

# On Pretraining Data Diversity for Self-Supervised Learning

Hasan Abed Al Kader Hammoud\*  
KAUST

Tuhin Das\*  
University of Oxford

Fabio Pizzati\*  
University of Oxford

Philip Torr  
University of Oxford

Adel Bibi  
University of Oxford

Bernard Ghanem  
KAUST

## Abstract

We explore the impact of training with more diverse datasets, characterized by the number of unique samples, on the performance of self-supervised learning (SSL) under a fixed computational budget. Our findings consistently demonstrate that increasing pretraining data diversity enhances SSL performance, albeit only when the distribution distance to the downstream data is minimal. Notably, even with an exceptionally large pretraining data diversity achieved through methods like web crawling or diffusion-generated data, among other ways, the inherent distribution shift remains a challenge. Our experiments are comprehensive with seven SSL methods using large-scale datasets such as ImageNet and YFCC100M amounting to over 200 GPU days. The code and trained models will be made available online.

## 1. Introduction

Self-supervised learning (SSL) has recently emerged as a new paradigm to pretrain large vision models at scale [25, 26, 44]. Leveraging the ability to learn from unlabelled data, pretraining on millions—or even billions [26, 50]—of images turned from an unachievable goal to a common practice. This exposure to extremely *diverse* datasets, *i.e.*, composed of a remarkable number of unique samples, granted impressive performance and unprecedented generalization capabilities to a growing number of vision models [44]. Large-scale datasets, in conjunction with substantial computational resources, have been the key driving forces for the success of SSL-based approaches. For instance, SEER [25], pretrained for approximately 11 GPU years on a billion images, exemplifies the massive computation and data resources employed for these models. It has

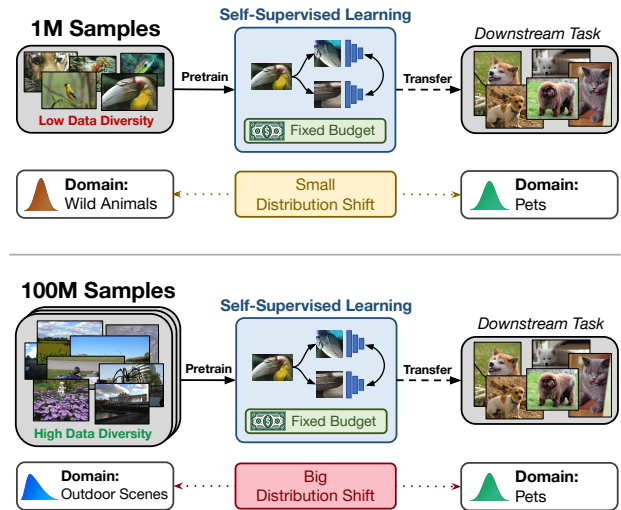


Figure 1. **Impact of Diversity on Pretraining:** Self-supervised learning (SSL) can be used to pretrain vision models on smaller datasets closely aligned to downstream task data, *e.g.*, pets classification, hence with a small distribution shift (top, wild animals pretraining). Conversely, we could pretrain on an extensively varied dataset, with wide distribution differences (outdoor scenes, bottom). We demystify the role of pretraining diversity in SSL under a fixed computational budget, and highlight its effects in relationship to the distribution shift.

become the implicit norm that increasing computation and data is beneficial, without any detailed analysis of how they separately impact SSL effectiveness. In particular, it is not clear to what extent large datasets are responsible for the impressive generalization capabilities of SSL models. Indeed, consider the example in Figure 1. Assuming a fixed monetary budget, in the form of computational expenses, for pretraining a vision model under the SSL paradigm, does iterating over a large set of images work best given a downstream task, or it is better to train repeatedly on a smaller set

\*Equal contribution as first authors. Please send any correspondence to [fabio@robots.ox.ac.uk](mailto:fabio@robots.ox.ac.uk)

of samples visually close to the ones of the downstream? In this context, we face the ever-lasting problem of distribution shift on vision models [55]. Without a proper understanding of the impact of data, there could be a wrong allocation of efforts in increasing data or computation, leading to inefficiencies in the deployment process of SSL models.

In this paper, we study the role of the diversity of the SSL pretraining data, under a fixed computational budget, when the pretraining data matches the distribution of the downstream data, as well as when they differ. Our experiments span various SSL methods [4, 8–10, 14, 27, 62] and multiple datasets [5, 18, 35–37, 41, 45, 52, 53, 66], and are tested under several computational budgets amounting to a total of 200 GPU days. We summarize our conclusions below:

1. Under a fixed computational budget, SSL pretraining on more diverse datasets can notably boost downstream model performance. This advantage diminishes, and might even become harmful, when the pretraining and target task data distributions are considerably dissimilar.
2. If the class labels of downstream tasks are known a priori, increasing pretraining data diversity with web-crawling or synthetic image generation is a reliable method to improve downstream performance. If class labels are unknown, augmenting the data with new samples might lead to reduced performance instead.
3. Under a fixed computational budget, significantly increasing the number of pretraining data samples does not adequately address the distribution shift between pretraining and downstream datasets.

This work offers an in-depth analysis of how data diversity and computational resources interact and of their combined effect on the performance of SSL models, with an aim to improve pretraining practices. Additionally, it underscores the importance of normalizing computational costs for a proper comprehension of the ways in which vision models utilize diverse datasets to achieve better generalization in downstream tasks.

## 2. Related Work

**Self-Supervised Learning** Early work in self-supervised learning used simple pretext tasks, such as relative patch prediction [19, 20], image colorization [64], image rotation prediction [23], or solving jigsaw puzzles [42] to train feature extractors in absence of annotations. Representations learned with those methods are limitedly effective, hence more recent literature has moved towards more sophisticated approaches, such as using stochastic image augmentations to generate correlated views of a training sample, and learning to extract augmentation-invariant representations for these correlated pairs. Among these multi-view methods, many exploit contrastive losses [1, 8, 10, 11, 13, 29, 32, 33, 39, 40, 43], enforcing similarity between views of the same image (positives) and

dissimilarity from other images (negatives). However, due to the need of many negative samples, contrastive methods often require large batch sizes to work effectively [10, 14]. Cluster-based methods such as SwAV [8], DINO [9] and DeepCluster v2 [6] learn generalizable representations by grouping samples into cluster prototypes. Others exploit prediction of features using siamese networks [12], learn features in a teacher-student learning fashion [27], or employ redundancy reduction techniques [4, 62]. Finally, masked image modeling [3, 30, 59] emerged as a scalable alternative for Vision Transformers [21], learning representations by predicting masked image patches as SSL task.

**Pretraining at Scale** SSL is most effective when pretraining is conducted at scale, benefiting from large datasets and great computational resources. Initial attempts at large-scale pretraining were made via combining contrastive learning and clustering [7, 54]. The SEER model [25] was trained on 1 billion images scraped from the internet using SwAV [8] applied on a ResNet [46] backbone. The importance of model scaling to leverage large-scale datasets is also shown in [24, 26, 60], as well as the need for increasing training duration [60]. Additional training strategies may be necessary to effectively achieve best performance at scale [17, 63]. All considered works focus on reaching the best representations, without further considerations about training costs, thus encouraging unfair comparisons about the efficacy of data. A preliminary work [15] found for SimCLR [10] that an increase in data from 500 thousand to 1 million images leads to a modest increase in downstream performance, but without limiting the training budget. While large uncurated datasets have been further explored [54], their efficacy in relationship with the distribution shift under normalized computation is still limitedly investigated.

