

Self-Distilled Pruning of Neural Networks

Anonymous ACL submission

Abstract

Pruning aims to reduce the number of parameters while maintaining performance close to the original network. This work proposes a novel *self-distillation* based pruning strategy, whereby the representational similarity between the pruned and unpruned versions of the same network is maximized. Unlike previous approaches that treat distillation and pruning separately, we use distillation to inform the pruning criteria, without requiring a separate student network as in knowledge distillation. We show that the proposed *cross-correlation objective for self-distilled pruning* implicitly encourages sparse solutions, naturally complementing magnitude-based pruning criteria. Experiments on the GLUE and XGLUE benchmarks show that self-distilled pruning increases mono- and cross-lingual language model performance. Self-distilled pruned models also outperform smaller Transformers with an equal number of parameters and are competitive against (6 times) larger distilled networks. We also observe that self-distillation (1) maximizes class separability, (2) increases the signal-to-noise ratio, and (3) converges faster after pruning steps, providing further insights into why self-distilled pruning improves generalization.

1 Introduction

Neural network pruning (Mozer and Smolensky, 1989; Karnin, 1990; Reed, 1993) zeros out weights of a pretrained model with the aim of reducing parameter count and storage requirements, while maintaining performance close to the original model. The criteria to choose which weights to prune has been an active research area over the past three decades (Karnin, 1990; LeCun et al., 1990; Han et al., 2015a; Anwar et al., 2017; Molchanov et al., 2017). Lately, there has been a focus on pruning models in the transfer learning setting whereby a self-supervised pretrained model trained on a large amount of unlabelled data is fine-tuned to

a downstream task while weights are simultaneously pruned. In this context, recent work proposes to learn important scores over weights with a continuous mask and prune away those that having the smallest scores (Mallya et al., 2018; Sanh et al., 2020). However, these learned masks double the number of parameters in the network, requiring twice the number of gradient updates to tune the original parameters *and* their continuous masks (Sanh et al., 2020). Ideally, we aim to perform task-dependent fine-pruning *without* adding more parameters to the network, or at least far fewer than twice the count. Additionally, we desire pruning methods that can recover from performance degradation directly after pruning steps, faster than current pruning methods while encoding task-dependent information into the pruning process. To this end, we hypothesize self-distillation may recover performance faster after consecutive pruning steps, which becomes more important with larger performance degradation at a higher compression regime. Additionally, self-distillation has shown to encourage sparsity as the training error tends to 0 (Mobahi et al., 2020). This implicit sparse regularization effect complements magnitude-based pruning, an efficient and well-established pruning approach.

Hence, this paper proposes to combine self-distillation and magnitude-based pruning to achieve task-dependent pruning efficiently. This is achieved by *maximizing the cross-correlation* between output representations of the fine-tuned pretrained network and a pruned version of the same network – referred to as *self-distilled pruning* (SDP). Cross-correlation maximization has shown to reduce redundancy and encourage sparse solutions (Zbontar et al., 2021), naturally fitting with magnitude-based pruning. This sets state of the art results for magnitude-based pruning. Unlike typical knowledge distillation (KD) where the student is a separate network trained from random

083 initialization, here the student is initially a masked
084 version of the teacher. We then provide three in-
085 sights as to why self-distillation leads to more gen-
086 eralizable pruned networks. We observe that self-
087 distilled pruning (1) *recovers performance faster*
088 after pruning steps (i.e., improves convergence), (2)
089 *maximizes the signal-to-noise ratio* (SNR), where
090 pruned weights are considered as noise, and (3)
091 *improves the fidelity* between pruned and unpruned
092 representations as measured by mutual information
093 of the respective penultimate layers. We focus on
094 pruning fine-tuned monolingual *and* cross-lingual
095 transformer models, namely BERT (Devlin et al.,
096 2018) and XLM-RoBERTa (Conneau et al., 2019).
097 To our knowledge, this is the first study that in-
098 troduces the concept of *self-distilled pruning*, an-
099 alyzes iterative pruning in the mono-lingual *and*
100 cross-lingual settings on the GLUE and XGLUE
101 benchmarks respectively and the only work to in-
102 clude an evaluation of pruned model performance
103 in the cross-lingual transfer setting.

104 2 Background and Related Work

105 **Regularization-based pruning** can be achieved
106 by using a weight regularizer that encourages net-
107 work sparsity. Three well-established regularizers
108 are L_1 , L_2 and L_0 weight regularization (Louizos
109 et al., 2017; Liu et al., 2017; Ye et al., 2018) for
110 weight sparsity (Han et al., 2015b,a). For struc-
111 tured pruning, Group-wise Brain Damage (Lebe-
112 dev and Lempitsky, 2016) and SSL (Wen et al.,
113 2016) propose to use Group LASSO (Yuan and
114 Lin, 2006) to prune whole structures (e.g., convo-
115 lution blocks or blocks within standard linear lay-
116 ers). Park et al. (2020) avoid pruning small weights
117 if they are connected to larger weights in consecu-
118 tive layers and vice-versa, by penalizing the Frobe-
119 nius norm between pruned and unpruned layers to
120 be small. **Importance-based pruning** assigns a
121 score for each weight in the network and removes
122 weights with the lowest importance score. The sim-
123 plest scoring criteria is magnitude-based pruning
124 (MBP), which uses the lowest absolute value (LAV)
125 as the criteria (Reed, 1993; Han et al., 2015b,a)
126 or L_1/L_2 -norm for structured pruning (Liu et al.,
127 2017). MBP can be seen as a zero-th order pruning
128 criteria. However higher order pruning methods ap-
129 proximate the difference in pruned and unpruned
130 model loss using a Taylor series expansion up until
131 1st order (LeCun et al., 1990; Hassibi and Stork,
132 1993) or the 2nd order, which requires approximat-
133 ing the Hessian matrix (Martens and Grosse, 2015;

134 Wang et al., 2019; Singh and Alistarh, 2020) for
135 scalability. Lastly, the regularization-based pruning
136 is commonly used with importance-based pruning
137 e.g using L_2 weight regularization alongside MBP.

Knowledge Distillation (KD) transfers the log-
138 its of an already trained network (Hinton et al.,
139 2015) and uses them as soft targets to optimize
140 a student network. The student network is typi-
141 cally smaller than the teacher network and benefits
142 from the additional information soft targets pro-
143 vide. There has been various extensions that involve
144 distilling intermediate representations (Romero
145 et al., 2014), distributions (Huang and Wang, 2017),
146 maximizing mutual information between student
147 and teacher representations (Ahn et al., 2019), us-
148 ing pairwise interactions for improved KD (Park
149 et al., 2019) and contrastive representation distilla-
150 tion (Tian et al., 2019; Neill and Bollegala, 2021).

