

SCAN: Bootstrapping Contrastive Pre-training for Data Efficiency

Yangyang Guo, Mohan Kankanhalli

National University of Singapore, Singapore

guoyang.eric@gmail.com, mohan@comp.nus.edu.sg

Abstract

While contrastive pre-training is widely employed, its data efficiency problem has remained relatively under-explored thus far. Existing methods often rely on static coreset selection algorithms to pre-identify important data for training. However, this static nature renders them unable to dynamically track the data utility throughout pre-training, leading to subpar pre-trained models. To address this challenge, our paper introduces a novel dynamic bootstrapping dataset pruning method. It involves pruning data preparation followed by dataset mutation operations, both of which undergo iterative and dynamic updates. We apply this method to two prevalent contrastive pre-training frameworks: **CLIP** and **MoCo**, representing vision-language and vision-centric domains, respectively. In particular, we individually pre-train seven CLIP models on two large-scale image-text pair datasets, and two MoCo models on the ImageNet dataset, resulting in a total of 16 pre-trained models. With a data pruning rate of 30-35% across all 16 models, our method exhibits only marginal performance degradation (less than 1% on average) compared to corresponding models trained on the full dataset counterparts across various downstream datasets, and also surpasses several baselines with a large performance margin. Additionally, the byproduct from our method, i.e., coresets derived from the original datasets after pre-training, also demonstrates significant superiority in terms of downstream performance over other static coreset selection approaches. Code is available at <https://github.com/guoyang9/SCAN>.

1. Introduction

Large models are heavily data-driven, particularly in the realm of pre-training [10, 11, 51]. This paradigm has been widely underpinned by the scaling law [22, 27, 30], which suggests that more data often leads to reduced generalization errors. However, using large quantities of data frequently results in a notable increase in the carbon footprint. Addressing this pressing issue requires substantial efforts to optimize the data training efficiency.

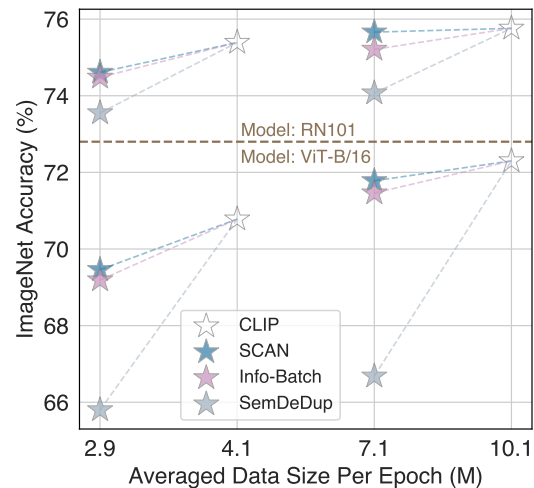


Figure 1. Interplay between averaged training data size per epoch and the downstream performance of CLIP, our method SCAN, and two SoTA baselines. **All models are trained for the same number of epochs.** Two CLIP models are pre-trained separately on two datasets containing 4.1M and 10.1M samples, respectively. The three data-efficient methods, including SCAN, utilize 30% less data, reducing the average data size to 2.9M and 7.1M, respectively.

This paper delves into the data efficiency problem for contrastive pre-training. Despite the pervasiveness of contrastive pre-training across both vision-centric [10, 11] and vision-language [29, 51] domains, nevertheless, the data efficiency issue has received scant attention in the existing literature. We attribute the reason for this fact to two challenges. **I-** Absence of reliable labels for self-supervised learning objectives. Unlike in supervised learning, where explicit labels aid in class prediction, self-supervised learning in contrastive pre-training operates without such guidance, making it unable to estimate the class probability of data samples such as EL2N [49]. **II-** Extensive data scale due to easy accessibility, e.g., the interleaved image-text data from the web [51]. Current datasets usually comprise millions [7, 56] or even billions of samples [54]. It is thus computationally very intensive for methods employing gradients [49] or second derivative (Influence Functions) [31] to evaluate individual data utility. Recent approaches in

the vision-language area have resorted to coreset selection algorithms [44] for a reduced pre-training dataset beforehand [1, 42, 43, 65, 68]. The crux of these methods lies in the semantic match/duplication that is quantified by some proxy metrics like CLIP matching score [51]. Consequently, a subset, namely a coreset, of the original dataset is filtered for pre-training from scratch.

Our motivation for this work is inspired by the advancement in dynamic sparse training (DST) [47, 73, 74], which dynamically prunes less influential learnable weights from models. Compared to its static sparse training counterparts [34, 61], DST demonstrates notable strengths in performance, robustness, and model compression without the need for over-parameterization [39, 47]. Intriguingly, we recognize that recent coreset selection algorithms predominantly adhere to the static approach, akin to the fixed weight masks employed in static sparse networks [34]. As a result, we argue that these coreset-based dataset pruning methods are subject to similar limitations as previous static sparse training ones, albeit with a shift in application scope from learnable weights to individual data samples. Given this context, approaching the dataset pruning challenge can be easily decomposed into two sub-problems: **1)- metric identification** and **2)- pruning strategy design**. Specifically, the proxy metric should meet several conditions: dynamic adaptability, quick obtainability (with minimal additional cost), and reflecting the learning status of each sample. Regarding the pruning strategy, we introduce a novel data bootstrapping algorithm named **SCAN**. Instead of employing a consistent pruning ratio throughout training, our SCAN approach identifies and eliminates data from less important subsets in a bootstrapping manner. These two operations from our SCAN method are performed iteratively for stable pre-training.

We validate the effectiveness of the proposed method with widely used contrastive pre-training frameworks in both vision-language (**CLIP** [28, 51]) and vision-centric (**MoCo** [10, 11]) domains. The pre-training datasets for CLIP include CC3M [56], MSCOCO [38], SBU-Captions [48], and CC12M [7], forming two groups of datasets with different scales. On the other hand, we pre-train MoCo models using the ImageNet Dataset [13]. Moreover, we employ various downstream datasets, including ImageNet [13], CIFAR-10, and CIFAR-100 [32], along with out-of-distribution datasets such as ImageNet-R [25] and ImageNet V2 [52]. Our evaluation protocols encompass full fine-tuning, linear probing, and zero-shot testing on ImageNet. Within the CLIP framework, we evaluate seven models covering ResNet [23], ViT [14], and Swin Transformer [40]. As for MoCo, due to resource constraints, we use two popular ViT models [62] in the experiments. Our experimental results, partially depicted in Fig. 1, exhibit that SCAN achieves a significant trade-off between training data size and downstream model performance as compared to

several baselines [1, 50, 71].