**Distribution Shift in SSL** Some studies have investigated how the pretraining data domain affects downstream performance on other domains. For various pretraining datasets, preliminary works [15, 24, 34, 38] observed that the best performing representations were learned on datasets that were similar to the downstream test task. Additionally, combining datasets before pretraining or combining self-supervised features learned from various datasets did not lead to significant improvements in [15]. In [49, 65], they showcase higher generalization capabilities of SSL models compared to their supervised counterparts, for several downstream tasks in the presence of a distribution shift. In [56], they pretrained on several datasets, and observed different generalization depending on the object-centric or scene-centric appearance of the pretraining dataset. Furthermore, initial considerations on the impact of external data on downstream tasks under distribution shift have

been proposed in [57]. Although partial information can be inferred by these works, there is still a lack of a fair, computation-normalized evaluation that allows to study the effects of the distribution shift in a controlled environment.

### 3. Preliminaries

**Pretraining** We first outline the general pretraining procedure common to state-of-the-art self-supervised learning methods. Specific differences among these methods are detailed in the supplementary material. The overall pretraining pipeline, common across many SSL approaches [4, 8–10, 27, 62], goes as follows: (1) **Data Sampling**: from a large-scale, unlabeled *upstream* pretraining dataset, denoted as  $\mathbb{D}_{\text{SSL}}$ , an image  $\mathbf{x}$  is randomly sampled; (2) **View Generation**: two correlated views,  $\tilde{\mathbf{x}}_A$  and  $\tilde{\mathbf{x}}_B$ , are generated from  $\mathbf{x}$  using two image augmentation functions sampled at random. These random augmentations include random resized cropping, horizontal flipping, blurring, and color adjustments [2], among others; (3) **Feature Extraction and Projection**: the correlated views undergo feature extraction through a network,  $f_{\theta_f}$ , parameterized by  $\theta_f$ , such as a ResNet [28] or a ViT [21], leading to representations  $\mathbf{h}_A = f_{\theta_f}(\tilde{\mathbf{x}}_A)$  and  $\mathbf{h}_B = f_{\theta_f}(\tilde{\mathbf{x}}_B)$ . A linear projection head,  $g_{\theta_g}$ , parameterized by  $\theta_g$ , then maps these representations to a latent space, resulting in  $\mathbf{z}_A = g_{\theta_g}(\mathbf{h}_A)$  and  $\mathbf{z}_B = g_{\theta_g}(\mathbf{h}_B)$ ; (4) **Joint Optimization**: The feature extractor  $f_{\theta_f}$  and the projection head  $g_{\theta_g}$  are jointly optimized according to the following objective:

$$\theta_f^*, \theta_g^* = \arg \min_{\theta_f, \theta_g} \mathbb{E}_{\mathbf{x} \sim \mathbb{D}_{\text{SSL}}} \mathcal{L}_{\text{SSL}}(\mathbf{z}_A, \mathbf{z}_B), \quad (1)$$

where  $\mathcal{L}_{\text{SSL}}$  is a loss function specific to the chosen SSL pretraining method.

After pretraining, the feature extractor  $f_{\theta_f^*}$  can be deployed for various *downstream* tasks such as image classification, object detection, or segmentation. This is typically achieved by training a task-specific head. Alternatively, the feature extractor  $f_{\theta_f^*}$  can either be fine-tuned or used together with a  $k$ -nearest neighbors (kNN) classifier.

**Linear Probing** There are several ways to evaluate the performance of a self-supervised learning method such as linear probing [10, 14, 27], kNN [8, 9, 58, 68], and few-shot evaluation [22, 25]. Consistent with the general protocol established in the literature [10, 15, 29, 64], we use linear probing to measure the quality of the features extracted for classification tasks. The procedure is as follows: (1) a labeled *downstream* dataset,  $\mathbb{D}_{\text{task}}$ , consisting of image-class pairs  $(\mathbf{x}, \mathbf{y}) \sim \mathbb{D}_{\text{task}}$ , is selected for evaluation. (2) For each image  $\mathbf{x}$ , its representation is extracted using the pretrained feature extractor  $f_{\theta_f^*}$ , after which the linear classification head  $t_{\theta_t}$ , parameterized by  $\theta_t$ , is then applied to obtain  $t_{\theta_t}(f_{\theta_f^*}(\mathbf{x}))$ . (3) The linear head  $t_{\theta_t}$  is optimized as

follows:

$$\theta_t^* = \arg \min_{\theta_t} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathbb{D}_{\text{task}}} [\mathcal{L}_{\text{task}}(t_{\theta_t}(f_{\theta_f^*}(\mathbf{x})), \mathbf{y})], \quad (2)$$

Note that only the parameters of the linear head  $\theta_t$  are optimized, while the parameters  $\theta_f^*$  of the feature extractor are kept frozen. The quality of features extracted by  $f_{\theta_f^*}$  is directly inferred from the classification accuracy achieved on the test set of  $\mathbb{D}_{\text{task}}$ , which serves as the primary indicator of the quality of the extracted features.

### 4. Normalized Evaluation

We stress how for a correct understanding of the impact of data diversity we need to analyze its effects isolating them from the impact of increased computation. To enable this, we introduce (1) a computational budget used to normalize computation across experiments, and (2) a quantification of the data diversity seen by models during pretraining.

**Computational Budget** Current progress in SSL pretraining simultaneously scales computational budget and dataset size to achieve the best performance [25]. This makes it difficult to assess the reasons behind the success of many SSL algorithms. Do these methods benefit mainly from the large amounts of computation, *i.e.*, running SSL pretraining for large numbers of epochs, or do they benefit from data diversity in larger datasets containing a vast variety of visual features? To perform a disentangled evaluation of these two factors, we first introduce  $\mathcal{C}$  as a measure of computational budget, which quantifies the total number of images an SSL method is allowed to process during pretraining. This is calculated as  $\mathcal{C} = N \cdot \mathcal{E}$ , where  $N$  is the number of unique samples in the pretraining dataset  $\mathbb{D}_{\text{SSL}}$ , hence the *data diversity* of the dataset, and  $\mathcal{E}$  is the number of pretraining epochs. Constraining multiple models pretrained with SSL to a fixed computational budget  $\mathcal{C}$  allows for meaningful comparison, as it is guaranteed that all SSL methods will have processed the same number of pretraining images.

**Quantifying Pretraining Diversity** Even under normalized computation  $\mathcal{C}$ , various SSL approaches may be exposed to different data diversity as they are trained with different datasets. For instance, a model trained for  $\mathcal{E} = 1000$  epochs on a dataset of size  $N = 1000$  will see less diversity than a model pretrained for  $\mathcal{E} = 1$  epoch on a dataset of size  $N = 10^6$ , despite that they are both pretrained under normalized computation, processing the same number of images. Hence, to capture this notion of exposure to diversity while pretraining under normalized computation  $\mathcal{C}$ , we define a *pretraining diversity*  $\mathcal{D}$ , which captures the number of unique samples encountered during training given a fixed  $\mathcal{C}$  as  $\mathcal{D} = N/\mathcal{C} = 1/\mathcal{E}$ . A model pretrained with larger  $\mathcal{D}$  indicates that the model is presented with a larger number of unique images during training with fixed  $\mathcal{C}$ , and hence is ex-

$N$ ( $\times 10^3$ )	$\mathcal{D}$ ( $\times 10^{-3}$ )	Accuracy $\uparrow$						
		SimCLR	B.T.	BYOL	SwAV	VICReg	MoCoV3	DINO
5	0.1	37.95	44.30	41.14	14.06	44.06	12.51	15.74
25	0.5	50.13	54.29	51.01	45.44	50.37	51.58	35.89
50	1.0	<b>58.59</b>	<b>58.42</b>	<b>58.28</b>	<b>56.37</b>	<b>55.83</b>	<b>55.61</b>	<b>40.93</b>

$\mathcal{C} = 50 \times 10^6$   
(a) CIFAR-100

$N$ ( $\times 10^3$ )	$\mathcal{D}$ ( $\times 10^{-3}$ )	Accuracy $\uparrow$						
		SimCLR	B.T.	BYOL	SwAV	VICReg	MoCoV3	DINO
10	0.2	36.91	40.98	35.62	34.23	37.58	36.56	34.54
50	1.0	48.77	52.01	48.05	43.63	48.75	46.45	44.39
100	2.0	<b>49.83</b>	<b>55.60</b>	<b>50.29</b>	<b>47.48</b>	<b>54.10</b>	<b>48.58</b>	<b>47.32</b>