Self-Distillation is a special case of KD whereby
152 the student and teacher networks have the same
153 capacity. Interestingly, self-distilled students often
154 generalize better than the teacher (Furlanello et al.,
155 2018; Yang et al., 2019), however the mechanisms
156 by which self-distillation leads to improved gener-
157 alization remains somewhat unclear. Recent works
158 have provided insightful observations of this phe-
159 nomena. For example, (Stanton et al., 2021) have
160 shown that soft targets make optimization easier
161 for the student when compared to the task-provided
162 one-hot targets. (Allen-Zhu and Li, 2020) view self-
163 distillation as implicitly combining ensemble learn-
164 ing and KD to explain the improvement in test
165 accuracy when dealing with multi-view data. The
166 core idea is that the self-distillation objective re-
167 sults in the network learning a unique set of fea-
168 tures that are distinct from the original model, similar
169 to features learned by combining the outputs of
170 independent models in an ensemble. Given this
171 background on pruning and distillation, we now
172 describe our proposed methodology for *SDP*.
173

174 3 Proposed Methodology

175 We begin by defining a dataset $\mathcal{D} := \{(X_i, y_i)\}_{i=1}^D$
176 with single samples $s_i = (X_i, y_i)$, where each
177 X_i (in the D training samples) consists of a se-
178 quence of vectors $X_i := (\mathbf{x}_1, \dots, \mathbf{x}_N)$ and $\mathbf{x}_i \in$
179 \mathbb{R}^d . For structured prediction (e.g., NER, POS)
180 $y_i \in \{0, 1\}^{N \times C}$, and for single and pairwise sen-
181 tence classification, $y_i \in \{0, 1\}^C$, where C is the
182 number of classes. Let $\mathbf{y}^S = f_\theta(X_i)$ be the out-
183 put prediction ($y^S \in \mathbb{R}^C$) from the student $f_\theta(\cdot)$
184 with pretrained parameters $\theta := \{\mathbf{W}_l, \mathbf{b}_l\}_{l=1}^L$ for

185 L layers. The input to each subsequent layer is de-
 186 noted as $z_l \in \mathbb{R}^{n_l}$ where $x := z_0$ for n_l number
 187 of units in layer l and the corresponding output
 188 activation as $A_l = g(z_l)$. The loss function for
 189 standard classification fine-tuning is defined as the
 190 cross $\ell_{CE}(\mathbf{y}^S, \mathbf{y}) := -\frac{1}{C} \sum_{i=1}^C y_c \log(\mathbf{y}_c^S)$.

191 For self-distilled pruning, we also require an al-
 192 ready fine-tuned teacher network f_Θ , that has been
 193 tuned from the pretrained state f_θ , to retrieve the
 194 soft teacher labels $\mathbf{y}^T := f_\Theta(x)$, where $\mathbf{y}^T \in \mathbb{R}^C$
 195 and $\sum_c y_c^T = 1$. The soft label \mathbf{y}^T can be more
 196 informative than the one-hot targets \mathbf{y} used for stan-
 197 dard classification as they implicitly approximate
 198 pairwise class similarities through logit probabili-
 199 ties. The Kullback-Leibler divergence ℓ_{KLD} is then
 200 used with the main task cross-entropy loss ℓ_{CE} to
 201 express $\ell_{SDP-KLD}$ as shown in Equation 1,

$$202 \ell_{SDP-KLD} = (1-\alpha)\ell_{CE}(\mathbf{y}^S, \mathbf{y}) + \alpha\tau^2 D_{KLD}(\mathbf{y}^S, \mathbf{y}^T) \quad (1)$$

203 where $D_{KLD}(\mathbf{y}^S, \mathbf{y}^T) = \mathbb{H}(\mathbf{y}^T) - \mathbf{y}^T \log(\mathbf{y}^S)$,
 204 $\mathbb{H}(\mathbf{y}^T) = -\mathbf{y}^T \log(\mathbf{y}^T)$ is the entropy of the teacher
 205 distribution and τ is the softmax temperature. Fol-
 206 lowing (Hinton et al., 2015), the weighted sum of
 207 cross-entropy loss and KLD loss shown in Equa-
 208 tion 1 is used as our main SDP-based KD loss
 209 baseline, where $\alpha \in [0, 1]$. After each pruning
 210 step during iterative pruning, we aim to recover
 211 the immediate performance degradation by min-
 212 imizing $\ell_{SDP-KLD}$. In our experiments, we use
 213 weight magnitude-based pruning as the criteria for
 214 SDP given MBP’s flexibility, scalability and minis-
 215 cule computation overhead (only requires a binary
 216 tensor multiplication to be applied for each linear
 217 layer at each pruning step). However, D_{KLD} only
 218 distills the knowledge from the soft targets which
 219 may not propagate enough information about the
 220 intermediate dynamics of the teacher, nor does it
 221 penalize representational redundancy. This brings
 222 us to our proposed cross-correlation SDP objective.

223 3.1 Maximizing Cross-Correlation Between 224 Pruned and Unpruned Embeddings

225 Iterative pruning can be viewed as progressively
 226 adding noise $\mathbf{M}_l \in \{0, 1\}^{n_{l-1} \times n_l}$ to the weights
 227 $\mathbf{W}_l \in \mathbb{R}^{n_{l-1} \times n_l}$. Thus, as the pruning steps in-
 228 crease, the outputs become noisier and the rela-
 229 tionship between the inputs and outputs becomes
 230 weaker. Hence, a correlation measure is a natural
 231 choice for dealing with such pruning-induced noise.
 232 To this end, we use a cross-correlation loss to maxi-
 233 mize the correlation between the output representa-
 234 tions of the last hidden state of the pruned network

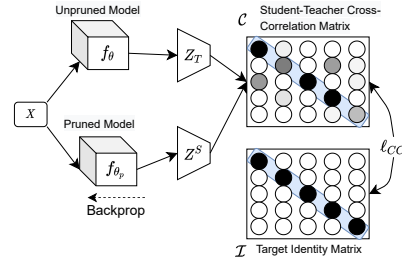


Figure 1: Self-Distilled Pruning with a Cross-Correlation Knowledge Distillation Loss.