To the best of our knowledge, we are the first to comprehensively study the data efficiency problem within the context of contrastive pre-training. Our work not only introduces an effective bootstrapping approach but also is able to produce a static coreset (a smaller dataset) that outperforms other static coreset selection methods [1, 71] by a large performance margin on diverse downstream image classification datasets. These contributions enable our work to hold a positive promise for the efficient utilization of data in contrastive pre-training, thereby potentially reducing both computational overhead and carbon footprint.

2. Related Work

Dataset Pruning and Distillation represent two common approaches to enhancing dataset efficiency during training. The former aims to synthesize a smaller dataset that achieves test performance similar to a full dataset when using the same model [12, 15, 55, 60]. Recent advancements in this area have leveraged techniques such as mutual information [55], frequency domain transformation [57], and multi-stage generation [12] to craft datasets that exhibit enhanced performance. Two notable limitations are inherent in these methods: I) The generalization capability is significantly constrained due to the reliance on distillation from a specific dataset and model, *e.g.*, the use of ResNet [23]. II) The dataset sizes utilized for comparison are frequently confined to small-scale datasets such as CIFAR [32] and Tiny-ImageNet [33]. In contrast, dataset pruning is employed to directly filter a smaller subset from the original dataset [35, 59]. Typically, previous methods involve initially learning an indication score, which serves as a basis for identifying and subsequently removing data samples falling below or above a certain threshold. For image classification tasks, prevailing methods often utilize gradient [49], loss value [50], and second derivative [31] as the indicator. More recently, efforts have emerged focusing on pruning vision-language pre-trained datasets [18, 36, 37, 65–67, 69, 70]. The key idea is to construct a coreset by identifying the semantic mismatch/duplication, a process facilitated by often using a pre-trained CLIP model [4, 51].

Contrastive Pre-Training has garnered wide attention as a technique for pre-training versatile models applicable to a range of downstream tasks [20]. Its fundamental principle involves bringing the embeddings of positive pairs closer while simultaneously pushing away negative ones. Traditional approaches within the vision-centric domain often build positive samples by leveraging alternative views of the anchor sample, as exemplified by approaches like SimCLR [8, 9] and MoCo [10, 11]. Benefiting from the advancement of transformer architectures [14, 64], recent endeavors have shifted towards patch masking followed by subsequent recovery [3, 6, 24]. Contrastive learning has achieved notable

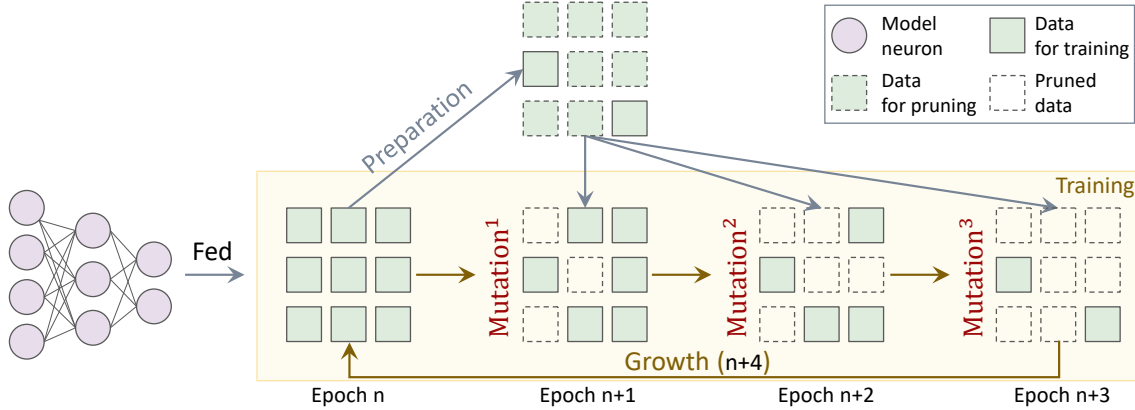


Figure 2. Overall pipeline of the proposed SCAN method. We begin by identifying a substantial portion of data samples as pruning candidates. Subsequently, a subset of these candidates is employed for pruning based on varying mutation ratios that are gradually increased (bootstrapping). After growing back to the original full dataset, the above two operations are iterated for another round.

success in vision-language pre-training as well [29, 51]. The alignment of modalities has propelled significant advancements in downstream multi-modal tasks, including visual question answering [2, 75] and cross-modal retrieval [5, 38]. Our study specifically targets the data efficiency challenge within CLIP and MoCo, which serve as prominent representatives of vision-language and vision-centric domains, respectively.

Dynamic Sparse Training (DST). Unlike earlier static sparse training methodologies [34, 61], DST learns a sparse neural network by pruning weights and growing them back throughout the training process. The weight importance is typically quantified using metrics such as magnitude [17, 45], gradients [73], or sensitivity [46]. The demonstrated superiority of DST over its static counterpart inspires us to design a dynamic approach for dataset pruning, especially considering that recent coreset-based methods predominantly adhere to fixed pruning strategies. Furthermore, the well-developed DST methods provide additional hints for devising our dataset pruning strategy.

3. Method

3.1. Background of Contrastive Pre-Training

Contrastive pre-training necessitates the utilization of both positive and negative pairs of samples, be it alternative views of an image [10, 11] or combinations of image and text [29, 51]. Its objective is to bring positive pairs closer in the embedding space while pushing negative ones away. At its core is the InfoNCE loss [63], defined as,

$$\mathcal{L}_{f \rightarrow g} = -\frac{1}{|\mathcal{D}_t|} \sum_{i=1}^{|\mathcal{D}_t|} \log \frac{\exp(f(I_i)^T g(T_i)/\tau)}{\sum_{j=1}^{|\mathcal{D}_t|} \exp(f(I_i)^T g(T_j)/\tau)}, \quad (1)$$

where τ is a learnable temperature, I_i and T_j respectively denote the sampled image and text from the batched examples

\mathcal{D}_t ; while f and g represent the image and text encoders, respectively. Likewise, we can obtain the loss from the other direction $\mathcal{L}_{g \rightarrow f}$ as well. The overall training loss is then computed as the mean of $\mathcal{L}_{f \rightarrow g} + \mathcal{L}_{g \rightarrow f}$. Without loss of generality, we take the contrastive learning utilized in CLIP as an example [51]. The application in other approaches such as MoCo [11] can be straightforwardly extrapolated.

Motivation. Contrastive pre-training often demands large-scale data to learn a versatile model. For instance, the original CLIP pre-training uses 400M image-text pairs sourced from the web [51], while recent studies push these limits to datasets containing billions of samples [54, 72]. Accordingly, the introduced footprint and storage cost present significant challenges for researchers. To address this issue, we propose a novel bootStrapping ContrAstive Pre-traiNing method (SCAN), to dynamically, efficiently, and effectively leverage smaller dataset for pre-training¹.