$\mathcal{C} = 50 \times 10^6$   
(b) Tiny ImageNet

Table 1. **Impact of Data Diversity on CIFAR-100 and Tiny ImageNet SSL Pretraining Performance:** We study the effects of diversity on CIFAR-100 (a) and Tiny ImageNet (b) across seven different methods and three data diversity settings for a ResNet-18 pretraining, where for all,  $\mathbb{D}_{\text{task}} = \mathbb{D}_{\text{SSL}}$ . This comparison includes analyzing classification accuracy through linear probing on the training set and evaluation on the test set of the respective datasets. Although performance fluctuates among different methods, a consistent trend is observed: higher data diversity typically leads to the generation of higher quality representations.

posed to more pretraining data diversity. In the next section, we explore the effects of variations in  $\mathcal{D}$  on SSL performance under a distribution shift, while keeping  $\mathcal{C}$  constant.

## 5. Fixed Budget SSL: In & Out-of-Distribution

**Training Configuration** We evaluate seven SSL methods: SimCLR [10], MoCoV3 [14], VICReg [4], DINO [9], BYOL [27], Barlow Twins [62], and SwAV [8]. We use different datasets for both pretraining  $\mathbb{D}_{\text{SSL}}$  and linear probing  $\mathbb{D}_{\text{task}}$  for different sections. We use `solo-learn` [16] as a codebase. For each method, we use the default parameters provided when available, otherwise, we conduct a hyperparameter search to ensure proper convergence.

### 5.1. Performance on the Same Distribution

We aim to answer a first question: *does collecting more samples from the same distribution help SSL pretraining with a fixed budget?* We conduct experiments to capture how the pretraining diversity  $\mathcal{D}$  influences SSL pretraining within a normalized  $\mathcal{C}$ , focusing on the simplest setting where the upstream and downstream datasets belong to the same distribution such that  $\mathbb{D}_{\text{SSL}} = \mathbb{D}_{\text{task}}$ . This serves as a fundamental ground for our further experiments.

**Setup** For pretraining, we use CIFAR-100 [36] and Tiny ImageNet [37] as  $\mathbb{D}_{\text{SSL}}$ , which contain  $50 \times 10^3$  and  $100 \times 10^3$  images, respectively. We set  $\mathcal{C} = 50 \times 10^6$ , chosen such that it allows all methods to converge during pretraining. In Section 6, we study the effect of varying the budget  $\mathcal{C}$ . We pretrain on subsets of  $\mathbb{D}_{\text{SSL}}$  with different sizes  $N$  (10%, 50%, and 100% of  $\mathbb{D}_{\text{SSL}}$ ), enabling us to observe the effects of pretraining diversity  $\mathcal{D}$  on training outcomes where  $\mathcal{E}$  is adjusted accordingly to match the budget  $\mathcal{C}$ . For example, using 100% of CIFAR-100 involves 1000 epochs of training, while 10% and 50% of the dataset lead to 10000 and 2000 epochs, respectively. These subsets of  $\mathbb{D}_{\text{SSL}}$  are created by shuffling the dataset and then selecting the first 10%, 50%, or 100% split. All models in this section use a ResNet-18 [28] backbone pretrained from scratch. For evaluation, we use in-distribution linear

probing, *i.e.*,  $\mathbb{D}_{\text{task}} = \mathbb{D}_{\text{SSL}}$ , where performance is measured by classification accuracy on the test set.

**Results** The results of linear probing are presented in Tables 1a and 1b for CIFAR-100 and Tiny ImageNet, respectively. While different SSL methods show varying levels of performance, we observe that, for all methods, an increase in the diversity  $\mathcal{D}$  consistently leads to higher linear probing accuracy, suggesting that including more in-distribution samples in  $\mathbb{D}_{\text{SSL}}$  helps under a normalized  $\mathcal{C}$ . For example, SimCLR achieves an accuracy of 37.95 when 10% of CIFAR-100 is provided, whereas this performance improve by around 12% when only 50% of the dataset is provided and by another 8% after pretraining on 100% of unique samples. This suggests that SSL methods substantially benefit from having a more diverse pretraining set in computationally budgeted in-distribution evaluation, a fundamental verification that allows us to proceed in our analysis.

### Conclusion (1)

When the distributions of the upstream and downstream tasks are the same, *i.e.*,  $\mathbb{D}_{\text{SSL}} = \mathbb{D}_{\text{task}}$ , in a normalized computation setup, increasing pretraining diversity  $\mathcal{D}$  proves to be an effective strategy for enhancing SSL pretraining performance.

### 5.2. Increasing Data Diversity

As observed in Section 5.1, if  $\mathbb{D}_{\text{SSL}} = \mathbb{D}_{\text{task}}$  having a more diverse pretraining set benefits SSL methods, under a normalized computation assumption. However, to increase diversity, sampling from the same distribution  $\mathbb{D}_{\text{task}}$  to extend  $\mathbb{D}_{\text{SSL}}$  is not always attainable. In fact, pretraining data is often sampled from a distribution different than  $\mathbb{D}_{\text{task}}$ . The added samples will then introduce a distribution shift between  $\mathbb{D}_{\text{SSL}}$  and  $\mathbb{D}_{\text{task}}$ . Then, we raise the following question: *is increasing pretraining diversity still beneficial when there is a distribution shift between  $\mathbb{D}_{\text{SSL}}$  and  $\mathbb{D}_{\text{task}}$ , *i.e.*,  $\mathbb{D}_{\text{SSL}} \neq \mathbb{D}_{\text{task}}$ ?* We explore strategies for acquiring new data for increasing  $\mathcal{D}$ , namely, including existing samples,



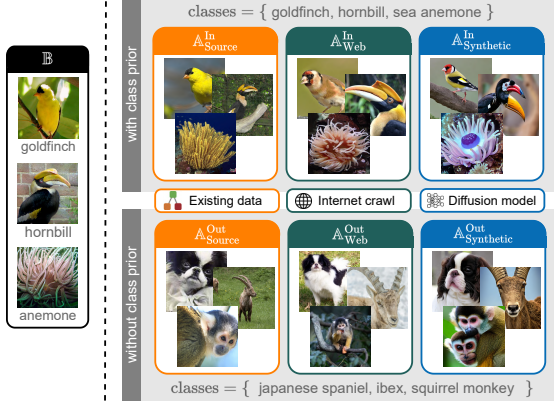


Figure 2. **Data Collection Strategies:** We analyze strategies for collecting additional data ( $\mathbb{A}$ ), *i.e.*, collecting more **source** data, crawling the **web** or using **synthetic** images. Using a class prior (top row) simulates **In**-distribution trainings. We also collect images without class prior (bottom row) to analyze the interactions between diversity and **Out**-of-distribution classes.

crawled internet data, and data synthetically generated by diffusion models. To evaluate the effects of the distribution shift in a controlled scenario, we analyze distributions closer to  $\mathbb{D}_{\text{task}}$  by using a class prior (in-distribution) and without a class prior (out-of-distribution).

