
SOSP: Efficiently Capturing Global Correlations by Second-Order Structured Pruning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Pruning neural networks reduces inference time and memory cost, as well as accel-
2 erates training when done at initialization. On standard hardware, these benefits
3 will be especially prominent if coarse-grained structures, like feature maps, are
4 pruned. We devise global saliency-based methods for second-order structured prun-
5 ing (SOSP) which include correlations among structures, whereas highest efficiency
6 is achieved by saliency approximations using fast Hessian-vector products. We
7 achieve state-of-the-art results for various object classification benchmarks, espe-
8 cially for large pruning rates highly relevant for resource-constrained applications.
9 We showcase that our approach scales to large-scale vision tasks, even though it
10 captures correlations across all layers of the network. Further, we highlight two
11 outstanding features of our methods. First, to reduce training costs our pruning
12 objectives can also be applied at initialization with no or only minor degradation in
13 accuracy compared to pruning after pretraining. Second, our structured pruning
14 methods allow to reveal architectural bottlenecks, which we remove to further
15 increase the accuracy of the networks.

16 1 Introduction

17 Deep neural networks have consistently grown in size over the last years with increasing performance.
18 However, this increase in size leads to slower inference, higher computational requirements and
19 higher cost. To reduce the size of the networks without affecting their performance, a large number
20 of pruning algorithms have been proposed (e.g., LeCun et al., 1990; Hassibi et al., 1993; Reed, 1993;
21 Han et al., 2015; Blalock et al., 2020). Pruning can either be *unstructured*, i.e. removing individual
22 weights, or *structured*, i.e. removing entire substructures like nodes or channels. Single-shot pruning
23 methods, as investigated in this work, usually consist of three steps: 1) training, 2) pruning, 3) another
24 training step often referred to as *fine-tuning*.

25 Unstructured pruning can significantly reduce the number of parameters of a neural network with
26 only little loss in the accuracy, but the resulting networks often show only a marginal improvement in
27 training and inference time, unless specialized hardware is used (He et al., 2017). In contrast, struc-
28 tured pruning can directly reduce inference time and even training time when applied at initialization
29 (Lee et al., 2018). To exploit these advantages, in this work, we focus on structured pruning.

30 Most sensitivity-based pruning methods such as OBD (e.g., LeCun et al., 1990) or C-OBD (Wang
31 et al., 2019a) evaluate the effect of removing a single weight or structure on the loss of the neural
32 network, while neglecting possible correlations between different structures and within the structures
33 themselves. This can significantly harm the estimation of the sensitivities. We take these correlations
34 into account by applying efficient second-order estimations that not only consider the diagonal terms
35 of the Hessian, but also all off-diagonal terms.

36 Global pruning removes structure by structure from all available structures of a network until a
 37 predefined percentage of pruned structures is reached. Recent examples for global structured pruning
 38 methods are NN Slimming (Liu et al., 2017), C-OBDD and EigenGamage (Wang et al., 2019a). Local
 39 pruning, on the other hand, first subdivides all global structures into subsets (e.g. layers) and removes
 40 a percentage of structures of each subset. Recent examples for local pruning methods are HRank
 41 (Lin et al., 2019), CCP (Peng et al., 2019), FPGM (He et al., 2019) and Variational Pruning (Zhao
 42 et al., 2019). Most local pruning schemes use a predefined layer-wise pruning ratio, which fixes the
 43 percentage of structures removed per layer. While this approach prevents the layers from collapsing,
 44 it also reduces some of the degrees of freedom, since some layers may be less important than others.

45 Our main goal in this work is to devise a simple and efficient second-order pruning method, which
 46 considers all global correlations for structured sensitivity pruning. In addition, we want to highlight
 47 the benefits that such methods may have over other structured global and local pruning schemes.

48 Our contributions are as follows:

- 49 • We introduce two novel saliency-based methods for second-order structured pruning (SOSP),
 50 which consider all correlations across structures and layers. We benchmark our SOSP
 51 methods against a variety of state-of-the-art pruning methods on several networks and
 52 datasets and achieve comparable or better results at low computational costs.
- 53 • We show that our pruning methods can also be applied at initialization almost matching the
 54 performance of pruning after training and significantly reducing the cost of network training.
- 55 • We exploit the structure of the pruning masks found by our SOSP methods to remove
 56 architectural bottlenecks, which further improves the performance of the pruned networks.
 57 In this work, we consider layers with disproportionately low pruning ratios architectural
 58 bottlenecks.

59 PyTorch code implementing our method is attached in the Supplementary Material and we will
 60 publish the code upon acceptance of this manuscript.

61 2 SOSP: Second-order structured pruning

62 A neural network (NN) maps an input $x \in \mathbb{R}^d$ to an output $f_\theta(x) \in \mathbb{R}^D$, where $\theta \in \mathbb{R}^P$ are its
 63 P parameters. NN training proceeds, after random initialization $\theta = \theta_0$ of the weights, by mini-
 64 batch stochastic gradient descent on the empirical loss $\mathcal{L}(\theta) := \frac{1}{N} \sum_{n=1}^N \ell(f_\theta(x_n), y_n)$, given the
 65 training dataset $\{(x_1, y_1), \dots, (x_N, y_N)\}$. In the classification case, $y \in \{1, \dots, D\}$ is a discrete
 66 ground-truth label and $\ell(f_\theta(x), y) := -\log \sigma(f_\theta(x))_y$ the cross-entropy loss, with $\sigma : \mathbb{R}^D \rightarrow \mathbb{R}^D$
 67 the softmax-function. For regression, $y \in \mathbb{R}^D$ and $\ell(f_\theta(x), y) = \frac{1}{2} \|f_\theta(x) - y\|^2$ is the squared loss.

68 Structured pruning aims to remove weights or rather entire structures from a NN f_θ with parameters
 69 θ . A structure can be a filter (channel) in a convolutional layer, a neuron in a fully-connected layer, or
 70 an entire layer in a parallel architecture. We assume the NN in question has been segmented into
 71 S structures $s = 1, \dots, S$, which can potentially be pruned. We define the notation $\theta_s \in \mathbb{R}^P$ as the
 72 vector whose only nonzero components are those weights from θ that belong to structure s .¹ Then, a
 73 *pruning mask* is a set $M = \{s_1, \dots, s_m\}$ of structures. Applying a mask M to a NN f_θ means to
 74 consider the NN with parameter vector $\theta_{\setminus M} := \theta - \sum_{s \in M} \theta_s$.²

75 We now develop our pruning methods that incorporate global correlations into their saliency as-
 76 sessment by efficiently including the second-order loss terms. The first method (SOSP-I) admits a
 77 direct interpretation in terms of individual loss sensitivities, while the second (SOSP-H) remains very
 78 efficient for the largest networks due to its Hessian-vector product approximation.