3.2. SCAN Overview

Our proposed SCAN method involves a two-step operation. Specifically, our **first** step entails identifying a proxy metric that is dynamically adaptable, easily obtainable, and capable of tracing the learning status of each sample. We abandon the use of gradients as done in related domains [49], given the substantially increased compute overhead incurred for individual samples². Instead, we opt for the loss value as the reliable indicator as it meets the above conditions. Under this context, we disentangle the loss in Eqn. 1 to obtain a loss set $\tilde{\mathcal{L}}_{f \rightarrow g} = \{\tilde{\mathcal{L}}_{f \rightarrow g}^{(1)}, \tilde{\mathcal{L}}_{f \rightarrow g}^{(2)}, \dots, \tilde{\mathcal{L}}_{f \rightarrow g}^{(|\mathcal{D}_t|)}\}$, with each element corresponding to the loss value of one data sample.

The **second** core step is to determine the pruning strategy. Addressing sparsity within the dataset pruning study

¹The name itself reflects our method’s capability to *scan* all data samples, identifying those that should be eliminated from further pre-training.

²The data size in contrastive pre-training is notably larger than in other related domains.

presents a substantial challenge. To approach this, we propose a pruned data preparation and then a dataset mutation approach. Specifically, in the pruned data preparation stage, candidate data samples that are potentially less important are selected. Thereafter, in the dataset mutation stage, samples are gradually bootstrapped for pruning across epochs. These two stages iterate through several rounds until the completion of pre-training.

3.3. Bootstrapping Dataset Pruning

Unlike conventional approaches that use the full dataset for pre-training, we vary the training data size for each epoch. As illustrated in Fig. 2, an example case utilizes 1.0, 6%, 4%, and 2% of the full dataset for four consecutive training epochs, resulting in an average dataset pruning ratio of $\sim 40\%$. We provide more details regarding the algorithm in the supplementary material.

3.3.1. Pruning Data Preparation

We opt to utilize the loss values of in-batch samples rather than those from the entire dataset for candidate selection. This decision is based on two facts: 1) Comparing InfoNCE losses across batches holds less significance compared to supervised learning, as the instance loss varies drastically with respect to the randomly selected batched samples. 2) Saving the losses for the entire dataset incurs more computational storage compared to in-batch ones.

Additionally, existing coreset selection approaches [26, 53] from the vision-language domain typically focus only on pruning *ill-matched* samples. These are characterized by image-text pairs indicating less semantic alignment, where, for example, the text inadequately describes the content of the paired image. However, we posit that beyond these ill-matched samples, there also exist samples that are *redundant*. These redundant samples are effectively memorized during the early stages of training and are less likely to be forgotten with further training iterations [19]. In view of this, we can safely eliminate these redundant data as training progresses.

To operationalize the above ideas, given a pruning ratio ρ , we separately identify the *ill-matched* and *redundant* set of data using:

$$\begin{cases} \mathcal{D}_t^{red} &= \mathcal{D}_{t:i}, \quad i \in \prec_{\rho} \bar{\mathcal{L}}_{f \rightarrow g}, \\ \mathcal{D}_t^{ill} &= \mathcal{D}_{t:j}, \quad j \in \succ_{\rho} \bar{\mathcal{L}}_{f \rightarrow g}, \end{cases} \quad (2)$$

where $\prec_{\rho} \bar{\mathcal{L}}_{f \rightarrow g}$ denotes the indices of the ρ smallest values of $\bar{\mathcal{L}}_{f \rightarrow g}$ (the loss set before loss summation and backpropagation). We then obtain the *redundant* subset \mathcal{D}_t^{red} by selecting data samples according to these indices from the original full in-batch set \mathcal{D}_t . This approach is driven by the intuition that small losses often denote effective data memorization by the given model. On the other hand, we can also identify the *ill-matched* subject \mathcal{D}_t^{ill} from the ρ largest loss values using $\succ_{\rho} \bar{\mathcal{L}}_{f \rightarrow g}$. This is because large losses are usually associated

with small cosine similarities [26], indicating a *poor match* between image and text pairs. The final candidate subset is formed by the union of these two: $\mathcal{D}'_t = \mathcal{D}_t^{red} \cup \mathcal{D}_t^{ill}$.

Thereafter, we merge the subset intersection from $\mathcal{L}_{f \rightarrow g}$ and $\mathcal{L}_{g \rightarrow f}$. In total, we obtain $2\rho|\mathcal{D}_t|$ candidates for the current batched samples, which amounts to twice the size of the expected pruned data, as will be explained in the next section. At last, we iterate through all training data to have the final candidate pruning subset \mathcal{D}' .

Preparation Warm-up. It is intuitive that the model’s learning capability may exhibit instability during the early training iterations. To mitigate this issue, we design a warm-up strategy [21], wherein the full dataset is utilized for training during the first several epochs. We track the average epoch-wise loss to determine the optimal timing for initiating pruning. Specifically, we calculate the difference between the loss from the previous epoch \mathcal{L}'_{pre} and the current epoch \mathcal{L}'_{cur} , and compare it against a pre-defined threshold value T_{td} . If the condition $(\mathcal{L}'_{pre} - \mathcal{L}'_{cur})/(\mathcal{L}'_{pre} + \epsilon) \geq T_{td}$ holds true, where ϵ is infinitesimal and is introduced to prevent overflow, it indicates relative stability in the pre-training process. Consequently, we can then start the *pruning data preparation* from this pre-training epoch.

3.3.2. Dataset Mutation

Rather than employing a static pruning ratio consistently throughout training, we advocate for a bootstrapping dataset mutation approach. The benefits of this methodology are shown in Sect. 4.2. Additionally, we refrain from performing the **pruning data preparation** solely once, as further training iterations may alter the matching and redundancy characteristics of data samples.

Specifically, we regenerate the candidate pruning data from scratch every $(\tau_{cos} + 1)$ epochs as detailed in Sect. 3.3.1. Subsequently, we adapt the cosine annealing strategy [41] to determine the current pruning ratio ρ_{cur} using:

$$\rho_{cur} = \frac{1}{2} (1 + \cos((\tau_{cos} - (\tau_{cur} \bmod (\tau_{cos} + 1))) \frac{\pi}{\tau_{cos}})). \quad (3)$$

It can be easily seen that the pruning ratio ρ_{cur} gradually and periodically increases with larger training epoch τ_{cur} . We can then randomly select $\rho_{cur}|\mathcal{D}'|$ samples from \mathcal{D}' for pruning – \mathcal{D}'_{ρ} . During this training epoch, batched instances \mathcal{D}_t are sampled from the reduced dataset $\mathcal{D} \setminus \mathcal{D}'_{\rho}$, which are then employed for pre-training using contrastive learning (Eqn. 1). This bootstrapping process provides us with robust estimates regarding data samples and enables us to refrain from making strict assumptions about the underlying distribution of the data [16, 58].