**Setup** We use ImageNet-100 [53] as  $\mathbb{D}_{\text{task}}$ , and we construct multiple  $\mathbb{D}_{\text{SSL}}$  to evaluate the effects of different data collection strategies and the distribution shift. We first introduce a set  $\mathbb{B}$  composed of  $65 \times 10^3$  images from ImageNet-100 (50% of the dataset), which we use as a baseline for  $\mathbb{D}_{\text{SSL}}$  with minimum diversity. We denote the 100 classes in ImageNet-100 as  $\mathcal{T}_{100}$ , and sample  $\mathbb{B}$  uniformly including images from all classes. Next, we compare with pretraining on more diverse datasets as  $\mathbb{D}_{\text{SSL}}$ , imposing  $\mathbb{D}_{\text{SSL}} = \mathbb{B} \cup \mathbb{A}$  where  $\mathbb{A}$  includes  $65 \times 10^3$  images sampled with one of three strategies. To highlight the effects of the distribution shift, we include in  $\mathbb{A}$  either images from **In**-distribution classes, *i.e.*, selecting images from classes overlapping with  $\mathcal{T}_{100}$ , or images from **Out**-of-distribution classes, which do not overlap with  $\mathcal{T}_{100}$ . This is to study the effects of the distribution shift, since we do *not* assume access to downstream classes in real scenarios. For the **Out**-of-distribution samples, we define  $\mathbb{A}$  as (1) random images sampled from the full ImageNet [18] dataset; (2) images crawled from Flickr, Bing, and DuckDuckGo; or (3) images generated with Stable Diffusion V2.1 [48]. We respectively call these sets  $\mathbb{A}_{\text{Source}}^{\text{Out}}$ ,  $\mathbb{A}_{\text{Web}}^{\text{Out}}$ , and  $\mathbb{A}_{\text{Synthetic}}^{\text{Out}}$ . We also define their **In**-distribution counterparts as  $\mathbb{A}_{\text{Source}}^{\text{In}}$ ,  $\mathbb{A}_{\text{Web}}^{\text{In}}$ , and  $\mathbb{A}_{\text{Synthetic}}^{\text{In}}$ , respectively. Note how  $\mathbb{B} \cup \mathbb{A}_{\text{Source}}^{\text{In}}$  is the full ImageNet-100, coherently with Section 5.1. Figure 2 shows each collection strategy, for which we provide implementation details in supplementary material.

Although many factors (such as the appearance of images) impact the distribution shift, using a class prior imposes that any strategy using it would still result in closer distribution with respect to the same strategy *without* class priors [15]. We pretrain a ResNet-18 from scratch, with the same settings of Section 5.1, and  $C = 50 \times 10^6$ . Note that this results in  $\mathcal{D} = 1.25 \times 10^{-3}$  for pretraining on  $\mathbb{B}$ , and  $\mathcal{D} = 2.5 \times 10^{-3}$  for pretraining on  $\mathbb{B} \cup \mathbb{A}$ , introducing a  $\mathcal{D}$  difference of a factor of 2.

**Results** We report the linear probing accuracies for pretraining on each  $\mathbb{B} \cup \mathbb{A}$  as bars in Figure 3, showing the  $C$ -normalized training on  $\mathbb{B}$  as a dashed line. Surprisingly, without class priors (**Out**), including  $\mathbb{A}_{\text{Source}}^{\text{Out}}$ ,  $\mathbb{A}_{\text{Web}}^{\text{Out}}$ , and  $\mathbb{A}_{\text{Synthetic}}^{\text{Out}}$  underperforms compared to pretraining on  $\mathbb{B}$  only. For instance, for SimCLR and while  $\mathbb{B}$  scores 72.3% accuracy, increasing diversity reduces the accuracy by 1% in all cases. This might appear to conflict with our findings in Section 5.1, however, the inclusion of **Out** samples leads to  $\mathbb{D}_{\text{SSL}} \neq \mathbb{D}_{\text{task}}$ , since we sample only classes not included in  $\mathbb{D}_{\text{task}}$ . We infer that, with normalized  $C$ , increasing  $\mathcal{D}$  without distribution priors may negatively impact performance. Conversely, when class priors are available (**In**), increasing pretraining diversity  $\mathcal{D}$  improves performance compared to  $\mathbb{B}$  pretraining. For instance, pretraining on  $\mathbb{A}_{\text{Web}}^{\text{In}}$  performs comparably to augmenting with additional ImageNet samples ( $\mathbb{A}_{\text{Source}}^{\text{In}}$ ), as in the case of SimCLR where the inclusion of  $\mathbb{A}_{\text{Web}}^{\text{In}}$  scores only 0.2% lower than  $\mathbb{A}_{\text{Source}}^{\text{In}}$ . Including  $\mathbb{A}_{\text{Synthetic}}^{\text{In}}$  data also helps, Synthetic  $\mathbb{A}_{\text{Synthetic}}^{\text{In}}$  data helps, although more limitedly due to the visual discrepancies between generated and real images. Ultimately, the effectiveness of diversity is linked to the distribution shift. These findings highlight the impact of the distribution shift on computationally-normalized SSL pretraining and help define evaluation practices for this task (see Section 7).

### Conclusion (2)

When the distributions of the upstream and downstream tasks differ,  $\mathbb{D}_{\text{SSL}} \neq \mathbb{D}_{\text{task}}$ , and in a normalized computation setup, increasing pretraining diversity  $\mathcal{D}$  may harm SSL pretraining performance, reflecting the adverse effects of distribution shift.

### 5.3. Scaling Pretraining Diversity

We showed that diversity improves pretraining performance when the training set and downstream task share the same data distribution ( $\mathbb{D}_{\text{SSL}} = \mathbb{D}_{\text{task}}$ ), as discussed in Section 5.1. This may change when a distribution shift is introduced, as explored in Section 5.2. However, it is still unclear from Section 5.2, whether including a larger number of samples, and thus increasing considerably the pretraining diversity, can compensate for the negative effects of the distribution

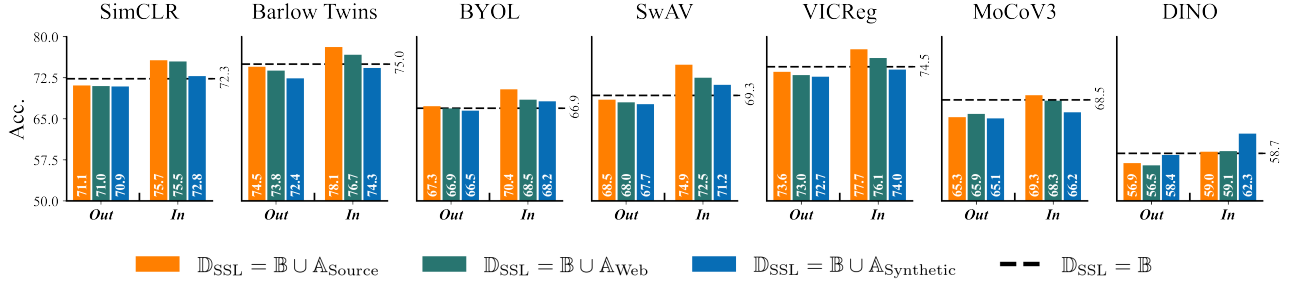


Figure 3. **Effect of Various Data Sources on SSL Pretraining:** We use a baseline set  $\mathbb{B}$  (black dashed line), comprising  $65 \times 10^3$  images from ImageNet-100, for pretraining a ResNet-18 with  $\mathcal{C} = 50 \times 10^6$ . Augmenting  $\mathbb{B}$  with **In**-distribution images enhances performance (above black line), while **Out**-of-distribution augmentations reduce it (below black line).