79 The basic idea behind both our pruning methods is to select the pruning mask M so as to (approx-
 80 imately) minimize the *joint* effect on the network loss

$$\lambda(M) := |\mathcal{L}(\theta) - \mathcal{L}(\theta_{\setminus M})|$$

¹We require that each weight is assigned to at most one structure. In practice, we associate with each structure those weights that go *into* the structure, rather than those that leave it.

²Note, a mask M learned on θ can be applied to a different $\theta' \neq \theta$ on the same NN architecture.

81 of removing all structures in M , subject to a constraint on the overall pruning ratio. To circumvent
 82 this exponentially large search space, we approximate the loss up to second order, so that

$$\lambda_2(M) = \left| \sum_{s \in M} \theta_s^T \frac{d\mathcal{L}(\theta)}{d\theta} - \frac{1}{2} \sum_{s, s' \in M} \theta_s^T \frac{d^2\mathcal{L}(\theta)}{d\theta d\theta^T} \theta_{s'} \right| \quad (1)$$

83 collapses to single-structure contributions plus pairwise correlations; note that the latter include
 84 interactions among the weights within a single $s = s'$, which can be sizeable for large structures.

85 The first-order terms $\lambda_1(s) := \theta_s \cdot d\mathcal{L}(\theta)/d\theta \in \mathbb{R}^P$ in (1) are efficient to evaluate by computing the
 86 gradient $d\mathcal{L}(\theta)/d\theta \in \mathbb{R}^P$ once and then a (sparse) dot product for every s . In contrast to these first-
 87 order terms, the network Hessian $H(\theta) := d^2\mathcal{L}(\theta)/d\theta^2 \in \mathbb{R}^{P \times P}$ in (1) is prohibitively expensive to
 88 compute or store in full. We therefore propose two different schemes to efficiently overcome this
 89 obstacle. Each scheme entails its own way to select the pruning mask. We name the full methods
 90 SOSP-I (individual sensitivities) and SOSP-H (Hessian-vector product).

91 2.1 SOSP-I: Saliency from individual sensitivities

92 SOSP-I approximates each individual term $\theta_s^T H(\theta) \theta_{s'}$ in (1) efficiently, as we will show in Eq. (6).
 93 We can therefore consider a modification of Eq. (1) in which the sensitivity is judged by considering
 94 all single and pairwise sensitivities individually:

$$\lambda_2^I(M) = \sum_{s \in M} \left| \theta_s^T \frac{d\mathcal{L}(\theta)}{d\theta} \right| + \frac{1}{2} \sum_{s, s' \in M} |\theta_s^T H(\theta) \theta_{s'}|. \quad (2)$$

95 To avoid cancellations between signed contributions, we take absolute values because this measures
 96 the strengths of the individual sensitivities $\lambda_1(s)$ and pairwise correlations $\theta_s^T H(\theta) \theta_{s'}$. While
 97 objectives other than λ_2^I are equally possible in the method, including λ_2 and modifications with the
 98 absolute value not pulled in all the way, we found empirically that λ_2^I performs best overall.

99 Then, SOSP-I iteratively selects the structures to prune, based on the objective (2): Starting from
 100 an empty pruning mask $M = \{\}$, we iteratively add to M the structure $s \notin M$ that minimizes the
 101 overall sensitivity $\lambda_2^I(M \cup \{s\})$. In practice, the algorithm pre-computes the matrix $Q \in \mathbb{R}^{S \times S}$,

$$Q_{s, s'} := \frac{1}{2} |\theta_s^T H(\theta) \theta_{s'}| + \left| \theta_s^T \frac{d\mathcal{L}(\theta)}{d\theta} \right| \cdot \delta_{s=s'}, \quad (3)$$

102 and selects at each iteration a structure $s \notin M$ to prune by

$$\arg \min_{s \notin M} \lambda_2^I(M \cup s) - \lambda_2^I(M) = \arg \min_{s \notin M} \left(Q_{s, s} + 2 \sum_{s' \in M} Q_{s, s'} \right), \quad (4)$$

103 terminating at the desired pruning ratio.

104 In order to approximate the Hessian terms $\theta_s^T H(\theta) \theta_{s'}$ efficiently, we omit from $H(\theta) =$
 105 $\frac{1}{N} \sum_n \nabla_{\theta}^2 \ell(f_{\theta}(x_n), y_n)$ those terms that involve the expensive second-order derivatives $\nabla_{\theta}^2 f_{\theta}(x_n)$
 106 of the NN outputs, while including second-order couplings due to ℓ . This is equivalent to approx-
 107 imating $H(\theta) \approx H(f_{\theta}^{lin}) := \frac{1}{N} \sum_n \nabla_{\theta}^2 \ell(f_{\theta}^{lin}(x_n), y_n)$ for the linearized $f_{\theta'}(x) \approx f_{\theta'}^{lin}(x) :=$
 108 $f_{\theta}(x) + \phi(x) \cdot (\theta' - \theta)$ with $\phi(x) := \nabla_{\theta} f_{\theta}(x) \in \mathbb{R}^{D \times P}$, which is well motivated by the NTK limit
 109 (Jacot et al., 2018) for large NNs at both initialization and after training. The terms in the sum become

$$\nabla_{\theta}^2 \ell(f_{\theta}^{lin}(x_n), y_n) = \phi(x_n)^T R_n \phi(x_n), \quad (5)$$

110 where $R_n \in \mathbb{R}^{D \times D}$ is diagonal for squared loss, and has an additional rank-1 contribution for
 111 cross-entropy (see App. B). Similar Hessian approximations were employed before in NNs (Hassibi
 112 et al., 1993; Wang et al., 2019a; Peng et al., 2019) and also in the Gauss-Newton optimization method
 113 (Fletcher, 2013). Our final approximation is to use a random subsample of $N' < N$ data points:

$$\theta_s^T H(\theta) \theta_{s'} \approx \frac{1}{N'} \sum_{n=1}^{N'} (\phi(x_n) \theta_s)^T R_n (\phi(x_n) \theta_{s'}). \quad (6)$$

114 In practice, one pre-computes all (sparse) products $\phi(x_n) \theta_s \in \mathbb{R}^D$ starting from the efficient gradient
 115 $\phi(x_n)$, before aggregating a batch onto the terms $\theta_s^T H(\theta) \theta_{s'}$. Eq. (6) also has an interpretation as
 116 output correlations between certain network modifications, without using derivatives (App. C).

