if they are unrelated to the end-task.

(2) task-similarity Both theoretical (Baxter, 2000; Maurer et al., 2016; Dery et al., 2022) and empirical works (Gururangan et al., 2020; Dery et al., 2022) have shown that transfer learning works best when the auxiliary task is similar or related to the end-task. As a proxy for similarity, we consider auxiliary tasks that are from the same *domain* as the end-task.

3.2.2 When should we introduce **A** ?

Structured pruning approaches like CoFi usually perform a two stage process. In the first stage, they generate a pruned model at the desired sparsity level; this involves learning both $\{\hat{\mathbf{z}}, \hat{\theta}\}_{\mathbf{T}}^{\gamma}$. In the second stage, pruned model is then fine-tuned on the end-task by updating only $\hat{\theta}_{\mathbf{T}}^{\gamma}$ keeping $\hat{\mathbf{z}}_{\mathbf{T}}^{\gamma}$ fixed. The auxiliary task can be introduced in either or both of these stages. We explore following choices:

Prune(A) \rightarrow FT(T): We do structural pruning to learn both the weights and structure for a small

model using *only* the transfer task, **A**: $\{\hat{\mathbf{z}}, \hat{\theta}\}_{\mathbf{A}}^{\gamma}$. We then fine-tune (FT) the pruned model on the target task (**T**) only to obtain $\hat{\theta}_{\mathbf{T}}^{\gamma}$. Here, the target task is used only in the final fine-tuning stage and is not involved in learning the pruned model structure.

Prune(T) \rightarrow FT(A, T): We learn both the weights and structure for a small model using the target task: $\{\hat{\mathbf{z}}, \hat{\theta}\}_{\mathbf{T}}^{\gamma}$. We share the pruned model parameters $\hat{\theta}^{\gamma}$ and fine-tune on both (**T**) and (**A**). Here, the auxiliary task is used only in the final fine-tuning stage and is not involved in learning the pruned model structure.

Prune(A, T) \rightarrow FT(T): We learn both the weights and structure for a small model using both the transfer and end-task: $\{\hat{\mathbf{z}}_{\mathbf{T}}, \hat{\mathbf{z}}_{\mathbf{A}}, \hat{\theta}\}^{\gamma}$ are learned jointly (we will explore how in Section 3.2.3). We then fine-tune (FT) the pruned model weights on the target task (**T**) only.

Prune(A, T) \rightarrow FT(T, A): We learn both the weights and structure for a small model using both the transfer and end-task: $\{\hat{\mathbf{z}}_{\mathbf{T}}, \hat{\mathbf{z}}_{\mathbf{A}}, \hat{\theta}\}^{\gamma}$ are learned jointly (we will explore how in Section 3.2.3). We then fine-tune (FT) the pruned model weights on **both** the target and auxiliary tasks.

3.2.3 How do we incorporate **A** when optimizing for $\{\hat{\mathbf{z}}_{\mathbf{T}}, \hat{\theta}\}^{\gamma}$

When using the auxiliary task directly during pruning, there is the question of what the best way to jointly optimize $\{\hat{\mathbf{z}}_{\mathbf{T}}, \hat{\mathbf{z}}_{\mathbf{A}}, \hat{\theta}\}^{\gamma}$ such that we achieve improved generalization for the final pruned model with respect to **T**. Note that we are assuming that the model parameter weights θ are shared between the two tasks but the structural variables are separate. This is because there are many more model parameters than structural variables $\|\theta\|_0 \gg \|\mathbf{z}_{\mathbf{T}}\|_0 + \|\mathbf{z}_{\mathbf{A}}\|_0$. And so introducing separate model weights for the auxiliary task presents a much more significant modelling overhead than introducing new structural variables.

We can explore different strategies for sharing variables across the two tasks such that the **T** benefits from **A**.

Single mask multi-task learns a single set of structural $\{\mathbf{z}\}$ and model $\{\theta\}$ parameters that are shared between both tasks. This choice tightly couples the two tasks. Whilst this allows maximal sharing of information between the target and

transfer task, poor choices of transfer tasks could cause this to perform worse than no transfer at all.