SimCLR									
	$\mathbb{D}_{SSL}$	$N$ ( $\times 10^6$ )	$\mathcal{D}$ ( $\times 10^{-3}$ )	Accuracy $\uparrow$					
				ImageNet	Cars	Flow.	Pets	Places	Food
ResNet-50	ImageNet	0.128	1.31	56.9	43.0	82.3	73.9	45.4	59.9
	ImageNet	0.256	2.61	61.1	45.5	84.0	76.0	47.0	64.4
	ImageNet	0.640	6.54	63.7	<b>46.8</b>	84.2	78.8	48.3	67.2
	ImageNet	1.281	13.0	<b>64.5</b>	46.4	<b>84.6</b>	<b>79.5</b>	48.8	<b>68.0</b>
	YFCC	98.17	1000	57.3	37.2	76.6	58.9	<b>50.1</b>	62.1
ViT-B/16	ImageNet	0.128	1.31	54.2	37.2	81.8	70.3	44.6	64.3
	ImageNet	0.256	2.61	61.3	39.5	83.8	77.7	47.2	69.4
	ImageNet	0.640	6.54	65.5	<b>39.9</b>	<b>84.6</b>	80.4	49.1	73.1
	ImageNet	1.281	13.0	<b>66.7</b>	39.5	83.7	<b>81.7</b>	<b>50.0</b>	<b>73.6</b>
	YFCC	98.17	1000	54.5	25.0	72.5	59.7	49.5	65.4

$\mathcal{C} = 98 \times 10^6$

MoCoV3									
	$\mathbb{D}_{SSL}$	$N$ ( $\times 10^6$ )	$\mathcal{D}$ ( $\times 10^{-3}$ )	Accuracy $\uparrow$					
				ImageNet	Cars	Flow.	Pets	Places	Food
ResNet-50	ImageNet	0.128	1.31	58.1	40.6	81.8	76.04	45.1	63.9
	ImageNet	0.256	2.61	62.9	45.3	85.0	79.8	47.6	68.7
	ImageNet	0.640	6.54	65.4	48.4	86.1	81.9	49.2	71.0
	ImageNet	1.281	13.0	<b>65.9</b>	<b>48.8</b>	<b>86.6</b>	<b>82.6</b>	49.5	<b>71.9</b>
	YFCC	98.17	1000	60.4	42.2	82.6	66.3	<b>50.7</b>	67.3
ViT-B/16	ImageNet	0.128	1.31	57.9	33.4	82.8	78.0	46.7	67.8
	ImageNet	0.256	2.61	63.7	35.0	85.3	82.9	48.4	71.0
	ImageNet	0.640	6.54	67.2	39.5	85.0	85.8	49.7	72.8
	ImageNet	1.281	13.0	<b>68.8</b>	<b>41.9</b>	<b>86.5</b>	<b>86.5</b>	<b>50.3</b>	<b>73.8</b>
	YFCC	98.17	1000	57.2	25.2	70.3	43.8	<b>50.3</b>	64.0

$\mathcal{C} = 98 \times 10^6$

Table 2. **Performance of ImageNet and YFCC100M SSL Pretraining on Various Downstream Tasks:** We train ResNet-50 and ViT-B/16 under normalized computation ( $\mathcal{C} = 98 \times 10^6$ ) using SimCLR (left) and MoCoV3 (right) on ImageNet and YFCC100M (YFCC) with multiple  $\mathcal{D}$ . Despite the significantly larger  $\mathcal{D}$  when trained on YFCC100M, these models cannot offset the effects of the distribution shift and are outperformed by models pretrained on ImageNet in the majority of downstream classification tasks.

shift. To address this, the following section presents larger-scale experiments, employing significantly varied  $\mathcal{D}$  values, aiming to explore the interplay between pretraining diversity and different distributions using millions of samples.

**Setup** For our large-scale pretraining experiments, we set  $\mathbb{D}_{SSL}$  to be two datasets of significantly different sizes: ImageNet [18] and YFCC100M [52], comprising 1.28 million and 98 million images, respectively. Following Section 5.1, we explore multiple  $\mathcal{D}$  values for pretraining. We set  $\mathcal{C} = 98 \times 10^6$ , which corresponds to one epoch on YFCC100M, effectively iterating once through each of its 98 million images to maximize diversity ( $\mathcal{D} = 1$ ). Normalizing  $\mathcal{C}$  (see Section 5.1), we pretrain on ImageNet for approximately 77 epochs, cumulatively utilizing 98 million samples. Due to the extensive cost of these experiments, we focus on SimCLR and MoCoV3 only. We also employ larger architectures, namely ResNet-50 [28] and ViT-B/16 [21], to leverage the extensive scale of the pretraining datasets. We also use multiple  $\mathbb{D}_{task}$  including ImageNet [18], Stanford Cars [35], Flowers-102 [41], Oxford-IIIT Pets [45], Places365 [66], and Food-101 [5].

**Results** The results of our large-scale experiments are detailed in Table 2. Consistently with findings in Sec-

tion 5.1, increasing  $\mathcal{D}$  leads to better pretraining efficacy when  $\mathbb{D}_{SSL} = \mathbb{D}_{task}$ . This is evident when ImageNet is used for both pretraining and downstream classification, reinforcing that our observations hold even at a larger scale. Instead, models pretrained on YFCC100M show substantially lower performance compared to those pretrained on ImageNet, although having much higher  $\mathcal{D}$ . This highlights the inefficiency of collecting data indiscriminately without considering distribution shifts. To stress this, note how the model pretrained on YFCC100M ( $\mathcal{D} = 1$ ) often performs similarly to those pretrained with drastically lower  $\mathcal{D}$  on ImageNet ( $\mathcal{D} = 1.31 \times 10^{-3}$ ). This aligns with our observations in Section 5.2, emphasizing that distribution differences remain a significant factor even when training with large datasets. However, the YFCC100M-pretrained model outperforms the ImageNet-pretrained model in Places365, suggesting a closer distribution relationship between YFCC100M and Places365. We explore this hypothesis further in Section 6, where we analyze distribution distances with existing metrics. Ultimately, our analysis highlights that scaling the data is not an effective solution for compensating the distribution shift, when computational resources are normalized.

### Conclusion (3)

Even an extremely large data diversity cannot mitigate the distribution shift under normalized computation. This emphasizes the importance of further research into how effectively SSL pretraining algorithms generalize.

## 6. Additional Analysis

Previously, we proposed a computationally-normalized evaluation to assess the role of  $\mathcal{D}$  with and without the distribution shift. We highlighted that, although pretraining diversity helps when  $\mathbb{D}_{\text{SSL}} = \mathbb{D}_{\text{task}}$ , vision models are unable to compensate for the distributional differences under normalized computation. Now, we analyze additional elements that support our previous observations.

**Distribution Distances Using FID & VisualDNA** In Section 5.3, we showed that pretraining on ImageNet typically outperforms YFCC100M on a variety of downstream tasks, with Places365 being the exception. We speculate that the distribution of ImageNet is more aligned with those of the downstream tasks compared to YFCC100M. To verify this, we evaluate the similarity between the datasets using FID [31] and VisualDNA [47] and report results in Table 3. With both FID or VisualDNA, the distribution of ImageNet is always closer to the downstream tasks, except for Places365 where YFCC100M is closer. This aligns with the lower performance of ImageNet on this dataset only (Table 2), further suggesting that the performance drop is caused by the distribution shift.

**Importance of Normalized Computation** We now study the impact of normalizing computation on the performance of SSL methods. We aim to understand if  $\mathcal{C}$  is not normalized across methods, this will lead to misleading conclusions, and thus motivating our normalized computation. In Table 4, we compare the performance of ResNet-50 and ViT-B/16 models pretrained using MoCoV3 with  $\mathcal{C} = 98 \times 10^6$  (as used in Section 5.3) against a tripled budget of  $\mathcal{C} = 294 \times 10^6$ . We show that inconsistencies arise when pretraining data diversity  $\mathcal{D}$  and computation  $\mathcal{C}$  are not fairly compared. For instance, in scenarios highlighted in red, we notice that a lower pretraining diversity coupled with a higher computational budget can yield better results. We observe that under pretraining diversity of only  $\mathcal{D} = 13 \times 10^{-3}$  and a computational budget of  $\mathcal{C} = 98 \times 10^6$  ResNet-50 only enjoys a 65.9% accuracy compared against with 69.8% with a smaller pretraining diversity of  $2.17 \times 10^{-3}$  but with a larger computational budget of  $294 \times 10^6$ . This shows that without normalized computation, it could be incorrectly concluded that pretraining diversity does not play a significant role.