117 **2.2 SOSp-H: Saliency from Hessian-vector product**

118 SOSp-H treats the second-order terms in (1) in a way that is motivated by the limit of large pruning
 119 ratios: At high pruning ratios, the sum $\sum_{s' \in M} \theta_{s'}$ in (1) can be approximated by $\sum_{s'=1}^S \theta_{s'} =:$
 120 θ_{struct} (this equals θ if every NN weight belongs to some structure s). The second-order term
 121 $\sum_{s, s' \in M} \theta_s^T H(\theta) \theta_{s'} \approx (\sum_{s \in M} \theta_s^T) (H(\theta) \theta_{struct})$ thus becomes tractable since the Hessian-vector
 122 product $H(\theta) \theta_{struct}$ is efficiently computable by a variant of the backpropagation algorithm. To
 123 account for each structure s and for the first- and second-order contributions separately, as above, we
 124 place absolute value signs in (1) so as to arrive at the final objective $\lambda_2^H(M) := \sum_{s \in M} \lambda_2^H(s)$ with

$$\lambda_2^H(s) := \left| \theta_s^T \frac{d\mathcal{L}(\theta)}{d\theta} \right| + \frac{1}{2} \left| \theta_s^T (H(\theta) \theta_{struct}) \right|. \quad (7)$$

125 The last term measures the correlations between one structure s and all other prunable structures,
 126 although some of these may cancel unlike for SOSp-I. To minimize $\lambda_2^H(M)$, SOSp-H starts from an
 127 empty pruning mask $M = \{\}$, and successively adds to M a structure $s \notin M$ with smallest $\lambda_2^H(s)$.

128 Unlike the Gauss-Newton approximation in SOSp-I, SOSp-H uses the exact Hessian $H(\theta)$, but can
 129 therefore not account for individual absolute s - s' -correlations, see Eq. (7) vs. (2). Both methods
 130 reduce to the same first-order pruning method when neglecting the second order (i.e. $H(\theta) := 0$).

131 **2.3 Computational complexity**

132 We detail here the computational complexities of our methods (for the experimental evaluation see
 133 Sec. 3.2). The approximation of Q in (3) requires complexity $O(N'D(F+P)) = O(N'DF)$ for
 134 computing all $\phi(x_n) \theta_s$, where $F \geq P$ denotes the cost of one forward pass through the network
 135 ($F \approx P$ for fully-connected NNs), plus $O(N'DS^2)$ for the sum in (6). This is tractable for modern
 136 NNs, while including the exact $H(\theta)$ would have complexity at least $O(N'DSF)$. Once Q has been
 137 computed, the selection procedure based on (4) has overall complexity $O(S^3)$, which is feasible for
 138 most modern convolutional NNs (Sec. 3.2). The total complexity of the SOSp-I method is thus

$$O(N'DF) + O(N'DS^2) + O(S^3). \quad (8)$$

139 SOSp-H has computational complexity $O(N'DF)$ to compute the sensitivities (7), which is com-
 140 parable to computing the sensitivities Q in SOSp-I when the number of structures is low ($S^2 \lesssim F$).
 141 Together with the sorting of the saliency values $\lambda_2^H(s)$, the overall complexity of SOSp-H is thus

$$O(N'DF) + O(S \log(S)). \quad (9)$$

142 Due to its weak dependency on S , in practice, SOSp-H efficiently scales to large modern networks
 143 and may even be used for *unstructured* second-order pruning, where $S = P$.

144 Both of our methods scale much better than naively including all off-diagonal Hessian terms, which
 145 is intractable for modern NNs due to its $O(N'DSF)$ scaling. Since SOSp-I builds on individual
 146 absolute sensitivities and the established Gauss-Newton approximation, we use SOSp-I in the
 147 following in particular to validate the more efficient SOSp-H method.

148 **3 Results**

149 To evaluate our methods, we train and prune VGGs (Simonyan & Zisserman, 2014), ResNets (He
 150 et al., 2016), and DenseNets (Huang et al., 2017) on the Cifar10/100 (Krizhevsky et al., 2009) and
 151 ImageNet (Deng et al., 2009) datasets. Stochastic gradient descent with an initial learning rate of 0.1,
 152 a momentum of 0.9 and weight decay of 10^{-4} is used to train these networks. For ResNet-32/56 and
 153 VGG-Net on Cifar10/100, we use a batch size of 128, train for 200 epochs and reduce the learning rate
 154 by a factor of 10 after 120 and 160 epochs. To fine-tune the network after pruning, we exactly repeat
 155 this learning rate schedule. For DenseNet-40 on Cifar10/100, we train for 300 epochs and reduce
 156 the learning rate after 150 and 225 epochs. For ResNets on ImageNet, we use a batch size of 256,
 157 train for 128 epochs and use a cosine learning rate decay. For all networks, we prune feature maps
 158 (i.e. channels) from all layers except the last fully-connected layer; for ResNets, we also exclude the
 159 downsampling-path from pruning. We approximate the Hessians by a subsample of size $N' = 1000$
 160 (see Sec. 2.1). We report the best or average final test accuracy over 3 trials if not noted otherwise.
 161 The experiments were run on an internal cluster with Nvidia Tesla V100 GPUs. Reproducing the
 162 results presented in this paper would take about 60 days of GPU run-time.

Table 1: Comparison of SOSP to other global pruning methods for high pruning ratios. The comparison for moderate pruning ratios is deferred to the appendix (see App. A.1). We tuned our pruning ratios to similar values as reported by the referred methods. To ensure identical implementations of the network models in PyTorch, reference numbers are taken from Wang et al. (2019a) and (Mingjie & Zhuang, 2018). In accordance with all referred methods, we report the mean and standard deviation of the best accuracies observed during fine-tuning. For final accuracies after fine-tuning see App. A.3. * denotes the baseline model. Both SOSP methods perform either on par or outperform the competing global pruning methods.