Multi-mask multi-task learns distinct structural parameters $\{z\}_T$ and $\{z\}_A$ for each task but a single set of model parameters $\{\theta\}$ is shared between both tasks. There is no transfer of structural information and only the shared model parameters provide a coupling of the two tasks.

Our δ -Formulation aims to leverage strength from both alternatives. We propose this method where both tasks share a base set of structural variables $\{z\}_{\text{base}}$ but also have task specific addends such that: $\{z\}_T = \{z_{\text{base}} + \delta_T\}$ and $\{z\}_A = \{z_{\text{base}} + \delta_A\}$. We regularize δ_* to encourage sharing between tasks via z_{base} whilst maintaining flexibility for task-specific modelling.

4 Experimental Setup

Our experimental framework is introduced to investigate the questions posed in the previous section.

Datasets We consider 3 pairs of tasks. One pair of classification tasks are from the computer science domain tasks – SCIIE (Luan et al., 2018) and ACL-ARC (Jurgens et al., 2018) with 3.2k and 3.7k training samples respectively. The second pair of tasks are biomedical domain tasks - RCT (Dernoncourt and Lee, 2017) (we artificially create a low-resource version of this task with 10k training samples) and CHEMPROT (Kringelum et al., 2016) which has 4.2k training samples. We use GLUE (Wang et al., 2018) tasks for our last pair: STSB and MRPC are sentence similarity and paraphrase detection tasks with 7k and 3.7k train examples respectively. For the GLUE tasks, we follow previous work (Jiao et al., 2019; Wang et al., 2019; Xia et al., 2022) and report results on the validation set. For Non-GLUE tasks, we report test set results. Please see Appendix A for more details about the tasks we investigate.

Model Details Since we use CoFi (Xia et al., 2022) as our representative structured pruning algorithm, we use the same model configuration. We use the BERT_{base} (Devlin et al., 2018) which has $\sim 110M$ parameters. We explore pruned model sparsities in the set $\{40\%, 70\%, 90\%, 95\%, 98\%\}$. $\gamma\%$ sparsity means that the model has been reduced to $(100-\gamma)\% \times 110M$ parameters. Similar to (Sanh et al., 2020) we also freeze the model embedding weights. See Appendix B for details about training

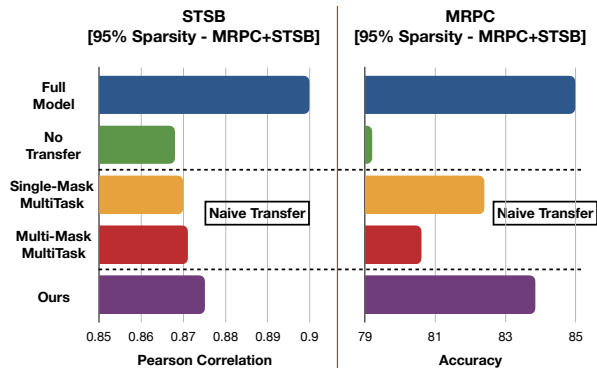


Figure 2: STSB and MRPC performance at 95% sparsity. Our proposed δ -Formulation outperforms all other methods on both tasks.

as well as hyper-parameter values.

Training details We mostly follow the training recipe from CoFi with a few minor changes. CoFi assumes that *one starts pruning after finetuning the full parent model on the target task* and so introduces a distillation loss as part of the pruning objective. In our case, we start directly from the pre-trained model without first fine-tuning on the target task. This is because of the risk of overfitting due to the smaller target task size. Due to this, we find that the distillation based losses from the original CoFi paper are unnecessary and we did not see significant performance differences with or without them. When multitasking, we explore a small set of weighting hyper-parameters $\{(1.0, 1.0), (1.0, 2.0), (2.0, 1.0)\}$ for any losses relating to the target and auxiliary tasks respectively. Table 5 has details of the hyper-parameters we cross validate against for all our experiments.

5 Empirical Recommendations for practitioners

In Section 3.2, we posed different design questions around how to perform transfer learning for structured pruning under limited data and presented different options for resolving said questions. In this section, we proceed to perform a sequence of experiments to validate which choices lead to superior end-task generalization after pruning, so we can make principled recommendations to practitioners.