$\mathbb{D}_{\text{SSL}}$	Backbone	VisualDNA↓					
		ImageNet	Cars	Flow.	Pets	Places	Food
ImageNet	MUGS	<b>0.00</b>	<b>11.49</b>	<b>12.46</b>	<b>6.09</b>	7.19	<b>9.08</b>
YFCC	MUGS	3.72	11.57	12.71	7.93	<b>6.20</b>	9.72
ImageNet	DINO v2	<b>0.00</b>	<b>7.05</b>	<b>6.95</b>	<b>4.43</b>	3.52	<b>7.26</b>
YFCC	DINO v2	2.37	7.12	7.46	5.75	<b>2.74</b>	7.90

$\mathbb{D}_{\text{SSL}}$	Backbone	FID↓					
		ImageNet	Cars	Flow.	Pets	Places	Food
ImageNet	Inception	<b>0.00</b>	<b>143.40</b>	<b>192.89</b>	<b>88.85</b>	64.91	<b>114.92</b>
YFCC	Inception	48.14	174.27	214.78	145.79	<b>38.86</b>	154.51

Table 3. **Comparing VisualDNA and FID Scores Across Datasets:** We assess the relationship between VisualDNA [47] and FID Score [31] for various large-scale  $\mathbb{D}_{\text{SSL}}$  and several downstream tasks. For VisualDNA, activations are extracted as suggested [47] with MUGS [67] with ViT-B/16 (MUGS) or DINO v2 [44] with ViT-B/16 (DINO v2), while FID activations are obtained via an Inception network [51]. Consistently, across all metrics, the ranking of dataset distances between  $\mathbb{D}_{\text{SSL}}$  and various  $\mathbb{D}_{\text{task}}$  aligns with the accuracy ranking in Table 2. Models exhibiting lower VisualDNA/FID scores benefit more from the diversity in pretraining data for self-supervised learning (SSL).

MoCoV3					
Network	$\mathbb{D}_{\text{SSL}}$	$\mathcal{C}$ ( $\times 10^6$ )	$N$ ( $\times 10^6$ )	$\mathcal{D}$ ( $\times 10^{-3}$ )	Accuracy↑ ImageNet
ResNet-50	ImageNet	98	0.128	1.31	58.1
	ImageNet	98	0.640	6.54	65.4
	ImageNet	98	1.281	13.0	<b>65.9</b>
	ImageNet	294	0.128	0.43	57.4
	ImageNet	294	0.640	2.17	<b>69.8</b>
	ImageNet	294	1.281	4.35	71.4
ViT-B/16	ImageNet	98	0.128	1.31	57.9
	ImageNet	98	0.640	6.54	67.2
	ImageNet	98	1.281	13.0	<b>68.8</b>
	ImageNet	294	0.128	0.43	56.9
	ImageNet	294	0.640	2.17	<b>71.9</b>
	ImageNet	294	1.281	4.35	74.9

Table 4. **Importance of Normalization.** We report the accuracy of MoCoV3 trained on ImageNet with different data diversity and variable  $\mathcal{C}$ , for the in-distribution assumption ( $\mathbb{D}_{\text{SSL}} = \mathbb{D}_{\text{task}}$ ). For a given network, cells in red highlight inconsistencies where although the model trained with  $\mathcal{C} = 294 \times 10^6$  has seen less samples, it outperforms the best model trained with one third of the computational budget ( $\mathcal{C} = 90 \times 10^6$ ), highlighting the importance of normalization for understanding how models exploit data.

**Model Saturation** In Section 5.3, we evaluated the effects of extreme differences in  $\mathcal{D}$ . We aim to understand the trend in pretraining on YFCC100M, and if adding more samples could compensate the distribution shift with ImageNet. Hence, we extend the YFCC100M experiments from Section 5.3, examining various subsets—specifically 0.1%, 1%, 10%, and 100% of the dataset. Again, we normalize the computational budget to  $\mathcal{C} = 98 \times 10^6$ , equivalent to one epoch on YFCC100M. The results of linear probing models pretrained using SimCLR and MoCoV3 with ResNet-50

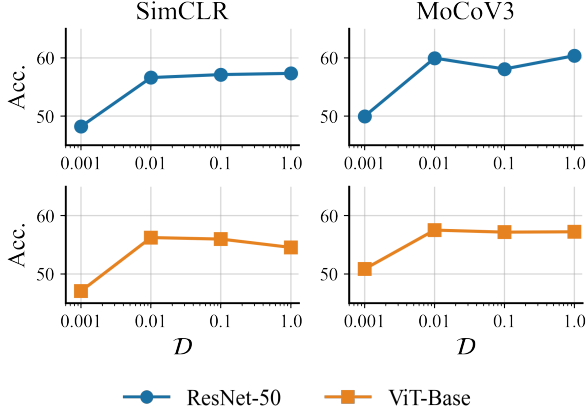


Figure 4. **Data Diversity Impact on YFCC100M Pretraining Performance:** Pretraining ( $\mathcal{C} = 98 \times 10^6$ ) of networks with  $\mathbb{D}_{\text{SSL}} = \text{YFCC100M}$  and  $\mathbb{D}_{\text{task}} = \text{ImageNet}$  for several dataset subsets. In the presence of a distribution shift, performance tends to saturate and does not benefit from additional data diversity.

and ViT-B/16 on ImageNet are shown in Figure 4. Interestingly, we obtain a performance plateau. This observation points to a saturation point in performance, showing that simply increasing the dataset further would not bridge the gap between the  $\mathbb{D}_{\text{SSL}}$  and  $\mathbb{D}_{\text{task}}$ . Hence, arbitrarily increasing pretraining diversity  $\mathcal{D}$  is not sufficient to bridge the distribution shift problem. It is worth noting that  $\mathbb{D}_{\text{SSL}} \neq \mathbb{D}_{\text{task}}$  here, hence still aligning with our findings in Section 5.1.

**Label Quantity** A question that arises in our setting is whether having a higher pretraining diversity leads to requiring less labeled samples during linear probing. In this section, we focus on understanding the impact of increased  $\mathcal{D}$  on the number of labeled samples required for effective linear probing in downstream tasks. We use the same trained models from Section 5.3, *i.e.*, SimCLR and MoCoV3, with ResNet-50 and ViT-B/16 architectures. In this setup, we set  $\mathbb{D}_{\text{task}} = \text{ImageNet}$ , with two upstream dataset scenarios: ImageNet (where  $\mathbb{D}_{\text{SSL}} = \mathbb{D}_{\text{task}}$ ) and YFCC100M (where  $\mathbb{D}_{\text{SSL}} \neq \mathbb{D}_{\text{task}}$ ). Our experiments are summarized in Table 5. We note that in the ViT-B/16 experiments, if ImageNet is used for pretraining, linear probing with just 5% labeled samples can surpass the performance achieved using 100% labeled data when YFCC100M serves as  $\mathbb{D}_{\text{SSL}}$ . This also implies that when  $\mathbb{D}_{\text{SSL}} \neq \mathbb{D}_{\text{task}}$ , a higher quantity of labeled data is necessary to attain competitive performance, and increasing  $\mathcal{D}$  does not bring considerable benefits on label quantity under distribution shift. This implies that that our findings in Section 5.3 hold regardless of the linear probing labeled set. We note that in scenarios where  $\mathbb{D}_{\text{SSL}} = \mathbb{D}_{\text{task}}$ , using only 50% of the labeled data can achieve similar performance as utilizing the full 100% of labeled samples, implying that increasing  $\mathcal{D}$  leads to reduced label requirement efforts for downstream tasks.