In each iterative round, where ρ_{cur} increases from 0 to 2ρ (where 0 corresponds to the pruned data preparation epoch), the average pruning ratio remains fixed at ρ as predefined. Finally, we grow back to the original full dataset for another round of pruning data preparation and mutation (Fig. 2).

Table 1. Performance comparison of CLIP models on the **CC12M+** pre-trained datasets. All methods utilize **30% fewer pre-trained data samples** than CLIP. Consequently, they also require approximately **30% less pre-training time**. The best results (excluding the original CLIP model) are highlighted in **bold**. A dash (-) indicates the collapse of pre-training, resulting in impaired evaluation of downstream tasks.

Architecture	Method	IN Zero-Shot		CIFAR10	CIFAR100	IN	IN-V2	IN-R
		Top-1	Top-5					
RN101	CLIP	18.78	41.14	95.96	82.13	75.76	64.31	40.57
	Random	14.05	30.60	95.02	78.34	73.99	60.27	36.13
	SemDeDup [1]	13.26	29.70	95.07	78.77	74.24	62.16	37.65
	D-Pruning [71]	12.59	28.62	94.94	78.89	74.07	61.30	37.07
	Info-Batch [50]	21.60	41.11	96.04	81.60	75.21	63.27	39.34
	SCAN	23.10	47.52	96.08	82.28	75.66	63.75	40.10
ViT-B/32	CLIP	24.62	49.10	95.62	82.11	63.40	49.97	31.09
	Random	09.12	21.09	90.13	69.98	51.99	41.01	20.08
	SemDeDup [1]	06.47	16.71	90.83	70.03	52.21	39.75	20.99
	D-Pruning [71]	06.27	15.88	90.11	69.69	51.72	38.66	20.42
	Info-Batch [50]	-	-	-	-	-	-	-
	SCAN	26.12	50.67	95.41	81.16	61.55	48.73	29.23
ViT-B/16	CLIP	23.43	47.71	96.76	84.25	72.30	59.39	33.50
	Random	14.45	32.41	94.35	76.45	67.09	54.38	27.07
	SemDeDup [1]	11.58	26.01	94.18	76.71	67.22	53.78	27.19
	D-Pruning [71]	10.00	23.72	93.82	75.96	66.68	53.13	26.23
	Info-Batch [50]	22.12	42.30	96.03	81.69	71.46	56.35	31.13
	SCAN	24.71	49.12	96.13	83.71	71.78	58.58	32.45

3.4. Time-Efficiency of SCAN

In fact, our proposed SCAN method introduces negligible additional pre-training time in comparison to each base model. The potentially increased time involves three major steps: metric selection, pruned data preparation, and pruned data retrieval. Since we directly utilize individual loss values, the metric selection step does not impose an additional time overhead. Second, we obtain the candidate pruned data from in-batch samples, typically containing only a few thousand data points, thereby enabling a fast sorting process. Last, retrieving from millions of pruned data sets also leads to negligible costs, as evidenced by our empirical observations.

More importantly, we maintain consistency in the number of training epochs for our SCAN model, aligning it with the epochs utilized by each respective base model. This approach ensures time efficiency within our proposed method, as demonstrated in Sect. 4.5.

4. Experiments

To ensure strict consistency across all data-efficient methods in terms of the number of training samples, we keep the total number of training epochs for our method and the baselines the same as that of the base model, such as CLIP. For clarity, we primarily report the average training data size per epoch in the following discussion.

4.1. Experimental Settings

Pre-Training Datasets. For CLIP models, we utilized two versions of pre-training datasets to examine the data-size scaling law as well. We employed the OpenCLIP repository [28] to conduct pre-training for all models (including CLIP) from scratch, ensuring a fair comparison between our proposed method and the baselines. Specifically, the smaller dataset, denoted as **CC3M+**, comprises CC3M [56], SBU-Captions [48], and MSCOCO [38], totaling 4.1 million

Table 2. Performance comparison of MoCo models. All methods utilize **35% fewer pre-trained data samples** than CLIP. Consequently, they also require approximately **35% less pre-training time**. The best results (excluding the original MoCo model) are highlighted in **bold**.

Arc	Method	ImageNet		CIFAR-100	
		Top-1	Top-5	Top-1	Top-5
ViT-S/16	MoCo [11]	78.48	94.17	86.02	97.85
	Random	75.38	91.18	84.00	95.59
	Info-Batch [50]	77.99	93.45	85.51	97.49
	SCAN	78.58	94.19	86.01	97.53
ViT-B/16	MoCo [11]	79.53	94.49	88.31	98.04
	Random	75.28	91.01	85.99	97.02
	Info-Batch [50]	78.46	94.18	87.69	97.70
	SCAN	79.15	94.33	88.11	97.78

Table 3. Performance comparison of CLIP models on the **CC3M+** pre-trained datasets. All methods utilize **30% fewer pre-trained data samples** than CLIP. Consequently, they also require approximately **30% less pre-training time**. The best results (excluding the original CLIP model) are highlighted in **bold**. A dash (-) indicates the collapse of pre-training, resulting in impaired evaluation of downstream tasks.

Architecture	Method	IN Zero-Shot		CIFAR10	CIFAR100	IN	IN-V2	IN-R
		Top-1	Top-5					
RN101	CLIP	15.72	35.19	96.17	81.78	75.39	63.42	39.69
	Random	12.35	29.03	95.01	78.99	73.73	61.01	36.11
	SemDeDup [1]	12.97	28.90	95.16	79.44	74.08	61.94	36.74
	D-Pruning [71]	12.77	28.03	94.85	78.23	73.55	61.56	36.48
	Info-Batch [50]	16.79	34.38	95.49	80.89	74.48	62.11	38.01
SCAN	15.59	34.77	95.77	81.95	74.61	62.92	38.05	
ViT-B/16	CLIP	18.17	37.62	96.58	82.47	70.78	57.28	30.13
	Random	13.26	31.27	91.62	73.53	50.60	40.55	21.80
	SemDeDup [1]	11.35	25.56	94.36	76.53	66.56	53.18	25.50
	D-Pruning [71]	10.00	23.31	93.46	75.37	65.80	52.50	24.39
	Info-Batch [50]	17.16	39.14	95.98	81.43	69.19	56.00	28.55
SCAN	18.21	37.20	96.00	81.49	69.46	56.04	28.60	
Swin-Base	CLIP	13.98	32.05	95.75	82.19	73.89	61.92	35.38
	Random	12.56	30.12	94.10	78.01	72.55	60.25	31.31
	SemDeDup [1]	-	-	-	-	-	-	-
	D-Pruning [71]	12.44	28.60	94.07	77.86	72.54	60.28	31.41
	Info-Batch [50]	13.82	32.05	95.80	81.25	72.69	59.47	33.13
SCAN	17.50	37.42	95.37	81.35	73.55	61.60	33.55	

image-text pairs. The larger dataset, denoted as **CC12M+**, includes CC12M [7], SBU-Captions [48], and MSCOCO [38], with a total of 10.1 million pairs.