Dataset	Cifar10			Cifar100		
	Test acc. (%)	Reduct. in weights (%)	Reduct. in MACs (%)	Test acc. (%)	Reduct. in weights (%)	Reduct. in MACs (%)
VGG-Net*	94.18	-	-	73.45	-	-
NN Slimming	85.01	97.85	97.89	58.69	97.76	94.09
NN Slim. + L_1	91.99	97.93	86.00	57.07	97.59	93.86
C-OBDD	92.34 ± 0.18	97.68 ± 0.02	77.39 ± 0.36	58.07 ± 0.60	97.97 ± 0.04	77.55 ± 0.25
EigenDamage	92.29 ± 0.21	97.15 ± 0.04	86.51 ± 0.26	65.18 ± 0.10	97.31 ± 0.01	88.63 ± 0.12
SOSP-I (ours)	92.62 ± 0.14	97.79 ± 0.02	83.52 ± 0.29	64.20 ± 0.23	97.83 ± 0.04	87.02 ± 0.20
SOSP-H (ours)	92.71 ± 0.19	97.81 ± 0.01	86.32 ± 0.29	64.59 ± 0.35	97.81 ± 0.01	86.32 ± 0.29
ResNet-32*	95.30	-	-	76.8	-	-
C-OBDD	91.75 ± 0.42	97.30 ± 0.06	93.50 ± 0.37	59.52 ± 0.24	97.74 ± 0.08	94.88 ± 0.08
EigenDamage	93.05 ± 0.23	96.05 ± 0.03	94.74 ± 0.02	65.72 ± 0.04	95.21 ± 0.04	94.62 ± 0.06
SOSP-I (ours)	92.43 ± 0.09	95.47 ± 0.33	94.07 ± 0.66	67.36 ± 0.46	92.69 ± 0.07	95.63 ± 0.13
SOSP-H (ours)	92.23 ± 0.12	95.26 ± 0.10	94.45 ± 0.40	68.42 ± 0.21	94.08 ± 0.21	95.06 ± 0.14
DenseNet-40*	94.58	-	-	74.11	-	-
NN Slim. + L_1	94.22	54.21	-	73.19	54.21	-
SOSP-I (ours)	94.21 ± 0.04	47.00 ± 0.10	36.35 ± 0.12	73.05 ± 0.11	45.22 ± 0.10	42.05 ± 1.16
SOSP-H (ours)	94.23 ± 0.05	49.39 ± 0.65	38.86 ± 0.70	73.05 ± 0.24	48.58 ± 0.22	42.05 ± 0.35

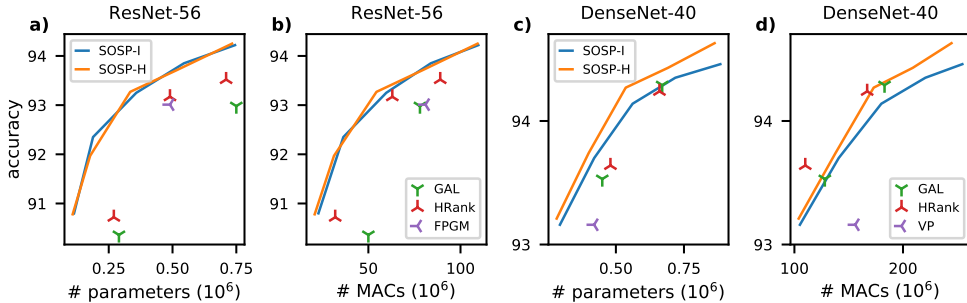


Figure 1: Comparison of SOSP to local, i.e. layer-wise, pruning methods on Cifar10. The best final test accuracy is plotted over the effective number of model parameters (a, c) and MACs (b, d). A tabular representation as well as statistics across trials are shown in App. A.4. SOSP outperforms all competing layer-wise pruning methods, especially over the number of effective parameters.

163 **3.1 Comparison to Literature**

164 Here we benchmark the performance of our SOSP methods against existing pruning algorithms on
 165 different datasets and networks. First, we compare against other recent global pruning methods,
 166 then against local structured pruning methods, i.e. with pre-specified layer-wise pruning rates. In all
 167 comparisons we report the achieved test accuracy, the number of parameters of the pruned network
 168 and the MACs (often referred to as FLOPs). To facilitate direct comparisons, we report test accuracies
 169 in the same way as the competing methods (e.g. best trial or average over trials), but additionally
 170 report mean and standard deviation of the test error for our models in App. A. Our count of the
 171 network parameters and MACs is based on the actual pruned network architecture (cf. App. D), even
 172 though our saliency measure associates with each structure only the weights *into* this structure.

173 We first compare our SOSP methods with global pruning methods on VGG-Net, ResNet-32 and
 174 DenseNet-40. We use the same variants and implementations of these networks as used by Neural
 175 Network Slimming (NN Slimming; Liu et al., 2017) as well as EigenDamage and C-OBDD (Wang
 176 et al., 2019a), e.g. capping the layer-wise ratio of removed structures at 95% for VGGs to prevent
 177 layer collapse and increasing the width of ResNet-32 by a factor 4. C-OBDD is a structured variant of

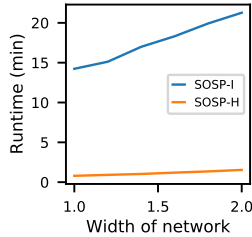


Figure 2: Runtime to calculate the pruning masks for ResNet-56 on Cifar10 over the width of the network for SOSP-I and SOSP-H. We vary the width of the network by increasing the width of each layer by a multiplicative factor.

Table 2: Best final test accuracies and pruning ratios (PR) across 2 trials on ImageNet. For comparison to CCP, we also provide their alternative MAC count (for details, see App. D). * denotes SOSP with kernel scaling (see main text). SOSP outperforms all three competing methods.

Model	Top-1% (Gap)	Parameters (PR)	MACs (PR)	Alt. MACs (PR)
ResNet-18	69.76 (0.0)	11.7M (0%)	1.82B (0%)	1.82B (0%)
SOSP (ours)	69.63 (0.13)	7.12M (39%)	1.37B (24%)	1.31B (28%)
FPGM	68.41 (1.87)	7.10M (39%)	1.06B (41%)	-
SOSP (ours)	68.78 (0.98)	6.42M (45%)	1.29B (29%)	1.20B (34%)
ResNet-50	76.15 (0.0)	25.5M (0%)	3.85B (0%)	3.85B (0%)
SOSP (ours)	76.56(-0.41)	19.9M(22%)	3.06B(21%)	2.72B(29%)
SOSP* (ours)	76.60 (-0.45)	17.9M (30%)	2.79B (28%)	2.47B (36%)
HRank	74.98 (1.17)	16.2M (36%)	2.30B (44%)	-
FPGM	75.59 (0.56)	15.9M (37%)	2.36B (42%)	-
SOSP (ours)	75.85 (0.30)	15.4M (40%)	2.44B (27%)	1.97B (49%)
HRank	71.98 (4.17)	13.8M (46%)	1.55B (62%)	-
CCP	75.21 (0.94)	-	-	1.77B (54%)
SOSP* (ours)	75.21 (0.94)	13.0M (49%)	2.13B (45%)	1.68B (56%)
SOSP (ours)	74.39 (1.76)	11.8M (54%)	1.89B (51%)	1.38B (64%)
SOSP* (ours)	73.38 (2.77)	9.9M (61%)	1.58B (59%)	1.10B (72%)