235 and the unpruned network to reduce the effects of
 236 this pruning noise. The proposed *cross-correlation*
 237 *SDP loss function*, ℓ_{CC} , is expressed in Equation 2,
 238 where λ controls the importance of minimizing the
 239 non-adjacent pairwise correlations between z^S and
 240 z^T in the correlation matrix C . Here, m denotes the
 241 sample index in a mini-batch of M samples. Unlike
 242 ℓ_{KLD} , this loss is applied to the outputs of the last
 243 hidden layer as opposed to the classification logit
 244 outputs. Thus, we have,

$$245 \ell_{CC} := \sum_i (1 - C_{ii})^2 + \lambda \sum_i \sum_{j \neq i} C_{ij}^2 \quad (2)$$

$$246 \text{ such that } C_{ij} := \frac{\sum_m z_{m,i}^S z_{m,j}^T}{\sqrt{\sum_m (z_{m,i}^S)^2} \sqrt{\sum_m (z_{m,j}^T)^2}}.$$

247 Maximizing correlation along the diagonal of
 248 C makes the representations invariant to pruning
 249 noise, while minimizing the off-diagonal term
 250 decorrelates the components of the representations
 251 that are batch normalized. To reiterate, z^S is ob-
 252 tained from the pruned version of the network (f_{θ_p})
 253 and z^T is obtained from the unpruned version (f_θ).
 254 Since the learned output representations should be
 255 similar if their inputs are similar, we aim to address
 256 the problem where a correlation measure may pro-
 257 duce representations that are instead *proportional*
 258 to their inputs. To address this, batch normaliza-
 259 tion is used across mini-batches to stabilize the
 260 optimization when using the cross-correlation loss,
 261 avoiding local optima that correspond to degen-
 262 erate representations that do not distinguish pro-
 263 portionalities. In our experiments, this is used with
 264 the classification loss and KLD distillation loss as
 265 shown in Equation 3.

$$266 \ell_{SDP-CC} = (1 - \alpha)\ell_{CE}(\mathbf{y}^S, \mathbf{y}) + \alpha\tau^2 D_{KLD}(\mathbf{y}^S, \mathbf{y}^T) + \beta\ell_{CC}(z^S, z^T) \quad (3)$$

267 Figure 1 illustrates the proposed framework of *Self-*
 268 *Distilled Pruning with cross-correlation loss* (SDP-
 269 CC), where I is the identity matrix. Additionally,
 270 we provide a PyTorch based pseudo-code for SDP-
 271 CC the supplementary material.

3.2 A Frobenius Distortion Perspective of Self-Distilled Pruning

To formalize the objective being minimized when using MBP with self-distillation, we take the view of *Frobenius distortion minimization* (FDM; Dong et al., 2017) which says that layer-wise MBP is equivalent to minimizing the Frobenius distortions of a single layer. This can be described as $\min_{\mathbf{M}: \|\mathbf{M}\|_0=p} \|\mathbf{W} - \mathbf{M} \odot \mathbf{W}\|_F$, where \odot is the Hadamard product and p is a constraint of the number of weights to remove as a percentage of the total number of weights for a layer. Therefore, the output distortion is approximately the product of single layer Frobenius distortions. However, this minimization only defines a 1st order approximation of pruning induced Frobenius distortions which is a loose approximation for deep networks. In contrast, the \mathbf{y}^T targets provide higher-order information outside of the l -th layer being pruned in this FDM framework because Θ encodes information of all neighboring layers. Hence, we reformulate the FDM problem for SDP as an approximately higher-order MBP method as in Equation 4 where \mathbf{W}^T are the weights in f_Θ .

$$\min_{\mathbf{M}: \|\mathbf{M}\|_0=p} \left[\|\mathbf{W} - \mathbf{M} \odot \mathbf{W}\|_F + \lambda \|\mathbf{W}^T - \mathbf{M} \odot \mathbf{W}\|_F \right] \quad (4)$$

As described in (Dong et al., 2017; Hassibi and Stork, 1993), the difference in error can be approximated with a Taylor Series (TS) expansion as $\delta \mathcal{E}_l \approx \left(\frac{\partial \mathcal{E}_l}{\partial \mathbf{W}_l} \right)^\top \delta \mathbf{W}_l + \frac{1}{2} \delta \mathbf{W}_l^\top \mathbf{H}_l \delta \mathbf{W}_l + O(\|\delta \mathbf{W}_l\|^3)$ where \mathbf{H} is the Hessian matrix. When using SDP with a 1st TS, we can further express the TS approximation for SDP as shown in Equation 5, where \mathcal{E}_l^S is the error of the pruned network for task provided targets and \mathcal{E}_l^T are the errors of the pruned network with distilled logits.

$$\begin{aligned} (\mathcal{E}_l - \mathcal{E}_l^S)^2 + \lambda (\mathcal{E}_l - \mathcal{E}_l^T)^2 &\approx \delta \mathcal{E}_l^S + \delta \mathcal{E}_l^T \\ &\approx \left(\frac{\partial \mathcal{E}_l^S}{\partial \theta_l} \right)^\top \delta \theta_l + \lambda \left(\frac{\partial \mathcal{E}_l^T}{\partial \theta_l} \right)^\top \delta \theta_l \end{aligned} \quad (5)$$

3.3 How Does Self-Distillation Improve Pruned Model Generalization ?

We put forth the following insights as to the advantages provided by self-distillation for better pruned model generalization, and later experimentally demonstrate their validity.

Recovering Faster From Performance Degradation After Pruning Steps. The first explanation for why self-distillation leads to better generalization in iterative pruning is that the soft targets bias

the optimization and smoothen the loss surface through implicit similarities between the classes encoded in the logits. We posit this too holds true for performance recovery after pruning steps, as the classification boundaries become distorted due to the removal of weights. Faster convergence is particularly important for high compression rates where the performance drops become larger.