5.1 How should you transfer?

In Section 3.2.3, we introduced various approaches for coupling the auxiliary task with the target task

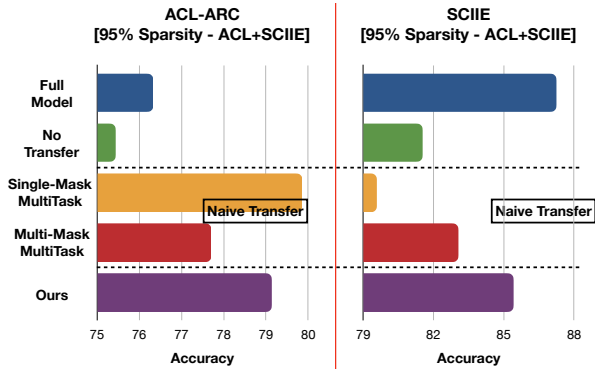


Figure 3: SCIIE and ACL_ARC performance at 95% sparsity. Our δ -Formulation produces the best average performance across the two tasks when either is used as the auxiliary task for the other.

during structured pruning. Figures 2, 3 and 4, show experimental results after implementing various options with different pairs of datasets. Across all dataset pairs, our δ -Formulation produces the best performance when averaged across the task pair.

For the SCIIE task, tightly coupling its structural variable with those of ACL_ARC (as an auxiliary task) under the single-mask multi-task approach can negatively impact performance compared to not doing transfer learning at all (Figure 3). Our δ -Formulation ensures that SCIIE actually benefits introducing transfer learning by outperforming the multi-mask multi-task approach that fully decouples the structural variables. For the ACL_ARC task, our formulation recovers close to the best performance (single mask multi-task). Note that in principle, our formulation can mimic the single-mask multitask setting by using a high enough regularization on the δ offsets but we used a default l_2 -regularization strength of $1e^{-2}$ for all experiments to exhibit robustness of our method. It is interesting to note that for the ACL_ARC task, all transfer learning approaches at 95% sparsity outperform training the full model on task data only. Note from Table 4 that ACL_ARC is our smallest dataset. We posit that training the full, large model on this task leads to overfitting, resulting in poor generalization compared to leveraging transfer-learning at a reduced model size.

Figure 4 presents an interesting scenario where Chemprot benefits from transfer but RCT does not. Whilst this could be due to the fact that we perform limited hyper-parameter tuning (mainly to exhibit the robustness of our method and to reflect compute constrained settings), it is encouraging to see that

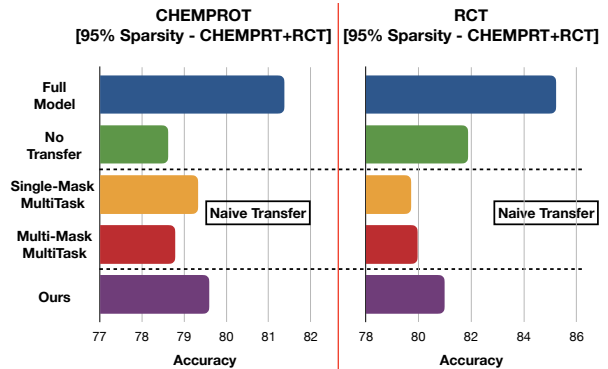


Figure 4: RCT and Chemprot performance at 95% sparsity. We see negative transfer from Chemprot to RCT across all transfer methods. Our proposed approach suffers least from performance degradation.

the δ -Formulation for coupling structural masks significantly helped dampen the impact of negative transfer in the case of RCT as the target task.

5.2 What should you transfer?

So far, we have discussed using the auxiliary task when learning both the structural variables and parameters of the pruned model. In this section, we investigate if transferring both is needed. We perform the following ablation at 95% sparsity to determine what is most important to transfer. For this, we assume that the *target task is not used during pruning* but is only introduced during the final fine-tuning of the smaller, pruned model.

Weights Only: We learn model weights and structural mask for the auxiliary task only. We then generate a random structural mask at the appropriate sparsity level (95%) and extract the model weights corresponding to this mask from the model trained on the transfer task. We then fine-tune this smaller, pruned model on the target task.