178 the original OBD algorithm (Hassibi et al., 1993), which neglects all cross-structure correlations that,
 179 in contrast, SOSP takes into account. The results over three trials for high pruning ratios are shown in
 180 Tab. 1 and for moderate pruning ratios in App. A.1. To enable the comparison to NN Slimming on
 181 an already pretrained VGG, we included the results of NN Slimming applied to a baseline network
 182 obtained without modifications to its initial training, i.e. without L_1 -regularization on the batch-
 183 normalization parameters. For moderate pruning rates, all pruning schemes approximately retain the
 184 baseline performance for VGG-Net and ResNet-32 on Cifar10 and VGG-Net on Cifar100 (see Tab.
 185 3). The only exception is the accuracy for C-OBD applied to VGG-Net on Cifar100, which drops by
 186 approximately 1%. For ResNet-32 on Cifar100 the accuracy after pruning is approximately 1% lower
 187 than the baseline, for all pruning schemes. In the regime of larger pruning ratios of approximately 97%,
 188 SOSP and EigenDamage significantly outperform NN Slimming and C-OBD. SOSP performs on par
 189 with EigenDamage, except for ResNet-32 on Cifar100, where SOSP outperforms EigenDamage by
 190 almost 3%. This result indicates that SOSP outperforms all other methods especially in the regime of
 191 baseline networks that have relatively few parameters already and on difficult datasets most relevant
 192 for applications. For DenseNet-40, we achieve similar results compared to NN Slimming. However,
 193 note that NN Slimming requires the modification of network pretraining.

194 Next, we compare our SOSP methods against four recently published local, i.e. layer-wise, pruning
 195 algorithms: FPGM He et al. (2019), GAL (Lin et al., 2019), CCP (Peng et al., 2019), Variational
 196 Pruning (VP; Zhao et al., 2019) and HRank (Lin et al., 2020). For ResNet-56, our SOSP methods
 197 outperform all other methods across all pruning ratios (see Fig. 1a and b). For DenseNet-40, SOSP
 198 achieves better accuracies when compared over parameters (Fig. 1c) and is on par with the best other
 199 methods over MACs (Fig. 1d). The reason for this discrepancy is probably that the SOSP objective is
 200 agnostic to the number of MACs (image size) in each individual structure.

201 3.2 Scalability and application to large-scale datasets

202 Before going to large datasets, we compare the scalability of our methods SOSP-I and SOSP-H.
 203 As the preceding section shows, both methods perform basically on par with each other in terms
 204 of accuracy. This confirms that SOSP-H is not degraded by the approximations leading to the
 205 efficient Hessian-vector product, or is helped by use of the exact Hessian. In terms of efficiency,
 206 however, SOSP-H shows clear advantages compared to SOSP-I, for which the algorithm to select the
 207 structures to be pruned scales with $O(S^3)$ (see Sec. 2.3) potentially dominating the overall runtime
 208 for large-scale networks. Measurements of the actual runtimes show that already for medium-sized
 209 networks SOSP-H is more efficient than SOSP-I (Fig. 2). Since SOSP-I becomes impractical for
 210 large-scale networks, for the ImageNet dataset we will only evaluate SOSP-H and refer to it as SOSP.

211 On ImageNet, we compare our results to literature for ResNet-18 and ResNet-50, see Tab. 2. Because
 212 SOSP assesses the sensitivity of each structure independently of the contributed MACs, it has a
 213 bias towards pruning small-scale structures. This tendency is strongest for ResNet-50, due to its
 214 1×1 -convolutions. Since these 1×1 -convolutions tend to contribute disproportionately little to the

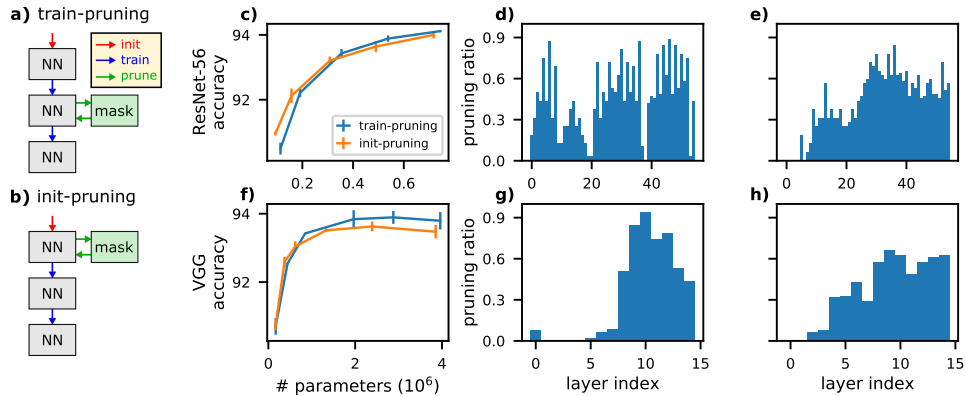


Figure 3: Comparison between pruning after training and at initialization on Cifar10. Both pruning schemes, train-pruning (a) and init-pruning (b), train the network for the same overall number of epochs, but generate and apply the pruning masks at different point in times. The average and standard deviation of the test accuracy across 3 trials is plotted against the number of model parameters for ResNet-56 (top row; c) and VGG (bottom row; f). For a single trial, in which overall 50% of the structures are pruned, we visualize the pruning masks of train-pruning and init-pruning by showing the layer-wise pruning ratios in (d, g) and (e, h), respectively.

215 overall number of MACs, we devised a scaled variant of SOSP, which divides the saliency of every
 216 structure by the kernel size (e.g. 1, 3 or 7). Compared to the vanilla SOSP, the scaled variant of SOSP
 217 is able to remove larger percentages of MACs with similar drops in accuracy (see Tab. 2).