Implicit Maximization of the Signal-to-Noise Ratio. One explanation for faster convergence is that optimizing for soft targets translates to maximizing the margin of class boundaries given the implicit class similarities provided by teacher logits. Intuitively, task provided one-hot targets do not inform SGD of how similar incorrect predictions are to the correct class, whereas the teacher logits do, to the extent they have learned on the same task. To measure this, we use a formulation of the signal-to-noise ratio¹ (SNR) to measure the class separability and compactness differences between pruned model representations trained with and without self-distillation. We formulate SNR as Equation 6, where for a batch of inputs \mathbf{X} , we obtain \mathbf{Z} output representations from the pruned network, which contain samples with C classes where each class has the same N number of samples. The numerator expresses the average ℓ_2 inter-class distance between instances of each class pair and the denominator expresses the intra-class distance between instances within the same class.

$$\frac{1/N(C-1)^2 \sum_n \sum_{c=1}^C \sum_{i \neq c}^C \|\sqrt{\mathbf{Z}_{c,n}} - \sqrt{\mathbf{Z}_{i,n}}\|_2}{1/C(N-1)^2 \sum_{c=1}^C \sum_n \sum_{j \neq n}^N \|\sqrt{\mathbf{Z}_{c,n}} - \sqrt{\mathbf{Z}_{c,j}}\|_2} \quad (6)$$

This estimation is $C - 1 \binom{C+1}{2}$ in the number of pairwise distances to be computed between the inter-class distances for the classes. For large output spaces (e.g., language modeling) we recommend defining the top k -NN classes for each class and estimate their distances on samples from them.

Quantifying Fidelity Between Pruned Models Trained With and Without Self-Distillation. A natural question to ask is *how much generalization power does the distilled soft targets provide when compared to the task provided one-hot targets ?* If best generalization is achieved when $\alpha = 1$ in Equation 1, this implies that the pruned network should have as high fidelity as possible with the unpruned network. However, as we will see there is a bias-variance trade-off between fidelity and

¹A measure typically used in signal processing to evaluate signal quality.

generalization performance, i.e., $\alpha = 1$ is not optimal in most cases. To measure fidelity between SDP representations and standard fine-tuned representations, we compute their *mutual information* (MI) and compare this to the MI between representations of pruned models without self-distillation and standard fine-tuned models. The MI between continuous variables can be expressed as,

$$\hat{I}(\mathbf{Z}^T; \mathbf{Z}^S) = H(\mathbf{Z}^T) - H(\mathbf{Z}^T | \mathbf{Z}^S) = -\mathbb{E}_{\mathbf{Z}^T}[\log p(\mathbf{Z}^T)] + \mathbb{E}_{\mathbf{Z}^T, \mathbf{Z}^S}[\log p(\mathbf{Z}^T | \mathbf{Z}^S)] \quad (7)$$

where $H(\mathbf{Z}^T)$ is the the entropy of the teacher representation and $H(\mathbf{Z}^T | \mathbf{Z}^S)$ is the conditional entropy that is derived from the joint distribution $p(\mathbf{Z}^T, \mathbf{Z}^S)$. This can also be expressed as the KL divergence between the joint probabilities and product of marginals as $I(\mathbf{Z}^T; \mathbf{Z}^S) = D_{\text{KLD}}[p(\mathbf{Z}^S, \mathbf{Z}^T) || p(\mathbf{Z}^S)p(\mathbf{Z}^T)]$. However, these theoretical quantities have to be estimated from test sample representations. We use a k -NN based MI estimator (Kraskov et al., 2004; Evans, 2008; Ver Steeg and Galstyan, 2013; Ver Steeg, 2000) which partitions the supports into a finite number of bins of equal size, forming a histogram that can be used to estimate $\hat{I}(\mathbf{Z}^S; \mathbf{Z}^T)$ based on discrete counts in each bin. This MI estimator is given as,

$$I(z^S; z^T) \approx \epsilon \left(\log \frac{\phi_{[z^S]}(i, k_{[z^S]}) \phi_{[z^T]}(i, k_{[z^T]})}{\phi_z(i, k)} \right) \quad (8)$$

where $\phi_{z^S}(i, k_{[z^S]})$ is the probability measure of the k -th nearest neighbour ball of $z^S \in \mathbb{R}^{n_L}$ and $\omega_{[z^T]}(i, k_{[z^T]})$ is the probability measure of the k_y -th nearest neighbour ball of $z^T \in \mathbb{R}^{n_L}$ where n_L is the dimension of the penultimate layer. In our experiments, we use 256 bins for the histogram with Gaussian smoothing and $k = 5$ (see (Kraskov et al., 2004) for further details).

4 Experimental Setup

Iterative Pruning Baselines. For XGLUE tasks, we perform 15 pruning steps on XLM-RoBERTA_{Base}, one per 15 epochs, while for the GLUE tasks, we perform 32 pruning steps on BERT_{Base}. The compression rate and number of pruning steps is higher for GLUE tasks compared to XGLUE, because GLUE tasks involve evaluation in the *supervised classification* setting; whereas in XGLUE we report in the more challenging *zero-shot cross-lingual transfer* setting with only a single language used for training (i.e., English). At each pruning step, we uniformly pruning

10% of the parameters for both the models. Although prior work suggests non-uniform pruning schedules (e.g., cubic schedule (Zhu and Gupta, 2017)), we did not see any major differences to uniform pruning. We compare the performance of the proposed SDP-CC method against the following baselines: **Random Pruning (MBP-Random)** - prunes weights uniformly at random across all layers. Random pruning can be considered as a lower bound on iterative pruning performance. **Layer-wise Magnitude Based Pruning (MBP)** - for each layer, prunes weights with the LAV. **Global Magnitude Pruning (Global-MBP)** - prunes the LAV of all weights in the network. **Layer-wise Gradient Magnitude Pruning (Gradient-MBP)** - for each layer, prunes the weights with the LAV of the accumulated gradients evaluated on a batch of inputs. **1st Taylor Series Pruning (TS)** - prunes weights based on the LAV of $|\text{gradient} \times \text{weight}|$. **L_0 norm MBP (Louizos et al., 2017)** - uses non-negative stochastic gates that choose which weights are set to zero as a smooth approximation to the non-differentiable L_0 -norm. **L_1 norm MBP (Li et al., 2016)** - applies L_1 weight regularization and uses MBP. **Lookahead pruning (LAP) (Park et al., 2020)** - prunes weight paths that have the smallest magnitude across blocks of layers, unlike MBP that does not consider neighboring layers. **Layer-Adaptive MBP (LAMP) (Lee et al., 2020)** - adaptively compute the pruning ratio per layer. For all above pruning methods we exclude weight pruning of the embeddings, layer normalization parameters and the last classification layer, as they play an important role for generalization and account for less than 1% of weights in both BERT and XLM-R_{Base}.