For the pre-training dataset of **MoCo** [10, 11], we adhered to the original implementation and utilized the ImageNet dataset [13].

Downstream Datasets. We conducted extensive downstream fine-tuning experiments across various datasets. Specifically, we utilized datasets such as ImageNet [13], CIFAR-10, CIFAR-100 [32], as well as out-of-distribution datasets including ImageNet V2 [52] and ImageNet-R [25] to validate the downstream performance of **CLIP** pre-trained models. For all these datasets, we explored diverse experimental settings, encompassing zero-shot transfer learning from ImageNet, linear probing, and full fine-tuning. Moreover, we employed both the ImageNet and CIFAR-100 datasets to conduct experiments on **MoCo** models.

Model Architectures. As OpenCLIP [28] offers a variety of model cards, we utilized model architectures from both ResNet [23] and ViT [14]. The model architectures employed for pre-training **CLIP** include RN50, RN101, ViT-S/32, ViT-S/16, ViT-B/32, ViT-B/16, and Swin-Base [40]. Due to resource constraints, we pre-trained **MoCo** using two model architectures: ViT-S/16 and ViT-B/16 [62]. *It is important to note that pre-training MoCo consumes approximately seven times more resources than pre-training a CLIP model. Therefore, we primarily conducted experiments on*

CLIP to assess the effectiveness of the proposed method.

Compared Baselines. Given the absence of common models for addressing the data efficiency problem in contrastive pre-training, we employed four different approaches in this study: Random, SemDeDup [1], D-Pruning [71], and Info-Batch [50]. *It is worth noting that Info-Batch is not specially designed for contrastive pre-training and does not support static coreset selection.* Unless otherwise specified, we utilized a pruning ratio of 30% for CLIP models and 35% for MoCo models.

4.2. Overall Experimental Results

We present the comprehensive results of CLIP in Tables 1 and 3, and the results of MoCo in Table 2. Additional experimental results can be found in the supplementary material. The insights from these tables can be summarized into four key observations:

- Our proposed SCAN method consistently outperforms the four compared baselines, indicating that our approach achieves a superior trade-off between performance and data efficiency compared to existing data-efficient methods.
- In comparison to the base CLIP models, SCAN achieves comparable performance while utilizing fewer data for pre-training. Specifically, our method preserves 99% of the original CLIP model’s performance in most cases, while using 30% less of the original training dataset. Take the IN result of the RN101 architecture in Table 3 as an ex-

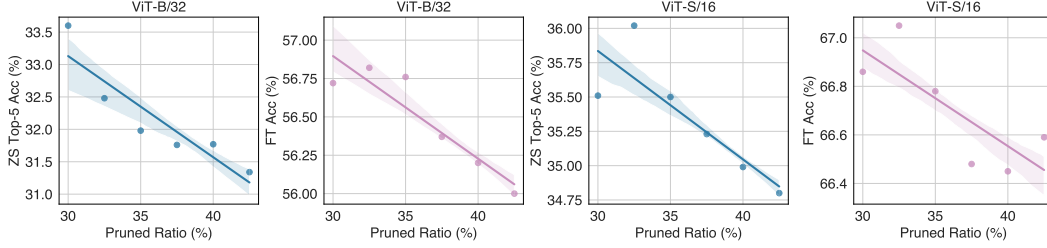


Figure 3. Downstream performance variation of two CLIP models *w.r.t.* different pruning ratios.

Table 4. Coreset selection model results on the ImageNet dataset. The coreset generated by our SCAN (static) method is derived from the intersection of the datasets obtained from the other two models. All models are pre-trained from scratch using each respective coreset.

Method	RN50			ViT-S/16			ViT-B/32		
	ZS@1	ZS@5	FT@1	ZS@1	ZS@5	FT@1	ZS@1	ZS@5	FT@1
SCAN	16.91	35.79	72.91	17.31	35.51	66.86	16.48	33.60	56.64
SemDeDup [1]	11.98	26.30	71.51	09.57	22.00	62.30	07.20	17.50	50.99
D-Pruning [71]	11.72	26.65	71.11	08.60	20.35	61.70	06.51	16.13	50.01
SCAN (static)	16.99	35.20	73.11	16.68	34.31	66.22	13.16	28.46	55.99

Table 5. Coreset overlap ratios. The coreset generated by our SCAN is derived from the intersection of the datasets obtained from the other two/three models.

Model			Overlap Ratio
RN50	ViT-S/16	ViT-B/32	
✓	✓		56.78%
✓		✓	55.77%
	✓	✓	61.73%
✓	✓	✓	47.89%

ample: 74.61 (SCAN) *v.s.* 75.39 (CLIP) - 99% of CLIP result. It is also worth noting that some of our methods even outperform the base CLIP models.

- Dynamic approaches such as Info-Batch and SCAN often outperform static coreset selection methods like SemDeDup and D-Pruning. This verifies the superiority of dynamic pruning approaches.
- A like-to-like comparison between Table 1 and Table 3 reveals that models trained with more data (CC12M+) consistently outperform those trained with less data (CC3M+), thereby verifying the dataset scaling law [27, 30].

4.3. Coreset Results from SCAN

Beyond the dynamic data pruning approach, we also investigated the coreset results obtained from our SCAN method³. To implement this, we first identify the pruned data using two pre-trained models from our method, such as RN50 and ViT-S/16. After that, we obtain the intersection of these two sets and ensure that the overall pruning ratio remains the same as ρ , thereby generating a coreset from the original full dataset.

Downstream Performance. We then performed pre-training for another model from scratch, *e.g.*, ViT-B/32, employing the exact same process as previous static coreset-based approaches [1, 71]. From the results in Table 4, we observe the following: I) Our selected coreset significantly outperforms other existing coreset-based baselines utilizing the exact same settings. II) Our static method achieves very competitive results compared to our dynamic variant. This observation further underscores the potential of our proposed method for identifying a subset that can be directly

and efficiently utilized in future research endeavors, thereby reducing more training overhead.