218 For both networks SOSP outperforms HRank and FPGM, especially when considering the main
 219 objective of SOSP, which is to reduce the number of parameters or structures. Since CCP uses a
 220 different way of calculating the MACs, which leads to consistently higher pruning ratios, we added
 221 an alternative MAC count to enable a fair comparison (for details, see App. D). Since HRank and
 222 FPGM do not mention their MAC counting convention, we assume they use the same convention
 223 as we do. Taking this into account, our scaled SOSP variant is able to prune more MACs than CCP,
 224 while having the same final accuracy.

225 3.3 Pruning at Initialization

226 Traditionally, pruning methods are applied to pretrained networks, as also done in the previous
 227 sections, but recently there has been growing attention on pruning at initialization following the
 228 works of Lee et al. (2018) and Frankle & Carbin (2018). Since SOSP employs the absolute value of
 229 the sensitivities, it can also be applied to a randomly initialized network without any modifications.
 230 Thus, SOSP can also be seen as an efficient second-order generalization of SNIP (Lee et al., 2018;
 231 van Amersfoort et al., 2020). While EigenDamage can in principle be modified and applied to a
 232 randomly initialized network, NN Slimming can not be applied at initialization.

233 Usually pruning at initialization leads to worse accuracies than pruning after training (Liu et al.,
 234 2018). However, pruning an already trained network is often followed by fine-tuning effectively
 235 training the network twice (Fig. 3a). For comparability with pruning at initialization, we unify the
 236 overall training schedule between these two settings and consequently apply two training cycles after
 237 pruning the randomly initialized network (Fig. 3b; for further discussion, see App. A.5).

238 We observe that applying SOSP at initialization performs almost equally well than applying SOSP
 239 after training (see ResNet-56 and VGG in Fig. 3c and f, respectively). In conclusion, applying SOSP
 240 at initialization can significantly reduce the time and resources required for network training with no
 241 or only minor degradation in accuracy.

242 3.4 Identifying & Removing Architectural Bottlenecks

243 Even though the previous section highlighted that applying SOSP after training and at initialization
 244 results in comparable accuracies, the pruning masks differ between these two scenarios (compare Fig.

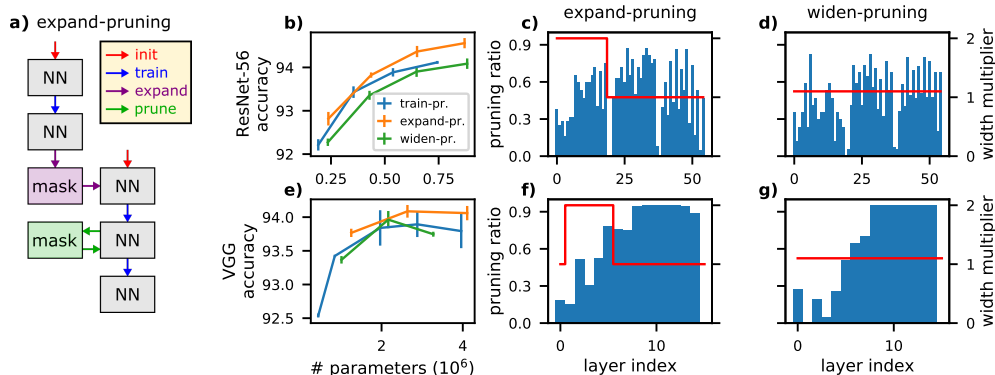


Figure 4: We remove architectural bottlenecks found by SOSF using the expand-pruning scheme (a) on Cifar10. The width of blocks and layers with low pruning ratios in the train-pruning scheme (Fig. 3d and g) are expanded by a width multiplier of 2 (c, f). As a baseline, we uniformly expand all layers in the network by a factor 1.1 (d, g). The layer-wise pruning ratios of the enlarged network models are shown as bar plots in (c, d, f, g). The average and standard deviation of the test accuracy across 3 trials are shown over the number of model parameters (b, e). Note that the full ResNet-56 and VGG models have $0.86 \cdot 10^6$ and $20 \cdot 10^6$ parameters, respectively.

245 3d and g to e and h, resp.). Despite this difference a common feature of all masks is that some layers
 246 are barely pruned or not pruned at all while others are pruned by up to 80%. This could indicate
 247 towards architectural bottlenecks. We consider a layer an architectural bottleneck if the respective
 248 layer has a considerably lower pruning ratio compared to the other layers. The low pruning ratio of
 249 bottleneck layers indicates that the substructures (e.g. filters) have a high sensitivity. Thus, widening
 250 these layers could improve the overall performance and allow for even smaller models with higher
 251 accuracies.

252 To utilize this insight we devise a procedure that we call *expand-pruning*. The idea is to first calculate
 253 the pruning ratios of a trained network and then to identify architectural bottlenecks, i.e. the layers
 254 with the lowest pruning ratios. Next, we widen these layers by a factor of two, which has been
 255 empirically shown to work well. Finally, we randomly initialize, train, prune, and fine-tune the
 256 expanded network (for a schematic, see Fig. 4a). As a naive baseline that we call the *widen-pruning*
 257 procedure, we widen all layers in the network with a constant factor instead of widening specific
 258 layers. We choose the constant factor such that the overall number of parameters matches that of the
 259 expand-prune procedure (see width multiplier in Fig. 4).

260 We evaluate the expand-pruning procedure for ResNet-56 and VGG on Cifar10, for which we expand
 261 the least pruned of the three main building blocks and the five least pruned layers, respectively (e.g.,
 262 see *width multipliers* in Fig. 4c and f selected on the basis of the pruning masks shown in Fig. 3d and
 263 g, resp.). Note that for ResNet-56 a more fine-grained removal of bottlenecks, e.g. on layer level,
 264 is not possible without changing the overall ResNet architecture. In summary, a selective removal
 265 of bottlenecks results in smaller network models with higher accuracy than pruning the vanilla
 266 network or unselectively increasing the network size (compare expand-pruning to train-pruning and
 267 widen-pruning in Fig. 4b and e). While in principle any global pruning method could be used for
 268 the expand-pruning procedure, SOSF is especially suited since it does not require to modify the
 269 network architecture like EigenDamage and can also be applied at initialization unlike NN Slimming,
 270 allowing for a similar expand scheme directly at initialization (see App. A.6).

271 4 Discussion

272 In this work we have demonstrated the effectiveness and scalability of our second-order structured
 273 pruning algorithms (SOSF). While both algorithms perform similarly well, SOSF-H is more easily
 274 scalable to large scale networks and datasets. We highlighted two major features of our method. Firstly,
 275 SOSF can be applied at initialization with only minor degradation in accuracy, which drastically
 276 reduces the required time and resources for training. Secondly, we showed that the pruning masks

277 found by SOSP can be used to systematically detect and remove architectural bottlenecks, further
278 improving the performance of pruned networks.