For **Knowledge Distillation** we also compare against a class of smaller knowledge distilled versions of BERT model with varying parameter sizes on the GLUE benchmark. We report prior results of *DistilBERT* (Sanh et al., 2019) and also mini-BERT models including *TinyBERT* (Jiao et al., 2019), *BERT-small* (Turc et al., 2019) and *BERT-medium* (Turc et al., 2019). In addition, we consider maximizing the cosine similarity between pruned and unpruned representations in the SDP loss, as $\ell_{\text{SDP-COS}} := \alpha \ell_{\text{CE}}(\mathbf{y}^S, \mathbf{y}) + \beta \left(1 - \frac{\mathbf{z}^S \cdot \mathbf{z}^T}{\|\mathbf{z}^S\| \|\mathbf{z}^T\|} \right)$. Unlike cross-correlation, there is no decorrelation of non-adjacent features in both representations for SDP-COS. This helps identify whether the redundancy reduction in cross-correlation is beneficial compared to the correlation loss that does not di-

Compression Method	Score (avg.)	Single Sentence		Similarity and Paraphrase			Natural Language Inference		
		CoLA (mcc)	SST-2 (acc)	MNLI (acc)	MRPC (f1/acc)	STS-B (pears./spear.)	QQP (f1/acc)	RTE (acc)	QNLI (acc)
BERT _{Base} (Ours)	84.06	53.24	90.71	80.27	80.9/77.7	83.5/83.8	83.9/88.0	68.59	86.91
Knowledge Distilled Baselines (% parameters w.r.t. original BERT)									
DistilBERT (60%)	82.85	51.3	91.3	82.2	87.5/-	86.9/-	-/-85.5	59.9	89.2
BERT-Medium (44.4%)	81.54	38.0	89.6	80.0	86.6/81.6	80.4/78.4	69.6/87.9	62.2	87.7
BERT-Small (20%)	79.02	27.8	89.7	77.6	83.4/76.2	78.8/77.0	68.1/87.0	61.8	86.4
BERT-Mini (10%)	76.97	0.0	85.9	75.1	74.8/74.3	75.4/73.3	66.4/86.2	57.9	84.1
BERT-Tiny (3.6%)	73.32	0.0	83.2	70.2	81.1/71.1	74.3/73.6	62.2/83.4	57.2	81.5
Pruning Baselines									
		20%	10%	10%	10%	10%	10%	10%	10%
Random	66.03	6.50	78.44	69.55	77.5/67.1	27.4/26.9	77.07/81.86	52.70	74.66
L_0 -MBP	77.25	31.68	83.37	75.61	78.4/68.2	75.9/75.7	81.56/86.49	64.26	82.62
L_2 -MBP	76.48	29.51	83.37	76.19	78.4/68.2	75.3/75.6	77.50/82.98	62.09	82.61
L_2 -Global-MBP	77.16	29.25	82.83	76.40	81.2/69.9	75.1/75.5	82.77/86.70	62.01	82.24
L_2 -Gradient-MBP	74.84	15.46	82.91	72.51	81.0/73.7	73.8/73.6	80.41/85.19	56.31	79.33
1 st -order Taylor	76.31	28.88	83.26	74.64	83.0/74.8	76.7/76.6	80.09/85.29	57.76	81.20
Lookahead	76.40	28.15	82.80	75.31	79.8/70.5	71.9/71.9	81.84/86.53	60.29	81.80
LAMP	74.03	20.31	83.26	74.27	72.3/63.7	73.7/74.1	79.32/85.07	58.84	81.09
Proposed Methodology									
L_2 -MBP + SDP-COS	77.83	31.80	86.00	75.68	81.6/72.2	76.4/76.3	81.39/86.68	61.73	83.07
L_2 -MBP + SDP-KLD	78.34	36.74	87.96	77.94	80.5/68.2	77.1/77.3	83.21/85.58	63.18	83.54
L_2 -MBP + SDP-CC	78.90	36.77	87.84	78.04	81.1/71.0	77.3/77.5	83.79/86.37	62.64	84.20

BERT- results reported from Sanh et al. (2019); Jiao et al. (2019); Turc et al. (2019) and MNLI results are for the matched dataset.

Table 1: GLUE benchmark results for pruned models @10% (or @20%) remaining weights.

rectly optimize this.

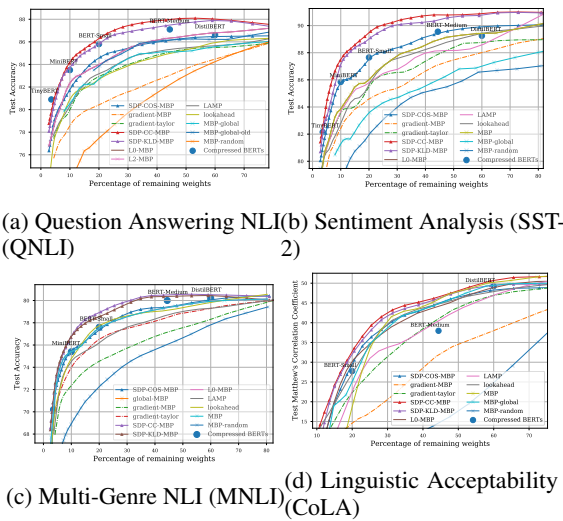


Figure 2: Iterative Pruning Test Performance on GLUE tasks.

5 Empirical Results

Pruning Results on GLUE. Table 1 shows the test performance across all GLUE tasks of the different models with varying pruning ratios, up to 10% remaining weights of original BERT_{Base} along with mini-BERT models (Sanh et al., 2019; Turc et al., 2019) of varying size. However, for the CoLA dataset, we report at 20% pruning as nearly all compression methods have an MCC score of 0, making the compressed method performance indistinguishable. For this reason, the GLUE score (Score) is

computed for all tasks and methods @10% apart from CoLA. The best performing compression method per task is marked in bold. We find that our proposed SDP approaches (all three variants) outperform against baseline pruning methods, with SDP-CC performing the best across all tasks. We note that for the tasks with fewer training samples (e.g., CoLA has 8.5k samples, STS-B has 7k samples and RTE has 3k samples), the performance gap is larger compared to BERT_{Base}, as the pruning step interval is shorter and less training data allows lesser time for the model to recover from pruning losses and also less data for teacher model to distil in the case of using SDP.