Masks Only: We learn model weights and structural mask for the auxiliary/transfer task only. We then reset the model weights to the pre-trained (not-yet-finetuned) state. Given the learned mask from the transfer task, and the untuned model weights, we then fine-tune this pruned model on the target task.

Masks and Weights: We use the transfer task to learn both the model weights and structural mask. We take weights and masks of this small model and fine-tune it on the target task.

Table 1 shows the results of this ablation. For both STSB \rightarrow MRPC and MRPC \rightarrow STSB, we see

that if we are only introducing the target task in the fine-tuning stage, it is beneficial to transfer both the weights and structure that are learned from the auxiliary task.

5.3 When should you transfer?

Table 2 shows results for the different choices presented in Section 3.2.2 relating to when to introduce the transfer task. These experiments are also conducted at a target sparsity of 95%.

We obtain the best performance with the $\text{Prune}(A, T) \rightarrow \text{FT}(T)$ and $\text{Prune}(A, T) \rightarrow \text{FT}(A, T)$ approaches. This matches intuition because we expect an appropriately chosen auxiliary task to be helpful in terms of learning both structure and parameters of the final pruned model. Thus, introducing it in the first (pruning) stage mitigates the challenge that is exacerbated by learning a larger set of variables from limited data.

5.4 How should we choose the transfer task?

Table 3 contains experimental results highlighting our investigation of different variables that can impact the quality of a transfer task.

We get the best improvements when we use a high resource auxiliary task from the same domain as the target task. As mentioned in Section 3.2.1, we use domain as a proxy for task relatedness. We see that even using a high resource task that is out-of-domain with respect to the end-task (RCT) can improve generalization over not introducing a transfer task at all (85.29 versus 79.2 for MRPC) and (0.873 versus 0.863 for STSB).

5.5 Does the learned structured sparsity translate to hardware speedups?

So far, we have only discussed the impact of transfer learning on generalization with respect to the end-task. However, when generating pruned models, we not only care about their generalization but also the degree of speedup that is achieved at the target sparsity.

Taking SCIIE as our primary task and ACL-ARC as the transfer task, we explore the accuracy-speedup tradeoff that is induced by leveraging transfer learning for structured pruning. We vary the degree of compression from 40% sparsity to 98%. To benchmark speed, we use the wall-clock time required to perform inference on the full SCIIE dataset through the model using at a batch-size of 128. All experiments were conducted

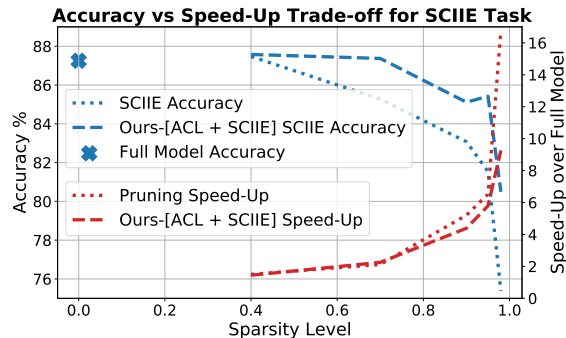


Figure 5: Accuracy vs speed tradeoff on SCIIE. Our δ -formulation, gives accuracy boosts on SCIIE at varying levels of compression.

on NVIDIA V100 GPUs. Figure 5 summarises our findings. For SCIIE+ACL, we fix the task weighting to the best performing configuration from our 95% sparsity experiments, the rest of the hyper-parameters are cross-validated from values in Table 5. At 50 \times compression (95%) sparsity, we are able to obtain a $\sim 5\%$ boost in accuracy over not using a transfer task, whilst achieving a $\sim 10\times$ speedup in inference. Transfer learning enables a more graceful degradation in accuracy (dashed blue line) whilst still finding pruned models with comparable speed-ups.

Another view of Figure 5 is to consider the model size required for a threshold level of accuracy for deployment. At a threshold of 84% accuracy, whilst naive pruning results would produce a model at 80% sparsity, we are able to produce one at 95% sparsity! This is a $\sim 1.2\times$ memory saving and $\sim 2.8\times$ inference speedup.