279 Compared to other global pruning methods, SOSP captures correlations between structures by a sim-
280 ple, effective and scalable algorithm that neither requires to modify the training nor the architecture of
281 the to be pruned network model and achieves comparable or better accuracies on benchmark datasets.
282 The C-OBDD algorithm (Wang et al., 2019a) is a structured generalization of the original unstructured
283 OBD algorithm (LeCun et al., 1990). In contrast to OBD, C-OBDD accounts for correlations within
284 each structure, but does not capture correlations between different structures within and across layers.
285 We show that considering these *global* correlations consistently improve the performance, especially
286 for large pruning ratios (Tab. 1). This observation is further confirmed by an ablation study in which
287 we neglect all cross-structure correlations significantly decreasing the performance under otherwise
288 identical experimental settings (App. A.2). The objective of EigenDamage (Wang et al., 2019a) to
289 include second order correlations is similar to ours, but the approaches are significantly different.
290 EigenDamage uses the Fisher-approximation, which is similar to the Gauss-Newton approximation
291 used for SOSP-I, and then, in addition to further approximations, apply low rank approximations
292 that require the substitution of each layer by a bottleneck-block structure. Our SOSP method is
293 simpler, easier to implement and does not require to modify the network architecture, but nevertheless
294 performs on par with EigenDamage. The approach of NN Slimming (Liu et al., 2017) is more
295 heuristic than SOSP and is easy to implement. However, networks need to be pretrained with L_1
296 regularization on the batch-normalization parameters, otherwise the performance is severely harmed
297 (Tab. 1 and 3). SOSP does not require any modifications to the network training and can be applied
298 to any pretrained network. A recent variant of NN Slimming was developed by Zhuang et al. (2020)
299 who optimize their hyperparameters to reduce the number of MACs. Using the number MACs as an
300 objective for SOSP is left for future studies.

301 In addition to the above comparison to other global pruning methods, we also compared our methods
302 to simpler local pruning methods that keep the pruning ratios constant for each layer and, consequently,
303 scale well to large-scale datasets. The pruning method closest to our SOSP-I method is the one by
304 Peng et al. (2019). While both works consider second-order correlations between structures, theirs is
305 based on a pruning objective different from our absolute sensitivities in λ_2^l and considers only intra-
306 layer correlations. Furthermore, they employ an auxiliary classifier with a hyperparameter, yielding
307 accuracy improvements that are difficult to disentangle from the effect of second-order pruning. Going
308 beyond a constant pruning ratio for each layer Su et al. (2020) discovered, for ResNets, that pruning
309 at initialization seems to preferentially prune initial layers and thus proposed a pruning scheme based
310 on a “keep-ratio” per layer which increases with the depth of the network. Our experiments confirm
311 some of the findings of Su et al. (2020), but we also show that the specific network architectures found
312 by pruning can drastically vary between different networks and especially between initialization and
313 after training (histograms in Fig. 3). While all local pruning methods specify pruning ratios for each
314 layer, our method performs automatic selection across layers (histograms in Fig. 3).

315 This automatic selection allows us to identify and remove architectural bottlenecks. However, our
316 global pruning method has a bias towards pruning small structures, absent from local pruning methods,
317 as the size of structures is usually identical within layers. We propose a simple solution by scaling
318 each structure by the inverse of the kernel size which helps to remove some of the bias. Alternatively,
319 to better reflect the computational costs in real-world applications, each structure could also be
320 normalized by the number of its required MACs (like done by van Amersfoort et al., 2020).

321 Recently, unstructured (Lee et al., 2018; Wang et al., 2019b; Tanaka et al., 2020) and structured
322 (van Amersfoort et al., 2020; Hayou et al., 2021) pruning schemes that are applicable to networks
323 at initialization were proposed. While these methods fail to achieve similar accuracies compared to
324 pruning after training, our SOSP method in the init-pruning setting achieves accuracies comparable
325 to pruning after training.

326 In accordance with Elsken et al. (2019), our results suggest that pruning can be used to optimize the
327 architectural hyperparameters of established networks (Liu et al., 2018) or super-graphs (Noy et al.,
328 2020). Instead of formulating this optimization as a pruning process, we envision our second-order
329 sensitivity analysis to be a valuable tool to identify and remove bottlenecks to find good neural
330 architectures more quickly. For example, whenever a building block of the neural network cannot
331 be compressed, this building block may be considered a bottleneck of the architecture and could be
332 inflated to improve the overall trade-off between accuracy and computational cost.