Smaller dense versions of BERT require more labelled data in order to compete with unstructured MBP and higher-order pruning methods such as 1st order Taylor series and Lookahead pruning. For example, we see BERT-Mini (@10%) shows competitive test accuracy with our proposed SDP-CC on QNLI, MNLI and QQP, the three datasets with the most training samples (105k, 393k and 364k respectively). Overall, L_2 -MBP + SDP-CC achieves the highest GLUE score for all models at 10% remaining weights when compared to BERT-Base parameter count. Moreover, we find that L_2 -MBP + SDP-CC achieves best performance for 5 of the 8 tasks, with 1 of the remaining 3 being from L_2 MBP+SDP-KLD. This suggests that redundancy reduction via a cross-correlation objective is use-

Prune Method	XNLI	NC	NER	PAWSX	POS	QAM	QADSM	WPR	Avg.
XLM-R _{Base}	73.48	80.10	82.60	89.24	80.34	68.56	68.06	73.32	76.96
Random	51.22	70.19	38.19	57.37	52.57	53.85	52.34	70.69	55.80
Global-Random	50.97	69.88	38.30	56.74	53.02	54.02	53.49	69.11	55.69
L_0 -MBP	64.75	78.98	56.22	72.09	71.38	59.31	53.35	71.70	65.97
L_2 -MBP	64.30	78.79	54.43	77.99	70.68	59.24	60.33	71.52	67.16
L_2 -Global-MBP	65.12	78.64	54.47	79.13	71.37	59.26	60.61	71.80	67.55
L_2 -Gradient-MBP	61.11	73.77	53.25	79.56	65.89	57.35	59.33	71.59	65.23
1 st -order Taylor	64.26	79.34	63.60	82.83	68.94	61.69	62.42	72.28	69.09
Lookahead	60.84	79.18	54.44	71.05	68.76	55.94	53.41	71.26	64.36
LAMP	58.04	63.64	51.92	66.05	67.43	55.36	52.42	71.09	60.74
L_2 -MBP + SDP-COS	64.96	79.02	62.77	78.70	72.88	60.21	60.94	72.04	68.94
L_2 -MBP + SDP-KLD	65.94	80.72	64.50	79.25	73.18	61.66	61.09	71.84	69.77
L_2 -MBP + SDP-CC	66.47	79.73	66.34	80.03	73.45	63.73	62.78	72.59	70.76

Table 2: XGLUE Iterative Pruning @ 30% Remaining Weights of XLM-R_{base} - Zero Shot Cross-Lingual Performance Per Task and Overall Average Score (Avg).

ful for SDP and clearly improve over SDP-COS which does not minimize correlations between off-diagonal terms. Figure 2 shows the performance across all pruning steps. Interestingly, for QNLI we observe the performance notably improves between 30-70% for SDP-CC and SDP-KLD. For SST-2, we observe a significant gap between SDP-KLD and SDP-CC compared to the pruning baselines and smaller versions of BERT, while TinyBERT becomes competitive at extreme compression (<4%). **Pruning Results on XGLUE.** We show the per task test performance and the *average task under-standing* score on XGLUE for pruning baselines and our proposed SDP approaches in Table 2. Our proposed cross-correlation objective for SDP again achieves the best average (Avg.) score and achieves the best task performance in 6 out of 8 tasks, while standard SDP-KLD achieves best performance on one (news classification) of the remaining two. Most notably, we outperform methods which use higher order gradient information (1st-order Taylor) at 30% remaining weights, which tends to be a point at which XLM-R_{Base} begins to degrade performance below 10% of the original fine-tuned test performance for SDP methods and competitive baselines. In Figure 3, we can observe this trend from the various tasks within XGLUE. We note that the number of training samples used for retraining plays an important role in the rate of performance degradation. For example, of the 6 presented XGLUE tasks, NER has the lowest number of training samples (15k) of all XGLUE tasks and also degrades the fastest in performance (from 90% to 50% Test F1 at 30% remaining weights). In comparison, XNLI has the most training samples for retraining (433k) and maintains performance relatively well, keeping within 10% of the original fine-tuned model at 30% remaining weights.

Summary of Results. From our experiments on GLUE and XGLUE task, we find that SDP consistently outperforms pruning, KD and smaller BERT baselines. SDP-KLD and SDP-CC both outperform larger sized BERT models (BERT-Small), somewhat surprisingly, given that BERT-Small (and the remaining BERT models) have the advantage of large-scale self-supervised pretraining, while pruning only has supervision from the downstream task. For NER in XGLUE, higher order pruning methods such as Taylor-Series pruning have an advantage at high compression rates mainly due to lack of training samples (only 15k). Apart from this low training sample regime, SDP with MBP dominates at high compression rates.

Measuring Fidelity To The Fine-Tuned Model. We now analyse the empirical evidence that soft targets used in SDP may force higher fidelity with the representations of the fine-tuned model when compared to using MBP without self-distillation. As described in subsection 3.3 we measure mutual dependencies between both representations of models with the best performing hyperparameter settings of α , β and the softmax temperature τ . We note that increasing the temperature τ translates to “peakier” teacher logit distributions, encouraging SGD to learn a student with high fidelity to the teacher. From the LHS of Figure 4, we can see that SDP models have higher mutual information (MI) with the teacher compared to MBP, which performs worse for PAWS-X (similar on remaining tasks, not shown for brevity). In fact, the rank order of the best performing pruned models at each pruning step has a direct correlation with MI, e.g., SDP-COS-MBP maintains highest MI and the highest test accuracy for PAWS-X for the same α . However, too high fidelity ($\alpha = 1.$) led to worse generalization compared to a balance between the task