Coreset Overlaps. We then calculated the intersection over union for each pair and trio of coresets. The results are presented in Table 5. The result reveals that I) The coresets obtained are highly related to the specific model architecture used. II) The coresets from ViT-B/32 and ViT-S/16 have a higher degree of overlap than the other two, as these two models share the same architecture family.

4.4. Ablation Study

Different Pruning Ratios. The performance change corresponding to different pruning ratios ρ is depicted in Fig. 3. Generally, we can observe that the performance tends to degrade with increasing pruning ratios. Further, selecting the optimal pruning ratio also involves a trade-off.

Different Mutation Epochs. The results for different mutation epochs are summarized in Table 6. It can be observed that employing a mutation epoch of three generally yields better results compared to other variants.

Different Pruning Modes. Our SCAN method combines samples categorized as both *redundant* and *ill-matched* for pruning purposes. Additionally, we conducted a separate analysis to evaluate the effectiveness of utilizing either *redundant* (**R.**) or *ill-matched* (**I.**) samples alone, and the results are presented in Table 7. It is evident from the table that employing *ill-matched* samples for pruning yields superior performance compared to using *redundant* samples alone. Moreover, the combination of these two categories results in further performance improvement.

Different Variants of the Same Pruning Ratio. In our experiments, we employed a pruning ratio ρ of 30% for CLIP models and divided it equally between redundant (**R.**)

³Info-Batch lacks such static coreset selection support.

Table 6. Results *w.r.t.* mutation epoch.

Mutation Epochs	ViT-S/32		ViT-S/16	
	ZS@5	FT@1	ZS@5	FT@1
2	28.72	56.38	36.28	66.81
3	33.60	56.72	35.51	67.04
4	28.31	56.25	33.80	66.94

Table 7. Results *w.r.t.* pruning modes.

Pruning Mode		ViT-S/32		ViT-S/16	
w/. R.	w/. I.	ZS@5	FT@1	ZS@5	FT@1
✓	✗	31.48	56.66	33.51	66.85
✗	✓	34.12	56.01	35.19	67.04
✓	✓	33.60	56.72	35.51	67.04

Table 8. Results *w.r.t.* ratio variants.

R. v.s. I. (%)	ViT-S/32		ViT-S/16	
	ZS@5	FT@1	ZS@5	FT@1
(30 : 10)	31.84	56.35	34.84	66.90
(20 : 20)	33.60	56.72	35.51	67.04
(10 : 30)	33.72	56.43	35.49	67.10

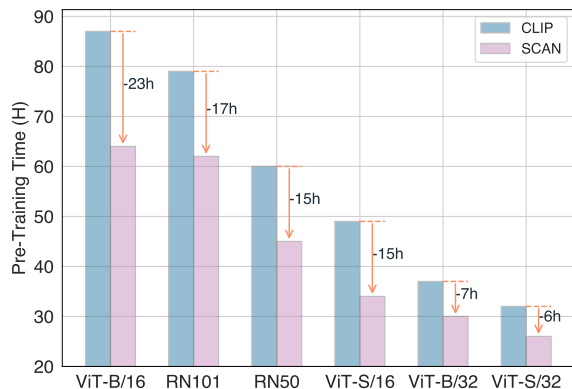


Figure 4. Comparison of pre-training time between the base CLIP model and our SCAN on the CC12M+ dataset. More numbers are provided in the supplementary material.

and ill-matched (**I.**) samples. Furthermore, we investigated other variants, and the outcomes are presented in Table 8. All the ratios are relative to the full dataset. It can be seen from the table that evenly distributing the pruning ratio leads to slightly better results.

4.5. In-depth Analysis

Pre-Training Time Comparison. The primary outcome of our method is the reduction in training time and, consequently, the decrease in carbon footprint. We illustrate the pre-training time in Fig. 4. It can be observed from the figure that our SCAN method contributes to a reduction of approximately 25% to 30% in the original pre-training computational cost (The additional negligible time cost can be attributed to the pruned data preparation and retrieval processes.). This advantage proves especially beneficial when training large models such as RN101 and ViT-B-16.

Visualization of Ill-Matched Samples. We also provide visualizations of some ill-matched samples identified by SCAN. Three examples are shown in Fig. 5. In the second example, an incorrect annotation is evident, as there is no

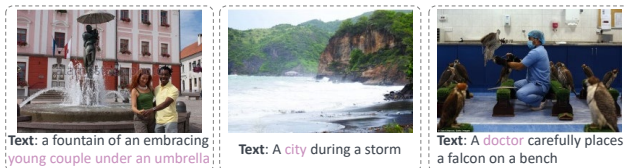


Figure 5. Examples of ill-matched samples identified by our SCAN.

city depicted in the image. Additionally, in the last example, the individual portrayed is identified as a *doctor* rather than a *stockman*.

5. Conclusion

Limitations. We acknowledge two potential limitations of this work. First, akin to other dataset pruning methodologies, our method necessitates the storage of the original large-scale dataset. This may pose storage challenges for researchers with limited computing resources. Second, *we are unable to validate the effectiveness of the proposed method on larger-scale datasets, such as those with hundreds of millions or billions of samples, due to our resource limitations in accessing storage capacity.* Lastly, our method may not seamlessly transfer to the recent popular large language model pre-training. Apart from differences in pre-training objectives, large language models (LLMs) often require only a few training epochs, typically one to three. In contrast, our iterative bootstrapping strategy requires more epochs to converge, the same as each corresponding contrastive pre-training model.

Summary. This work sets an initial effort to comprehensively address the data efficiency challenge in contrastive pre-training. We propose a novel dataset bootstrapping approach, applying it to a range of contrastive pre-training model architectures and evaluating it using various protocols. Our experiments demonstrate that, across all experimental settings, the proposed method achieves a superior balance between downstream model performance and data efficiency compared to both the base models and several existing data efficiency approaches. Additionally, it also helps obtain an effective coreset dataset that significantly outperforms other coreset-based baselines, thereby further reducing time costs and training overhead.

Future Work. In the future, we plan to validate the effectiveness of our method 1) in more domains such as language-centric contrastive pre-training, 2) with larger pre-training datasets for vision-language like LAION-400M [53]. Moreover, our method is orthogonal to other efficiency studies, such as model compression. As such, by integrating strategies from these related domains, we aim to build a more efficient training pipeline framework, thus contributing to substantial reductions in carbon footprints.