SimCLR							
Network	$\mathbb{D}_{\text{SSL}}$	$N$ ( $\times 10^6$ )	$\mathcal{D}$ ( $\times 10^{-3}$ )	Accuracy $\uparrow$			
ResNet-50	ImageNet	1.28	1.31	<b>50.3</b>	<b>54.2</b>	<b>61.7</b>	<b>64.5</b>
	YFCC	98.2	1000	39.4	44.3	54.0	57.3
ViT-B/16	ImageNet	1.28	1.31	<b>55.2</b>	<b>59.4</b>	<b>65.6</b>	<b>66.7</b>
	YFCC	98.2	1000	40.7	45.7	53.3	54.5

MoCoV3							
Network	$\mathbb{D}_{\text{SSL}}$	$N$ ( $\times 10^6$ )	$\mathcal{D}$ ( $\times 10^{-3}$ )	Accuracy $\uparrow$			
ResNet-50	ImageNet	1.28	1.31	<b>52.1</b>	<b>56.2</b>	<b>63.8</b>	<b>65.9</b>
	YFCC	98.2	1000	42.8	48.2	57.7	60.4
ViT-B/16	ImageNet	1.28	1.31	<b>59.3</b>	<b>63.0</b>	<b>68.3</b>	<b>68.8</b>
	YFCC	98.2	1000	43.9	49.0	56.0	57.2

$\mathcal{C} = 98 \times 10^6$

Table 5. **Evaluating Network Accuracy With Varied Label Quantity:** We evaluate the accuracy of networks trained on ImageNet and YFCC with different labeling percentages of  $\mathbb{D}_{\text{task}} = \text{ImageNet}$ . The increased diversity still does not compensate for the distribution shift. However, for in-distribution data, one can get away with fewer labels with more diverse pretraining data.

## 7. Discussion

In this section, we interpret our findings and discuss their implication for practical pretraining aspects towards good practices for future works on SSL.

**Training Takeaways** Coherently with our findings in Section 5.1 and similarly to prior art [15, 24, 34], we find that aligned distributions benefit performance, in particular increasing  $\mathbb{D}_{\text{SSL}}$  diversity helps in learning better feature extractors as long as the distribution of the new samples match those of the downstream task data. Differently from the state-of-the-art, we demonstrate that this holds also in a computationally-normalized setting, implying that collecting large scale in-distribution data matching the downstream task could be an effective and efficient approach to improving SSL. This is useful, as it also allows for learning downstream tasks with fewer labels.

**Evaluation Takeaways** Our analysis reveals that to permit a fair evaluation of SSL methods, computationally normalized tests are necessary to avoid inconsistencies, as revealed in Section 6. Moreover, it is crucial to identify out-of-distribution downstream tasks for a correct evaluation of generalization performance. By evaluating only on  $\mathbb{D}_{\text{task}}$  with a low distribution shift, there is a risk of reporting inflated performance, not representative of a real gain in generalization. This is particularly the case, as new SSL approaches may be reporting higher downstream performance when pretrained on a different dataset. This is since, as per Sections 5 and 6, increasing the computation and the in-distribution data, respectively, can improve performance. This may result in incorrectly concluding that the SSL algorithm is performing better.

**Future Research** In Section 5.3 we showed that even very



diverse datasets, such as YFCC100M, fall short in satisfactory generalization performance. These findings pave the way for further exploration into efficient domain generalization for SSL pretraining. Furthermore, they open doors to investigating why language models enjoy enhanced generalization when exposed to a diverse SSL pretraining diversity compared to vision models. This relates also with scaling performance of vision-language models [61].

## References

- [1] Philip Bachman, R Devon Hjelm, and William Buchwalter. Learning representations by maximizing mutual information across views. In *NeurIPS*, 2019. 2
- [2] Randall Balestriero, Mark Ibrahim, Vlad Sobal, Ari Morcos, Shashank Shekhar, Tom Goldstein, Florian Bordes, Adrien Bardes, Gregoire Mialon, Yuandong Tian, et al. A cookbook of self-supervised learning. *arXiv preprint arXiv:2304.12210*, 2023. 3
- [3] Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. BEiT: BERT pre-training of image transformers. In *ICML*, 2022. 2
- [4] Adrien Bardes, Jean Ponce, and Yann LeCun. VICReg: Variance-invariance-covariance regularization for self-supervised learning. In *ICLR*, 2022. 2, 3, 4
- [5] Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101—Mining discriminative components with random forests. In *ECCV*, 2014. 2, 6
- [6] Mathilde Caron, Piotr Bojanowski, Armand Joulin, and Matthijs Douze. Deep clustering for unsupervised learning of visual features. In *ECCV*, 2018. 2
- [7] Mathilde Caron, Piotr Bojanowski, Julien Mairal, and Armand Joulin. Unsupervised pre-training of image features on non-curated data. In *ICCV*, 2019. 2
- [8] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. In *NeurIPS*, 2020. 2, 3, 4
- [9] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *ICCV*, 2021. 2, 3, 4
- [10] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *ICML*, 2020. 2, 3, 4, 1
- [11] Ting Chen, Simon Kornblith, Kevin Swersky, Mohammad Norouzi, and Geoffrey Hinton. Big self-supervised models are strong semi-supervised learners. In *NeurIPS*, 2020. 2
- [12] Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *CVPR*, 2021. 2
- [13] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020. 2
- [14] Xinlei Chen, Saining Xie, and Kaiming He. An empirical study of training self-supervised vision transformers. In *ICCV*, 2021. 2, 3, 4, 1
- [15] Elijah Cole, Xuan Yang, Kimberly Wilber, Oisín Mac Aodha, and Serge Belongie. When does contrastive visual representation learning work? In *CVPR*, 2022. 2, 3, 5, 8
- [16] Victor Guilherme Turrissi da Costa, Enrico Fini, Moin Nabi, Nicu Sebe, and Elisa Ricci. solo-learn: A library of self-supervised methods for visual representation learning. In *JMLR*, 2022. 4, 5
- [17] Mostafa Dehghani, Josip Djolonga, Basil Mustafa, Piotr Padlewski, Jonathan Heek, Justin Gilmer, Andreas Peter Steiner, Mathilde Caron, Robert Geirhos, Ibrahim Alabdulmohsin, et al. Scaling vision transformers to 22 billion parameters. In *ICML*, 2023. 2
- [18] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009. 2, 5, 6
- [19] Carl Doersch and Andrew Zisserman. Multi-task self-supervised visual learning. In *ICCV*, 2017. 2
- [20] Carl Doersch, Abhinav Gupta, and Alexei A Efros. Unsupervised visual representation learning by context prediction. In *ICCV*, 2015. 2
- [21] 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*, 2021. 2, 3, 6
- [22] Linus Ericsson, Henry Gouk, and Timothy M Hospedales. How well do self-supervised models transfer? In *CVPR*, 2021. 3
- [23] Spyros Gidaris, Praveer Singh, and Nikos Komodakis. Unsupervised representation learning by predicting image rotations. In *ICLR*, 2018. 2
- [24] Priya Goyal, Dhruv Mahajan, Abhinav Gupta, and Ishan Misra. Scaling and benchmarking self-supervised visual representation learning. In *ICCV*, 2019. 2, 8
- [25] Priya Goyal, Mathilde Caron, Benjamin Lefaudeux, Min Xu, Pengchao Wang, Vivek Pai, Mannat Singh, Vitaliy Liptchinsky, Ishan Misra, Armand Joulin, et al. Self-supervised pretraining of visual features in the wild. *arXiv preprint arXiv:2103.01988*, 2021. 1, 2, 3
- [26] Priya Goyal, Quentin Duval, Isaac Seessel, Mathilde Caron, Ishan Misra, Levent Sagun, Armand Joulin, and Piotr Bojanowski. Vision models are more robust and fair when pretrained on uncurated images without supervision. *arXiv preprint arXiv:2202.08360*, 2022. 1, 2
- [27] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent: A new approach to self-supervised learning. In *NeurIPS*, 2020. 2, 3, 4
- [28] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 3, 4, 6
- [29] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *CVPR*, 2020. 2, 3
- [30] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *CVPR*, 2022. 2, 3

