

Selectivity Drives Productivity: Efficient Dataset Pruning for Enhanced Transfer Learning





Website



Yihua Zhang^{1*}, Yimeng Zhang^{1*}, Aochuan Chen^{1*}, Jinghan Jia¹, Jiancheng Liu¹, Gaowen Liu², Mingyi Hong³, Shiyu Chang⁴, Sijia Liu¹ ¹Michigan State University, ²Cisco Research, ³University of Minnesota, Twin Cities, ⁴University of California, Santa Barbara



Dataset Pruning for Transfer Learning

- ♦ Not all sources classes are necessary or beneficial^[1]
- Some source data could be harmful to downstream performance.
- **Removing** specific source classes can **improve** transfer learning.

✤ Conventional dataset pruning lacks effectiveness on transfer learning.

• Conventional SOTA DP methods do NOT yield significant improvement over random pruning on transfer learning.



Figure 1. Transfer learning accuracy of existing DP methods on ImageNet at different pruning ratios, where ResNet-101 is the source model.

An Overview of Our Proposal



Research Question

Existing method: brute-force based, effective but not affordable.

(Q) How can we extend DP to transfer learning with high efficiency, broad applicability, and lossless or even improved target performance?

Label Mapping for Supervised Pretraining

- Rationale behind the design: source data similar to downstream data intend to contribute more during the transfer process.
- Dataset selection as a voting process: each downstream training data can vote for its most similar/relevant source training class.
- Label mapping through a pretrained small surrogate model.

Pretrained



Training Data

👍 Source Class 1 🔿 Source Class 3 Δ Source Class 2 Source Class 4

0

00

Ο

00

 \circ

Feature Space of a Surrogate Model

0 0

000

000

18 8

Classification

Inference Result

- Feature Mapping for Self-Supervised Pretraining
- Data cluster as the basic pruning unit when labels are unavailable.
 - Training Data of a Downstream Task
- O Data Cluster Center of the Pretrained Dataset
- → Voting based on the Distance in the Feature Space
- Retrained source Data Cluster
- Pruned source Data Cluster

[1] S. Jain et al. "A data-based perspective on transfer learning." CVPR 2023.

[2] S. Kim et al. "Coreset sampling from open-set for fine-grained self-supervised learning " CVPR 2023

Experiment Result Highlights

NEURAL INFORMATION PROCESSING SYSTEMS



Paper

Source Data Pruning Ratio Figure 2. Unstructured pruning trajectory given by test accuracy (%) vs. sparsity (%) on various (dataset, model) pairs. The performance of dense model and that of the best winning ticket are marked using dashed lines in each plot. The solid line and shaded area of each pruning method represent the mean and variance of test accuracies over 3 trials.

Dataset	OxfordPets				SUN397				Flowers102						
Pruning Ratio	0%	50%	60%	70%	80%	0%	50%	60%	70%	80%	0%	50%	60%	70%	80%
Random Moderate GraNd FM (ours)	69.26	62.32 63.37 64.42 69.92	61.27 62.45 63.34 69.99	59.09 63.31 61.14 70.29	53.75 57.42 56.42 70.21	47.36	45.63 45.73 45.72 48.46	45.08 45.14 45.58 48.58	43.54 44.23 45.24 47.90	39.81 40.82 41.72 46.00	85.17	82.23 82.45 82.85 85.22	82.60 81.45 82.44 85.42	81.03 81.69 82.14 84.37	80.02 81.32 81.73 84.61

Table 1. The downstream performance with different source data pruning ratios in the SSL pretraining setting. A randomly initialized RN-101 is self-supervised pretrained using MOCO V2 on each full/pruned source dataset and finetuned on the downstream task through LP.

Pruning Ratio	0%	20%	40%	60%	80%
Time	5.4 (4.6	3.5	2.4	1.3
Consumption (h)		15%↓)	(35%↓)	(56%↓)	(76%↓)
Table 2. Time consu The reported time co LM/FM dataset prun	imption o	f LM/FM	I to obtain	the pretra	ined model
	nsumption	n covers :	surrogate n	nodel (RN-	18) training
	ing, and s	ource mo	odel pretrai	ning (RN-1	01).
	SUN39	7		DTD	

		SUN.	397		DID					
Method	Pruning	Acc	.(%)	Time	Pruning	Acc.(%)		Time		
	Ratio	LP	FF	(h)	Ratio	LP	FF	(h)		
DENSE	N/A	51.45	54.21	5.4	N/A	65.91	67.21	5.4		
DENSE-ADV	N/A	52.97	55.67	13.7	N/A	67.23	68.92	13.7		
LM	70	50.95	54.28	1.9	50	66.25	67.22	2.9		
LM-ADV	70	52.07	55.49	4.2	50	67.02	68.54	6.7		

Table 4. Downstream performance of models pretrained on full/pruned source dataset (Dense/LM) w/wo adversarial pretraining (Adv).



Source Data Pruning Ratio

igure 3. Source dataset pruning trajectory given the downstream task OxfordPets using different surrogate models.

Acknowledgement. The work of Y. Zhang, Y. Zhang, A. Chen, J. Jia, J Liu, and S. Liu was supported in part by the Cisco Research Award, the NSF Grant IIS-2207052, and the ARO Award W911NF2310343. The work of S. Chang was supported by the Cisco Research Award and the NSF Grant IIS-2207052. The work of M. Hong was supported by NSF grants CIF-1910385 and EPCN-2311007