References

- [1] Amro Abbas, Kushal Tirumala, Daniel Simig, Surya Ganguli, and Ari S. Morcos. Semdedup: Data-efficient learning at web-scale through semantic deduplication. *CoRR*, 2023. 2, 5, 6, 7
- [2] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C. Lawrence Zitnick, and Devi Parikh. VQA: visual question answering. In *ICCV*, pages 2425–2433. IEEE, 2015. 3
- [3] Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. Beit: BERT pre-training of image transformers. In *ICLR*. OpenReview.net, 2022. 2
- [4] Romain Beaumont. Clip retrieval: Easily compute clip embeddings and build a clip retrieval system with them, 2022. 2
- [5] Kevin W. Bowyer and Patrick J. Flynn. A 20th anniversary survey: Introduction to ‘content-based image retrieval at the end of the early years’. *TPAMI*, 22(12):1348, 2000. 3
- [6] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *ICCV*, pages 9630–9640. IEEE, 2021. 2
- [7] Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12m: Pushing web-scale image-text pre-training to recognize long-tail visual concepts. In *CVPR*, pages 3558–3568. IEEE, 2021. 1, 2, 6
- [8] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey E. Hinton. A simple framework for contrastive learning of visual representations. In *ICML*, pages 1597–1607. PMLR, 2020. 2
- [9] Ting Chen, Simon Kornblith, Kevin Swersky, Mohammad Norouzi, and Geoffrey E. Hinton. Big self-supervised models are strong semi-supervised learners. In *NeurIPS*, 2020. 2
- [10] Xinlei Chen, Haoqi Fan, Ross B. Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *CoRR*, 2020. 1, 2, 3, 6
- [11] Xinlei Chen, Saining Xie, and Kaiming He. An empirical study of training self-supervised vision transformers. In *ICCV*, pages 9620–9629. IEEE, 2021. 1, 2, 3, 5, 6
- [12] Xuxi Chen, Yu Yang, Zhangyang Wang, and Baharan Mirzsoleiman. Data distillation can be like vodka: Distilling more times for better quality. In *ICLR*, 2024. 2
- [13] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, pages 248–255. IEEE, 2009. 2, 6
- [14] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*. OpenReview.net, 2021. 2, 6
- [15] Jiawei Du, Qin Shi, and Joey Tianyi Zhou. Sequential subset matching for dataset distillation. In *NeurIPS*, 2023. 2
- [16] Bradley Efron. Second thoughts on the bootstrap. *Statistical science*, pages 135–140, 2003. 4
- [17] Utku Evci, Trevor Gale, Jacob Menick, Pablo Samuel Castro, and Erich Elsen. Rigging the lottery: Making all tickets winners. In *ICML*, pages 2943–2952. PMLR, 2020. 3
- [18] Alex Fang, Albin Madappally Jose, Amit Jain, Ludwig Schmidt, Alexander T. Toshev, and Vaishaal Shankar. Data filtering networks. In *ICLR*. OpenReview.net, 2024. 2
- [19] Vitaly Feldman and Chiyuan Zhang. What neural networks memorize and why: Discovering the long tail via influence estimation. In *NeurIPS*, 2020. 4
- [20] Tianyu Gao, Xingcheng Yao, and Danqi Chen. Simcse: Simple contrastive learning of sentence embeddings. In *EMNLP*, pages 6894–6910. ACL, 2021. 2
- [21] Akhilesh Gotmare, Nitish Shirish Keskar, Caiming Xiong, and Richard Socher. A closer look at deep learning heuristics: Learning rate restarts, warmup and distillation. In *ICLR*. OpenReview.net, 2019. 4
- [22] Sachin Goyal, Pratyush Maini, Zachary C. Lipton, Aditi Raghunathan, and J. Zico Kolter. Scaling laws for data filtering - data curation cannot be compute agnostic. In *CVPR*, pages 22702–22711. IEEE, 2024. 1
- [23] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778. IEEE, 2016. 2, 6
- [24] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross B. Girshick. Masked autoencoders are scalable vision learners. In *CVPR*, pages 15979–15988. IEEE, 2022. 2
- [25] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, Dawn Song, Jacob Steinhardt, and Justin Gilmer. The many faces of robustness: A critical analysis of out-of-distribution generalization. In *ICCV*, pages 8320–8329. IEEE, 2021. 2, 6
- [26] Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference-free evaluation metric for image captioning. In *EMNLP*, pages 7514–7528. ACL, 2021. 4
- [27] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katherine Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Oriol Vinyals, Jack W. Rae, and Laurent Sifre. An empirical analysis of compute-optimal large language model training. In *NeurIPS*, 2022. 1, 7
- [28] Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi, Ali Farhadi, and Ludwig Schmidt. Openclip, 2021. 2, 5, 6
- [29] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc V. Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *ICML*, pages 4904–4916. PMLR, 2021. 1, 3
- [30] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec

- Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *CoRR*, abs/2001.08361, 2020. 1, 7
- [31] Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In *ICML*, pages 1885–1894. PMLR, 2017. 1, 2
- [32] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. 2, 6
- [33] Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. *CS 231N*, 7(7):3, 2015. 2
- [34] Namhoon Lee, Thalaiyasingam Ajanthan, and Philip H. S. Torr. Snip: single-shot network pruning based on connection sensitivity. In *ICLR*. OpenReview.net, 2019. 2, 3
- [35] Tianjian Li, Haoran Xu, Philipp Koehn, Daniel Khashabi, and Kenton Murray. Error norm truncation: Robust training in the presence of data noise for text generation models. In *ICLR*, 2024. 2
- [36] Xin Li, Sima Behpour, Thang Long Doan, Wenbin He, Liang Gou, and Liu Ren. UP-DP: unsupervised prompt learning for data pre-selection with vision-language models. In *NeurIPS*, 2023. 2
- [37] Xianhang Li, Zeyu Wang, and Cihang Xie. An inverse scaling law for CLIP training. In *NeurIPS*, 2023. 2
- [38] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: common objects in context. In *ECCV*, pages 740–755. Springer, 2014. 2, 3, 5, 6
- [39] Shiwei Liu, Lu Yin, Decebal Constantin Mocanu, and Mykola Pechenizkiy. Do we actually need dense over-parameterization? in-time over-parameterization in sparse training. In *ICML*, pages 6989–7000. PMLR, 2021. 2
- [40] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *ICCV*, pages 9992–10002. IEEE, 2021. 2, 6
- [41] Ilya Loshchilov and Frank Hutter. SGDR: stochastic gradient descent with warm restarts. In *ICLR*. OpenReview.net, 2017. 4
- [42] Adyasha Maharana, Prateek Yadav, and Mohit Bansal. D2 pruning: Message passing for balancing diversity and difficulty in data pruning. In *ICLR*, 2024. 2
- [43] Anas Mahmoud, Mostafa Elhoushi, Amro Abbas, Yu Yang, Newsha Ardalani, Hugh Leather, and Ari Morcos. SIEVE: multimodal dataset pruning using image captioning models. In *CVPR*. IEEE, 2024. 2
- [44] Baharan Mirzasoleiman, Jeff A. Bilmes, and Jure Leskovec. Coresets for data-efficient training of machine learning models. In *ICML*, pages 6950–6960. PMLR, 2020. 2
- [45] Decebal Constantin Mocanu, Elena Mocanu, Peter Stone, Phuong H Nguyen, Madeleine Gibescu, and Antonio Liotta. Scalable training of artificial neural networks with adaptive sparse connectivity inspired by network science. *Nature communications*, 9(1):2383, 2018. 3
- [46] Michael C Mozer and Paul Smolensky. Using relevance to reduce network size automatically. *Connection Science*, 1(1): 3–16, 1989. 3
- [47] Aleksandra Nowak, Bram Grooten, Decebal Constantin Mocanu, and Jacek Tabor. Fantastic weights and how to find them: Where to prune in dynamic sparse training. In *NeurIPS*, 2023. 2
- [48] Vicente Ordonez, Girish Kulkarni, and Tamara L. Berg. Im2text: Describing images using 1 million captioned photographs. In *NIPS*, pages 1143–1151, 2011. 2, 5, 6
- [49] Mansheej Paul, Surya Ganguli, and Gintare Karolina Dzirgaitė. Deep learning on a data diet: Finding important examples early in training. In *NeurIPS*, pages 20596–20607, 2021. 1, 2, 3
- [50] Ziheng Qin, Kai Wang, Zangwei Zheng, Jianyang Gu, Xianguyu Peng, Daquan Zhou, and Yang You. Infobatch: Lossless training speed up by unbiased dynamic data pruning. In *ICLR*, 2024. 2, 5, 6
- [51] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *ICML*, pages 8748–8763. PMLR, 2021. 1, 2, 3
- [52] Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do imagenet classifiers generalize to imagenet? In *ICML*, pages 5389–5400. PMLR, 2019. 2, 6
- [53] Christoph Schuhmann, Robert Kaczmarczyk, Aran Komatsuzaki, Aarush Katta, Richard Vencu, Romain Beaumont, Jenia Jitsev, Theo Coombes, and Clayton Mullis. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. In *NeurIPS Workshop Datacentric AI*. Jülich Supercomputing Center, 2021. 4, 8
- [54] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski, Srivatsa Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev. LAION-5B: an open large-scale dataset for training next generation image-text models. In *NeurIPS*, 2022. 1, 3
- [55] Yuzhang Shang, Zhihang Yuan, and Yan Yan. MIM4DD: mutual information maximization for dataset distillation. In *NeurIPS*, 2023. 2
- [56] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *ACL*, pages 2556–2565. ACL, 2018. 1, 2, 5
- [57] DongHyeok Shin, Seungjae Shin, and Il-Chul Moon. Frequency domain-based dataset distillation. In *NeurIPS*, 2023. 2
- [58] Siva Sivaganesan. An introduction to the bootstrap (bradley efron and robert j. tibshirani). *SIAM Rev.*, 36(4):677–678, 1994. 4
- [59] Ben Sorscher, Robert Geirhos, Shashank Shekhar, Surya Ganguli, and Ari Morcos. Beyond neural scaling laws: beating power law scaling via data pruning. In *NeurIPS*, 2022. 2
- [60] Peng Sun, Bei Shi, Daiwei Yu, and Tao Lin. On the diversity and realism of distilled dataset: An efficient dataset distillation paradigm. In *CVPR*, 2024. 2
- [61] Hidenori Tanaka, Daniel Kunin, Daniel L. K. Yamins, and Surya Ganguli. Pruning neural networks without any data by iteratively conserving synaptic flow. In *NeurIPS*, 2020. 2, 3

- [62] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *ICML*, pages 10347–10357. PMLR, 2021. [2](#), [6](#)
- [63] Aäron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. 2018. [3](#)
- [64] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NIPS*, pages 5998–6008, 2017. [2](#)
- [65] Alex Jinpeng Wang, Kevin Qinghong Lin, David Junhao Zhang, Stan Weixian Lei, and Mike Zheng Shou. Too large; data reduction for vision-language pre-training. In *ICCV*, pages 3124–3134. IEEE, 2023. [2](#)
- [66] Weizhi Wang, Khalil Mrini, Linjie Yang, Sateesh Kumar, Yu Tian, Xifeng Yan, and Heng Wang. Finetuned multimodal language models are high-quality image-text data filters. *CoRR*, abs/2403.02677, 2024.
- [67] Yiping Wang, Yifang Chen, Wendan Yan, Alex Fang, Wenjing Zhou, Kevin Jamieson, and Simon Shaolei Du. Cliploss and norm-based data selection methods for multimodal contrastive learning. *CoRR*, abs/2405.19547, 2024. [2](#)
- [68] Ryan Webster, Julien Rabin, Loïc Simon, and Frédéric Jurie. On the de-duplication of LAION-2B. *CoRR*, 2023. [2](#)
- [69] Hu Xu, Saining Xie, Po-Yao Huang, Licheng Yu, Russell Howes, Gargi Ghosh, Luke Zettlemoyer, and Christoph Feichtenhofer. Cit: Curation in training for effective vision-language data. In *ICCV*, pages 15134–15143. IEEE, 2023. [2](#)
- [70] Hu Xu, Saining Xie, Xiaoqing Ellen Tan, Po-Yao Huang, Russell Howes, Vasu Sharma, Shang-Wen Li, Gargi Ghosh, Luke Zettlemoyer, and Christoph Feichtenhofer. Demystifying CLIP data. 2024. [2](#)
- [71] Shuo Yang, Zeke Xie, Hanyu Peng, Min Xu, Mingming Sun, and Ping Li. Dataset pruning: Reducing training data by examining generalization influence. In *ICLR*, 2023. [2](#), [5](#), [6](#), [7](#)
- [72] Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui Wu. Coca: Contrastive captioners are image-text foundation models. *TMLPR*, 2022. [3](#)
- [73] Geng Yuan, Xiaolong Ma, Wei Niu, Zhengang Li, Zhenglun Kong, Ning Liu, Yifan Gong, Zheng Zhan, Chaoyang He, Qing Jin, Siyue Wang, Minghai Qin, Bin Ren, Yanzhi Wang, Sijia Liu, and Xue Lin. MEST: accurate and fast memory-economic sparse training framework on the edge. In *NeurIPS*, pages 20838–20850, 2021. [2](#), [3](#)
- [74] Yuxin Zhang, Lirui Zhao, Mingbao Lin, Yunyun Sun, Yiwu Yao, Xingjia Han, Jared Tanner, Shiwei Liu, and Rongrong Ji. Dynamic sparse no training: Training-free fine-tuning for sparse llms. In *ICLR*, 2024. [2](#)
- [75] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for vision-language models. *IJCV*, 130(9):2337–2348, 2022. [3](#)