5.6 What are the structural differences between a pruned model using transfer learning and without?

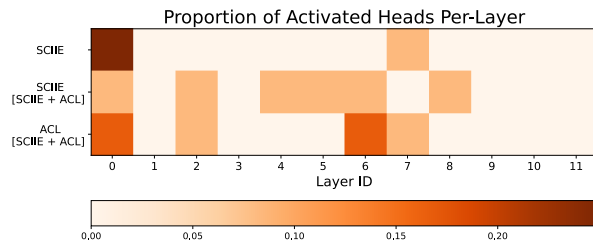


Figure 6: Structural visualization at 98% sparsity. Qualitatively, using a transfer task changes the pruned model structure significantly. The ACL transfer task in this case induces the learned SCIIE structure to be more diffuse across the layers of the model.

Table 1: Transferring structure, weights or both on STSB and MRPC? It is most beneficial to transfer both the learned weights and structural variables (masks)

	Metric	No Transfer	Weights Only	Structure Only	Both
STSB \rightarrow MRPC	Accuracy %	79.2	68.4 (\downarrow)	76.96 (\downarrow)	79.7 (\uparrow)
MRPC \rightarrow STSB	Pearson C.	0.868	0.23 (\downarrow)	0.8527 (\downarrow)	0.871 (\uparrow)

Table 2: When to introduce each task on MRPC and STSB? We find that it is optimal to prune with both the auxiliary and target task jointly.

	Metric	No Transfer	Prune(A) \rightarrow FT(T)	Prune(T) \rightarrow FT(A, T)	Prune(T, A) \rightarrow FT(T)	Prune(T, A) \rightarrow FT(A, T)
STSB \rightarrow MRPC	Accuracy %	79.2	79.7 (\uparrow)	83.09 (\uparrow)	83.82 (\uparrow)	84.56 (\uparrow)
MRPC \rightarrow STSB	Pearson C.	0.868	0.871 (\uparrow)	0.861 (\downarrow)	0.8751 (\uparrow)	0.872 (\uparrow)

Table 3: Selecting the auxiliary task. A high-resource, in-domain task leads to the best result. For all experiments, best results from hyper-parameter search are reported. All models (except Full BERT) are pruned to 95% sparsity.

Target	Full BERT	No Transfer	Domain	Resourced-ness	Transfer Task	Performance
MRPC	83.48	79.2	In-Domain	High (364k)	QQP	85.78
			In-Domain	Low (7k)	STSB	83.82
			Out-of-Domain	High (180k)	RCT	85.29
STSB	0.901	0.868	In-Domain	High (364k)	QQP	0.877
			In-Domain	Low (3.7k)	MRPC	0.875
			Out-of-Domain	High (180k)	RCT	0.873

At extreme sparsity levels, the differences in speedup from learning a pruned model with and without transfer learning (Figure 5) suggest that the models discovered have different structures.

Figures 6 and 7 show the fraction of attention heads and MLP intermediate dimensions respectively, that are preserved across each layer with respect to the original BERT_{base} model. Pruning with the target task alone results most of the preserved parameters coming from earlier in the network. With an auxiliary task however, the pattern of preserved modules is more diffuse across layers. We posit that this results from the the two tasks being multi-tasked preferring different layers thus resulting in a more diffuse distribution of preserved modules as a compromise in order to perform reasonably well on both tasks.

6 Conclusion

As coined in Ahia et al. (2021), the *low-resource double bind* describes the challenge of producing compressed models to serve compute-starved (memory and latency limits) tasks under a setting where these tasks also have limited data for pruning. In this work, we have explored adapting transfer

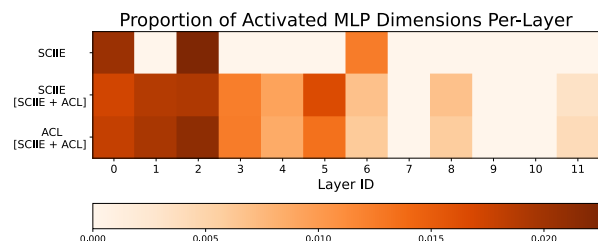


Figure 7: Structural visualization at 98% sparsity. Qualitatively, using a transfer task changes the pruned model structure significantly. The ACL transfer task in this case induces the learned SCIE structure to be more diffuse across the layers of the model.

learning, which has traditionally been leveraged only for learning model weights, to robustly prune models when the target task is data-limited.