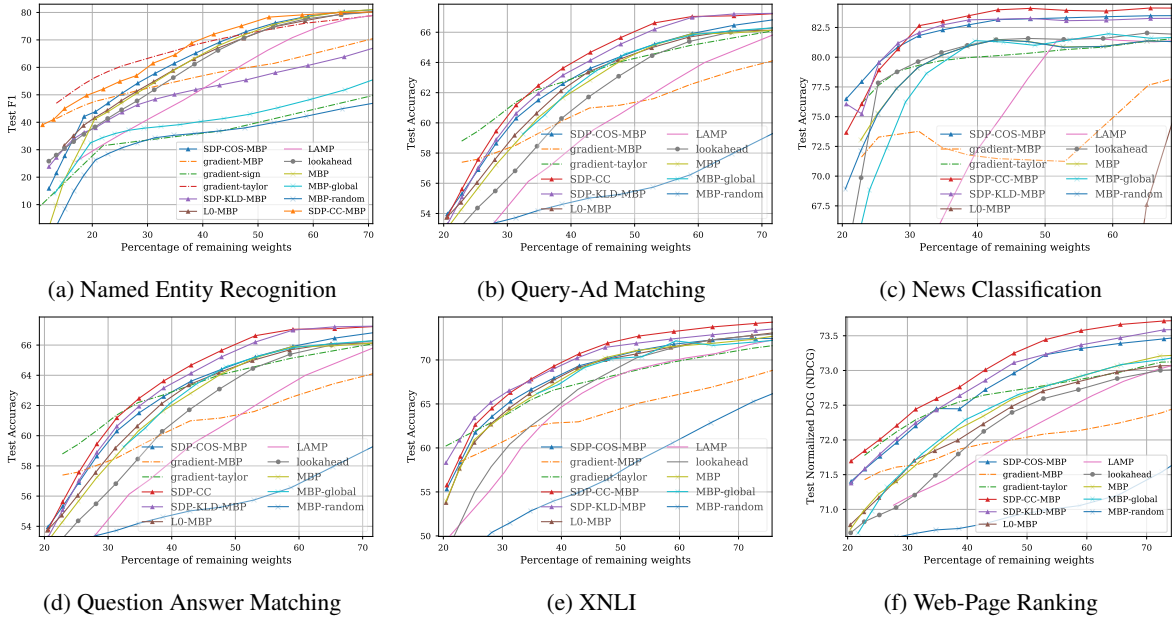


Figure 3: Zero-Shot Results After Iteratively Fine-Pruning XLM-R_{Base} on XGLUE tasks.

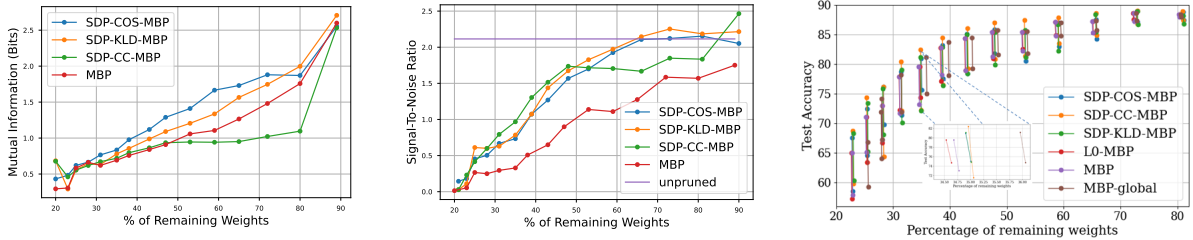


Figure 4: Mutual Information Between Unpruned and Pruned Representations (left) and Signal-To-Noise Ratio (middle) from PAWS-X Development Set Representations and (right) Pruning Performance Recovery with Self-Distilled Pruning.

580 provided targets and the teacher logits ($\alpha = 0.5$).
 581 **Self-Distilled Pruning Increases Class Separability and The Signal-To-Noise Ratio (SNR).** We
 582 also find that the SNR is increased at each pruning
 583 step as formulated in subsection 3.3. From this ob-
 584 servation, we find that *SDP-CC-MBP* using cross-
 585 correlation loss does particularly well in the 30%-
 586 50% remaining weights range. More generally, all
 587 3 SDP losses clearly lead to better class separabil-
 588 ity and class compactness across all pruning steps
 589 compared to MBP (i.e., no self-distillation).
 590

591 **Self-Distilled Pruning Recovers Faster Per-**
 592 **formance Degradation Directly After Pruning**
 593 **Steps.** In the right plot of Figure 4 we show how
 594 SDP with Magnitude pruning (SDP-MBP) recovers
 595 during training in between pruning steps. The top
 596 of each vertical bar is the recovery development
 597 accuracy and the bottom is the initial performance
 598 degradation prior to retraining. We see that SDP
 599 pruned models degrade in performance more than
 600 magnitude pruning without self-distillation. This
 601 suggests that SDP-MBP may force weights to be
 602 closer, as there is more initial performance degra-

603 dation if weights are not driven to zero. However,
 604 the recovery is faster. This may be explained by re-
 605 cent work that suggests the stability generalization
 606 tradeoff (Bartoldson et al., 2019).

6 Conclusion

607 In this paper, we proposed a novel *self-distillation*
 608 based pruning technique based on a *cross-*
 609 *correlation* objective. We extensively studied the
 610 confluence between pruning and self-distillation
 611 for masked language models and its enhanced utility
 612 on downstream tasks in both monolingual and
 613 multi-lingual settings. We find that self-distillation
 614 aids in recovering directly after pruning in iterative
 615 magnitude-based pruning, increases representa-
 616 tional fidelity with the unpruned model and
 617 implicitly maximize the signal-to-noise ratio. Ad-
 618 ditionally, we find our cross-correlation based
 619 self-distillation pruning objective minimizes neu-
 620 ronal redundancy and achieves state-of-the-art in
 621 magnitude-based pruning baselines, and even out-
 622 performs KD based smaller BERT models with
 623 more parameters.
 624

625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678

References

Sungsoo Ahn, Shell Xu Hu, Andreas Damianou, Neil D Lawrence, and Zhenwen Dai. 2019. Variational information distillation for knowledge transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9163–9171.

Zeyuan Allen-Zhu and Yuanzhi Li. 2020. Towards understanding ensemble, knowledge distillation and self-distillation in deep learning. *arXiv preprint arXiv:2012.09816*.

Sajid Anwar, Kyu Yeon Hwang, and Wonyong Sung. 2017. Structured pruning of deep convolutional neural networks. *ACM Journal on Emerging Technologies in Computing Systems (JETC)*, 13(3):1–18.

Brian R Bartoldson, Ari S Morcos, Adrian Barbu, and Gordon Erlebacher. 2019. The generalization-stability tradeoff in neural network pruning. *arXiv preprint arXiv:1906.03728*.

Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Unsupervised cross-lingual representation learning at scale. *arXiv preprint arXiv:1911.02116*.

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*.

Xin Dong, Shangyu Chen, and Sinno Jialin Pan. 2017. Learning to prune deep neural networks via layer-wise optimal brain surgeon. *arXiv preprint arXiv:1705.07565*.

Dafydd Evans. 2008. A computationally efficient estimator for mutual information. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 464(2093):1203–1215.

Tommaso Furlanello, Zachary Lipton, Michael Tschanen, Laurent Itti, and Anima Anandkumar. 2018. Born again neural networks. In *International Conference on Machine Learning*, pages 1607–1616. PMLR.

S Han, H Mao, and WJ Dally. 2015a. Compressing deep neural networks with pruning, trained quantization and huffman coding. *arXiv preprint*.

Song Han, Jeff Pool, John Tran, and William J Dally. 2015b. Learning both weights and connections for efficient neural networks. *arXiv preprint arXiv:1506.02626*.

Babak Hassibi and David G Stork. 1993. *Second order derivatives for network pruning: Optimal brain surgeon*. Morgan Kaufmann.

Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*.