333 References

- 334 Blalock, D., Ortiz, J. J. G., Frankle, J., and Gutttag, J. What is the state of neural network pruning?
335 *arXiv preprint arXiv:2003.03033*, 2020.
- 336 Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. Imagenet: A large-scale hierarchical
337 image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255.
338 Ieee, 2009.
- 339 Elsken, T., Metzen, J. H., and Hutter, F. Neural architecture search: A survey. *Journal of Machine*
340 *Learning Research*, 20(55):1–21, 2019.
- 341 Fletcher, R. *Practical methods of optimization*. John Wiley & Sons, 2013.
- 342 Frankle, J. and Carbin, M. The lottery ticket hypothesis: Finding sparse, trainable neural networks.
343 In *International Conference on Learning Representations*, 2018.
- 344 Han, S., Pool, J., Tran, J., and Dally, W. Learning both weights and connections for efficient neural
345 network. *Advances in neural information processing systems*, 28:1135–1143, 2015.
- 346 Hassibi, B., Stork, D. G., and Wolff, G. J. Optimal brain surgeon and general network pruning. In
347 *IEEE international conference on neural networks*, pp. 293–299. IEEE, 1993.
- 348 Hayou, S., Ton, J.-F., Doucet, A., and Teh, Y. W. Robust pruning at initialization. In *International*
349 *Conference on Learning Representations*, 2021. URL https://openreview.net/forum?id=vXj_ucZQ4hA.
- 351 He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings*
352 *of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- 353 He, Y., Zhang, X., and Sun, J. Channel pruning for accelerating very deep neural networks. In
354 *Proceedings of the IEEE International Conference on Computer Vision*, pp. 1389–1397, 2017.
- 355 He, Y., Liu, P., Wang, Z., Hu, Z., and Yang, Y. Filter pruning via geometric median for deep
356 convolutional neural networks acceleration. In *Proceedings of the IEEE/CVF Conference on*
357 *Computer Vision and Pattern Recognition*, pp. 4340–4349, 2019.
- 358 Huang, G., Liu, Z., Van Der Maaten, L., and Weinberger, K. Q. Densely connected convolutional
359 networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.
360 4700–4708, 2017.
- 361 Jacot, A., Gabriel, F., and Hongler, C. Neural tangent kernel: convergence and generalization in
362 neural networks. In *Proceedings of the 32nd International Conference on Neural Information*
363 *Processing Systems*, pp. 8580–8589, 2018.
- 364 Krizhevsky, A., Hinton, G., et al. Learning multiple layers of features from tiny images. 2009.
- 365 LeCun, Y., Denker, J. S., and Solla, S. A. Optimal brain damage. In *Advances in neural information*
366 *processing systems*, pp. 598–605, 1990.
- 367 Lee, N., Ajanthan, T., and Torr, P. Snip: Single-shot network pruning based on connection sensitivity.
368 In *International Conference on Learning Representations*, 2018.
- 369 Lin, M., Ji, R., Wang, Y., Zhang, Y., Zhang, B., Tian, Y., and Shao, L. Hrank: Filter pruning using
370 high-rank feature map. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*
371 *Pattern Recognition*, pp. 1529–1538, 2020.
- 372 Lin, S., Ji, R., Yan, C., Zhang, B., Cao, L., Ye, Q., Huang, F., and Doermann, D. Towards optimal
373 structured cnn pruning via generative adversarial learning. In *Proceedings of the IEEE/CVF*
374 *Conference on Computer Vision and Pattern Recognition*, pp. 2790–2799, 2019.
- 375 Liu, Z., Li, J., Shen, Z., Huang, G., Yan, S., and Zhang, C. Learning efficient convolutional networks
376 through network slimming. In *Proceedings of the IEEE International Conference on Computer*
377 *Vision (ICCV)*, Oct 2017.

- 378 Liu, Z., Sun, M., Zhou, T., Huang, G., and Darrell, T. Rethinking the value of network pruning. In
379 *International Conference on Learning Representations*, 2018.
- 380 Mingjie, S. and Zhuang, L. Network slimming. [https://github.com/Eric-mingjie/
381 network-slimming](https://github.com/Eric-mingjie/network-slimming), 2018.
- 382 Noy, A., Nayman, N., Ridnik, T., Zamir, N., Doveh, S., Friedman, I., Giryas, R., and Zelnik, L.
383 ASAP: Architecture search, anneal and prune. In Chiappa, S. and Calandra, R. (eds.), *Proceedings
384 of the Twenty Third International Conference on Artificial Intelligence and Statistics*, volume 108
385 of *Proceedings of Machine Learning Research*, pp. 493–503, 2020.
- 386 Peng, H., Wu, J., Chen, S., and Huang, J. Collaborative channel pruning for deep networks. In
387 *International Conference on Machine Learning*, pp. 5113–5122. PMLR, 2019.
- 388 Reed, R. Pruning algorithms—a survey. *IEEE transactions on Neural Networks*, 4(5):740–747, 1993.
- 389 Simonyan, K. and Zisserman, A. Very deep convolutional networks for large-scale image recognition.
390 *arXiv preprint arXiv:1409.1556*, 2014.
- 391 Su, J., Chen, Y., Cai, T., Wu, T., Gao, R., Wang, L., and Lee, J. D. Sanity-checking pruning methods:
392 Random tickets can win the jackpot. In *Advances in Neural Information Processing Systems*, 2020.
- 393 Tanaka, H., Kunin, D., Yamins, D. L., and Ganguli, S. Pruning neural networks without any data
394 by iteratively conserving synaptic flow. *Advances in Neural Information Processing Systems*, 33,
395 2020.
- 396 Tang, Y., Wang, Y., Xu, Y., Tao, D., Xu, C., Xu, C., and Xu, C. Scop: Scientific control for reliable
397 neural network pruning. *Advances in Neural Information Processing Systems*, 2020.
- 398 van Amersfoort, J., Alizadeh, M., Farquhar, S., Lane, N., and Gal, Y. Single shot structured pruning
399 before training. *arXiv preprint arXiv:2007.00389*, 2020.
- 400 Wang, C., Grosse, R., Fidler, S., and Zhang, G. Eigendamage: Structured pruning in the kronecker-
401 factored eigenbasis. In *International Conference on Machine Learning*, pp. 6566–6575, 2019a.
- 402 Wang, C., Zhang, G., and Grosse, R. Picking winning tickets before training by preserving gradient
403 flow. In *International Conference on Learning Representations*, 2019b.
- 404 Zhao, C., Ni, B., Zhang, J., Zhao, Q., Zhang, W., and Tian, Q. Variational convolutional neural
405 network pruning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
406 Recognition*, pp. 2780–2789, 2019.
- 407 Zhuang, T., Zhang, Z., Huang, Y., Zeng, X., Shuang, K., and Li, X. Neuron-level structured pruning
408 using polarization regularizer. *Advances in Neural Information Processing Systems*, 33, 2020.

409 Checklist

- 410 1. For all authors...
- 411 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
412 contributions and scope? [Yes]
- 413 (b) Did you describe the limitations of your work? [Yes]
- 414 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 415 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
416 them? [Yes]
- 417 2. If you are including theoretical results...
- 418 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 419 (b) Did you include complete proofs of all theoretical results? [Yes]
- 420 3. If you ran experiments...
- 421 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
422 mental results (either in the supplemental material or as a URL)? [Yes]

- 423 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
424 were chosen)? [Yes]
- 425 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
426 ments multiple times)? [Yes]
- 427 (d) Did you include the total amount of compute and the type of resources used (e.g., type
428 of GPUs, internal cluster, or cloud provider)? [Yes]
- 429 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 430 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 431 (b) Did you mention the license of the assets? [Yes]
- 432 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 433 (d) Did you discuss whether and how consent was obtained from people whose data you're
434 using/curating? [N/A]
- 435 (e) Did you discuss whether the data you are using/curating contains personally identifiable
436 information or offensive content? [N/A]
- 437 5. If you used crowdsourcing or conducted research with human subjects...
- 438 (a) Did you include the full text of instructions given to participants and screenshots, if
439 applicable? [N/A]
- 440 (b) Did you describe any potential participant risks, with links to Institutional Review
441 Board (IRB) approvals, if applicable? [N/A]
- 442 (c) Did you include the estimated hourly wage paid to participants and the total amount
443 spent on participant compensation? [N/A]