We have provided practitioners with recommendations on how to choose a transfer task, when and how to incorporate it into the pruning optimization procedure and what elements to transfer from the auxiliary task to the target. Equipped with this knowledge, we plan to explore the problem of structured pruning under limited target data for larger scale models.

References

- Orevaoghene Ahia, Julia Kreutzer, and Sara Hooker. 2021. The low-resource double bind: An empirical study of pruning for low-resource machine translation. *arXiv preprint arXiv:2110.03036*.
- Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. 2023. Palm 2 technical report. *arXiv preprint arXiv:2305.10403*.
- Jonathan Baxter. 2000. A model of inductive bias learning. *Journal of artificial intelligence research*, 12:149–198.
- Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. 2021. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*.
- Rich Caruana. 1997. Multitask learning. *Machine learning*, 28:41–75.
- Franck Dernoncourt and Ji Young Lee. 2017. Pubmed 200k rct: a dataset for sequential sentence classification in medical abstracts. *arXiv preprint arXiv:1710.06071*.
- Lucio M Dery, Yann Dauphin, and David Grangier. 2021a. Auxiliary task update decomposition: The good, the bad and the neutral. *arXiv preprint arXiv:2108.11346*.
- Lucio M Dery, Paul Michel, Mikhail Khodak, Graham Neubig, and Ameet Talwalkar. 2022. Aang: Automating auxiliary learning. *arXiv preprint arXiv:2205.14082*.
- Lucio M Dery, Paul Michel, Ameet Talwalkar, and Graham Neubig. 2021b. Should we be pre-training? an argument for end-task aware training as an alternative. *arXiv preprint arXiv:2109.07437*.
- Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. 2022. Llm.int8(): 8-bit matrix multiplication for transformers at scale. *arXiv preprint arXiv:2208.07339*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Dumitru Erhan, Aaron Courville, Yoshua Bengio, and Pascal Vincent. 2010. Why does unsupervised pre-training help deep learning? In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 201–208. JMLR Workshop and Conference Proceedings.
- Angela Fan, Edouard Grave, and Armand Joulin. 2019. Reducing transformer depth on demand with structured dropout. *arXiv preprint arXiv:1909.11556*.
- Jonathan Frankle and Michael Carbin. 2018. The lottery ticket hypothesis: Finding sparse, trainable neural networks. *arXiv preprint arXiv:1803.03635*.
- Siddhant Garg, Lijun Zhang, and Hui Guan. 2023. Structured pruning for multi-task deep neural networks.
- Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A Smith. 2020. Don’t stop pretraining: Adapt language models to domains and tasks. *arXiv preprint arXiv:2004.10964*.
- Song Han, Jeff Pool, John Tran, and William Dally. 2015. Learning both weights and connections for efficient neural network. *Advances in neural information processing systems*, 28.
- Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and Qun Liu. 2019. Tinybert: Distilling bert for natural language understanding. *arXiv preprint arXiv:1909.10351*.
- David Jurgens, Srijan Kumar, Raine Hoover, Dan McFarland, and Dan Jurafsky. 2018. Measuring the evolution of a scientific field through citation frames. *Transactions of the Association for Computational Linguistics*, 6:391–406.
- Jens Kringelum, Sonny Kim Kjaerulff, Søren Brunak, Ole Lund, Tudor I Oprea, and Olivier Taboureau. 2016. Chemprot-3.0: a global chemical biology diseases mapping. *Database*, 2016:bav123.
- Zhuang Liu, Mingjie Sun, Tinghui Zhou, Gao Huang, and Trevor Darrell. 2018. Rethinking the value of network pruning. *arXiv preprint arXiv:1810.05270*.
- Christos Louizos, Max Welling, and Diederik P. Kingma. 2018. Learning sparse neural networks through l_0 regularization.
- Yi Luan, Luheng He, Mari Ostendorf, and Hannaneh Hajishirzi. 2018. Multi-task identification of entities, relations, and coreference for scientific knowledge graph construction. *arXiv preprint arXiv:1808.09602*.
- Xiaolong Ma, Sheng Lin, Shaokai Ye, Zhezhi He, Linfeng Zhang, Geng Yuan, Sia Huat Tan, Zhengang Li, Deliang Fan, Xuehai Qian, et al. 2021. Non-structured dnn weight pruning—is it beneficial in any platform? *IEEE transactions on neural networks and learning systems*, 33(9):4930–4944.
- Andreas Maurer, Massimiliano Pontil, and Bernardino Romera-Paredes. 2016. The benefit of multitask representation learning. *Journal of Machine Learning Research*, 17(81):1–32.
- Paul Michel, Omer Levy, and Graham Neubig. 2019. Are sixteen heads really better than one? *Advances in neural information processing systems*, 32.