Zehao Huang and Naiyan Wang. 2017. Like what you like: Knowledge distill via neuron selectivity transfer. *arXiv preprint arXiv:1707.01219*.

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*.

Ehud D Karnin. 1990. A simple procedure for pruning back-propagation trained neural networks. *IEEE transactions on neural networks*, 1(2):239–242.

Alexander Kraskov, Harald Stögbauer, and Peter Grassberger. 2004. Estimating mutual information. *Physical review E*, 69(6):066138.

Vadim Lebedev and Victor Lempitsky. 2016. Fast convnets using group-wise brain damage. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2554–2564.

Yann LeCun, John S Denker, and Sara A Solla. 1990. Optimal brain damage. In *Advances in neural information processing systems*, pages 598–605.

Jaeho Lee, Sejun Park, Sangwoo Mo, Sungsoo Ahn, and Jinwoo Shin. 2020. Layer-adaptive sparsity for the magnitude-based pruning. In *International Conference on Learning Representations*.

Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. 2016. Pruning filters for efficient convnets. *arXiv preprint arXiv:1608.08710*.

Zhuang Liu, Jianguo Li, Zhiqiang Shen, Gao Huang, Shoumeng Yan, and Changshui Zhang. 2017. Learning efficient convolutional networks through network slimming. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2736–2744.

Christos Louizos, Max Welling, and Diederik P Kingma. 2017. Learning sparse neural networks through l_0 regularization. *arXiv preprint arXiv:1712.01312*.

Arun Mallya, Dillon Davis, and Svetlana Lazebnik. 2018. Piggyback: Adapting a single network to multiple tasks by learning to mask weights. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 67–82.

James Martens and Roger Grosse. 2015. Optimizing neural networks with kronecker-factored approximate curvature. In *International conference on machine learning*, pages 2408–2417. PMLR.

Hossein Mobahi, Mehrdad Farajtabar, and Peter L Bartlett. 2020. Self-distillation amplifies regularization in hilbert space. *arXiv preprint arXiv:2002.05715*.

Dmitry Molchanov, Arsenii Ashukha, and Dmitry Vetrov. 2017. Variational dropout sparsifies deep neural networks. In *International Conference on Machine Learning*, pages 2498–2507. PMLR.

731	Michael C Mozer and Paul Smolensky. 1989. Skeletonization: A technique for trimming the fat from a network via relevance assessment. In <i>Advances in neural information processing systems</i> , pages 107–115.	784
732		785
733		786
734		787
735		788
736	James O’ Neill and Danushka Bollegala. 2021. Semantically-conditioned negative samples for efficient contrastive learning. <i>arXiv preprint arXiv:2102.06603</i> .	789
737		790
738		791
739		792
740	Sejun Park, Jaeho Lee, Sangwoo Mo, and Jinwoo Shin. 2020. Lookahead: A far-sighted alternative of magnitude-based pruning. <i>arXiv preprint arXiv:2002.04809</i> .	793
741		794
742		795
743		796
744	Wonpyo Park, Dongju Kim, Yan Lu, and Minsu Cho. 2019. Relational knowledge distillation. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 3967–3976.	797
745		798
746		799
747		800
748	Russell Reed. 1993. Pruning algorithms-a survey. <i>IEEE transactions on Neural Networks</i> , 4(5):740–747.	801
749		802
750	Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. 2014. Fitnets: Hints for thin deep nets. <i>arXiv preprint arXiv:1412.6550</i> .	803
751		804
752		805
753		806
754	Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. <i>arXiv preprint arXiv:1910.01108</i> .	807
755		808
756		809
757		810
758	Victor Sanh, Thomas Wolf, and Alexander M Rush. 2020. Movement pruning: Adaptive sparsity by fine-tuning. <i>arXiv preprint arXiv:2005.07683</i> .	811
759		812
760		
761	Sidak Pal Singh and Dan Alistarh. 2020. Woodfisher: Efficient second-order approximations for model compression. <i>arXiv preprint arXiv:2004.14340</i> .	
762		
763		
764	Samuel Stanton, Pavel Izmailov, Polina Kirichenko, Alexander A Alemi, and Andrew Gordon Wilson. 2021. Does knowledge distillation really work? <i>arXiv preprint arXiv:2106.05945</i> .	
765		
766		
767		
768	Yonglong Tian, Dilip Krishnan, and Phillip Isola. 2019. Contrastive representation distillation. <i>arXiv preprint arXiv:1910.10699</i> .	
769		
770		
771	Iulia Turc, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Well-read students learn better: On the importance of pre-training compact models. <i>arXiv preprint arXiv:1908.08962</i> .	
772		
773		
774		
775	Greg Ver Steeg. 2000. Non-parametric entropy estimation toolbox (npeet). Technical report, Technical Report. 2000. Available online: https://www.isi.edu/~gregv	
776		
777		
778		
779	Greg Ver Steeg and Aram Galstyan. 2013. Information-theoretic measures of influence based on content dynamics. In <i>Proceedings of the sixth ACM international conference on Web search and data mining</i> , pages 3–12.	
780		
781		
782		
783		
	Chaoqi Wang, Roger Grosse, Sanja Fidler, and Guodong Zhang. 2019. Eigendamage: Structured pruning in the kronecker-factored eigenbasis. In <i>International Conference on Machine Learning</i> , pages 6566–6575. PMLR.	
	Wei Wen, Chunpeng Wu, Yandan Wang, Yiran Chen, and Hai Li. 2016. Learning structured sparsity in deep neural networks. <i>arXiv preprint arXiv:1608.03665</i> .	
	Chenglin Yang, Lingxi Xie, Siyuan Qiao, and Alan L Yuille. 2019. Training deep neural networks in generations: A more tolerant teacher educates better students. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 33, pages 5628–5635.	
	Jianbo Ye, Xin Lu, Zhe Lin, and James Z Wang. 2018. Rethinking the smaller-norm-less-informative assumption in channel pruning of convolution layers. <i>arXiv preprint arXiv:1802.00124</i> .	
	Ming Yuan and Yi Lin. 2006. Model selection and estimation in regression with grouped variables. <i>Journal of the Royal Statistical Society: Series B (Statistical Methodology)</i> , 68(1):49–67.	
	Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and Stéphane Deny. 2021. Barlow twins: Self-supervised learning via redundancy reduction. <i>arXiv preprint arXiv:2103.03230</i> .	
	Michael Zhu and Suyog Gupta. 2017. To prune, or not to prune: exploring the efficacy of pruning for model compression. <i>arXiv preprint arXiv:1710.01878</i> .	