Antonio Polino, Razvan Pascanu, and Dan Alistarh. 2018. Model compression via distillation and quantization. *arXiv preprint arXiv:1802.05668*.

Vivek Ramanujan, Mitchell Wortsman, Aniruddha Kembhavi, Ali Farhadi, and Mohammad Rastegari. 2020. What’s hidden in a randomly weighted neural network? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.

Victor Sanh, Thomas Wolf, and Alexander Rush. 2020. Movement pruning: Adaptive sparsity by fine-tuning. *Advances in Neural Information Processing Systems*, 33:20378–20389.

Mingjie Sun, Zhuang Liu, Anna Bair, and J Zico Kolter. 2023. A simple and effective pruning approach for large language models. *arXiv preprint arXiv:2306.11695*.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. [Attention is all you need](#). In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.

Elena Voita, David Talbot, Fedor Moiseev, Rico Sennrich, and Ivan Titov. 2019. Analyzing multi-head self-attention: Specialized heads do the heavy lifting, the rest can be pruned. *arXiv preprint arXiv:1905.09418*.

Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. 2018. Glue: A multi-task benchmark and analysis platform for natural language understanding. *arXiv preprint arXiv:1804.07461*.

Ziheng Wang, Jeremy Wohlwend, and Tao Lei. 2019. Structured pruning of large language models. *arXiv preprint arXiv:1910.04732*.

Mengzhou Xia, Zexuan Zhong, and Danqi Chen. 2022. Structured pruning learns compact and accurate models. *arXiv preprint arXiv:2204.00408*.

Nakyeong Yang, Yunah Jang, Hwanhee Lee, Seohyeong Jung, and Kyomin Jung. 2023. [Task-specific compression for multi-task language models using attribution-based pruning](#).

Unlike the original CoFi, we turn off output prediction distillation for all experiments. ie – we do not distill the predictions from the pre-trained models since unlike in the original CoFi paper, we are not starting from a model that has already been fine-tuned on the target task but rather we are starting from the pre-trained model itself.

During pruning, we perform 10K gradient descent steps to learn both the structural and parameter variables of the model. We perform 20 epochs of post-pruning finetuning on the target task.

A Datasets

Table 4 describes the tasks and datasets used in our experiments.

B Training Details

We follow as closely as possible the hyper-parameters that are used in the original CoFi code base. Table 5 reports CoFi-specific hyper-parameter settings.

Table 4: Specifications of datasets used to evaluate our methods.

Domain	Task	Task-Type	Train Size	Metric
BIOMED	CHEMPROT (Kringelum et al., 2016)	Classification	4169	Accuracy
	RCT (Dernoncourt and Lee, 2017)	Classification	10K*	Accuracy
CS	SCIIE (Luan et al., 2018)	Classification	3219	Accuracy
	ACL-ARC (Jurgens et al., 2018)	Classification	1688	Accuracy
GLUE	STS-B (Wang et al., 2018)	Sentence Similarity	7K	Pearson’s Correlation
	MRPC (Wang et al., 2018)	Paraphrase Detection	3.7K	Accuracy

Table 5: Hyper-parameter choices

Hyper-parameter	Values	Description
Task pair weightings	(1, 1), (1, 2), (2, 1)	Weightings applied to transfer task vs target task during training.
Model LR - Pruning	1e-4, 2e-5	Learning rate used for model parameters during pruning.
Model LR - Finetuning	1e-4, 2e-5	Learning rate used for finetuning pruned model.
Structure LR	0.1, 0.01	Learning rate used for learning structural parameters.
δ - l_2 Reg Weight	1e-2	Regularization weight used in δ -formulation.