- [31] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *NeurIPS*, 30, 2017. 7, 5
- [32] R Devon Hjelm, Alex Fedorov, Samuel Lavoie-Marchildon, Karan Grewal, Phil Bachman, Adam Trischler, and Yoshua Bengio. Learning deep representations by mutual information estimation and maximization. In *ICLR*, 2018. 2
- [33] Olivier J. Hénaff, Aravind Srinivas, Jeffrey De Fauw, Ali Razavi, Carl Doersch, S. M. Ali Eslami, and Aaron van den Oord. Data-efficient image recognition with contrastive predictive coding. In *ICML*, 2019. 2
- [34] Klemen Kotar, Gabriel Ilharco, Ludwig Schmidt, Kiana Ehsani, and Roozbeh Mottaghi. Contrasting contrastive self-supervised representation learning pipelines. In *ICCV*, 2021. 2, 8
- [35] Jonathan Krause, Jia Deng, Michael Stark, and Li Fei-Fei. Collecting a large-scale dataset of fine-grained cars. In *FGVC*, 2013. 2, 6
- [36] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. 4
- [37] Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. *CS 231N*, 2015. 2, 4
- [38] Alexander Cong Li, Ellis Langham Brown, Alexei A Efros, and Deepak Pathak. Internet explorer: Targeted representation learning on the open web. In *ICML*, 2023. 2
- [39] Junnan Li, Pan Zhou, Caiming Xiong, and Steven C. H. Hoi. Prototypical contrastive learning of unsupervised representations. In *ICLR*, 2021. 2
- [40] Ishan Misra and Laurens van der Maaten. Self-supervised learning of pretext-invariant representations. In *CVPR*, 2020. 2
- [41] Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number of classes. In *ICVGIP*, 2008. 2, 6
- [42] Mehdi Noroozi and Paolo Favaro. Unsupervised learning of visual representations by solving jigsaw puzzles. In *ECCV*, 2016. 2
- [43] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. 2, 4
- [44] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. DINOv2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023. 1, 7
- [45] Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, and CV Jawahar. Cats and dogs. In *CVPR*, 2012. 2, 6
- [46] Ilija Radosavovic, Raj Prateek Kosaraju, Ross Girshick, Kaiming He, and Piotr Dollár. Designing network design spaces. In *CVPR*, 2020. 2
- [47] Benjamin Ramtoula, Matthew Gadd, Paul Newman, and Daniele De Martini. Visual DNA: Representing and comparing images using distributions of neuron activations. In *CVPR*, 2023. 7, 5
- [48] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022. 5
- [49] Yuge Shi, Imant Daunhawer, Julia E Vogt, Philip Torr, and Amartya Sanyal. How robust is unsupervised representation learning to distribution shift? In *ICLR*, 2022. 2
- [50] Chen Sun, Abhinav Shrivastava, Saurabh Singh, and Abhinav Gupta. Revisiting unreasonable effectiveness of data in deep learning era. In *ICCV*, 2017. 1
- [51] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *CVPR*, 2016. 7
- [52] Bart Thomee, David A. Shamma, Gerald Friedland, Benjamin Elizalde, Karl Ni, Douglas Poland, Damian Borth, and Li-Jia Li. YFCC100M: The new data in multimedia research. In *Commun. ACM*, 2016. 2, 6
- [53] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive multiview coding. In *ECCV*, 2020. 2, 5
- [54] Yonglong Tian, Olivier J. Henaff, and Aaron Van Den Oord. Divide and contrast: Self-supervised learning from uncured data. In *ICCV*, 2021. 2
- [55] Antonio Torralba and Alexei A Efros. Unbiased look at dataset bias. In *CVPR*, 2011. 2
- [56] Wouter Van Gansbeke, Simon Vandenhende, Stamatios Georgoulis, and Luc V Gool. Revisiting contrastive methods for unsupervised learning of visual representations. In *NeurIPS*, 2021. 2
- [57] Grant Van Horn, Elijah Cole, Sara Beery, Kimberly Wilber, Serge Belongie, and Oisín Mac Aodha. Benchmarking representation learning for natural world image collections. In *CVPR*, 2021. 3
- [58] Zhirong Wu, Yuanjun Xiong, Stella X Yu, and Dahua Lin. Unsupervised feature learning via non-parametric instance discrimination. In *CVPR*, 2018. 3
- [59] Zhenda Xie, Zheng Zhang, Yue Cao, Yutong Lin, Jianmin Bao, Zhuliang Yao, Qi Dai, and Han Hu. SimMIM: A simple framework for masked image modeling. In *CVPR*, 2022. 2
- [60] Zhenda Xie, Zheng Zhang, Yue Cao, Yutong Lin, Yixuan Wei, Qi Dai, and Han Hu. On data scaling in masked image modeling. In *CVPR*, 2023. 2
- [61] 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. *arXiv preprint arXiv:2309.16671*, 2023. 9
- [62] Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and Stéphane Deny. Barlow Twins: Self-supervised learning via redundancy reduction. In *ICML*, 2021. 2, 3, 4
- [63] Xiaohua Zhai, Alexander Kolesnikov, Neil Houlsby, and Lucas Beyer. Scaling vision transformers. In *CVPR*, 2022. 2
- [64] Richard Zhang, Phillip Isola, and Alexei A. Efros. Colorful image colorization. In *ECCV*, 2016. 2, 3
- [65] Nanxuan Zhao, Zhirong Wu, Rynson W. H. Lau, and Stephen Lin. What makes instance discrimination good for transfer learning? In *ICML*, 2021. 2
- [66] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. *T-PAMI*, 2017. 2, 6
- [67] Pan Zhou, Yichen Zhou, Chenyang Si, Weihao Yu, Teck Khim Ng, and Shuicheng Yan. Mugs: A multi-granular self-supervised learning framework. *arXiv preprint arXiv:2203.14415*, 2022. 7

- [68] Chengxu Zhuang, Alex Lin Zhai, and Daniel Yamins. Local aggregation for unsupervised learning of visual embeddings. In *ICCV*, 2019. 3

# On Pretraining Data Diversity for Self-Supervised Learning

## Supplementary Material

In this supplementary material, we present additional experiments and insights on our findings presented in the main paper. First, we further develop the inconsistencies found within an incorrectly-normalized framework (Section A). Then, we propose different settings for our analysis at scale (Section B). We extend our analysis on label quantity in Section C. Finally, we introduce additional details about our settings and implementations (Section D). For reproducibility of our results, we will share on GitHub our codebase, along the fine-tuned parameters and the data ordering files.

For improved readability, all references to the main paper Sections and Tables in blue (e.g. “Section 1”).

### A. Importance of Normalization of Computation

In this section, we aim to complement our experiment in Section 6, providing further proof highlighting the importance of having a normalized computational budget when evaluating the performance of SSL methods. In the following experiments, we show that if computation is not normalized properly, one might fall into unfair comparisons and misleading conclusions.

#### A.1. Increasing Total Computation

In our first setup, we pretrain SimCLR [10] and MoCoV3 [14] on Tiny ImageNet with various  $\mathcal{D}$ , over a range of increasing amounts of budget  $\mathcal{C}$ . We assume  $\mathbb{D}_{\text{SSL}} = \mathbb{D}_{\text{task}}$ . We take the same subsets of Tiny ImageNet as in Section 5.1, which consist of 10%, 50% and 100% of the training data. We vary  $\mathcal{C}$  from  $5 \times 10^6$  to  $100 \times 10^6$ , and we measure the accuracy of the pretrained models on the  $\mathbb{D}_{\text{task}}$  test set, results of which are shown in Table 6. Note that we refer to models in this section using the data diversity  $N$  instead of pretraining diversity  $\mathcal{D}$ , as each cell in Table 6 has a different  $\mathcal{D}$ . Following our previous experiments, we argue that comparison between different diversities only hold as long as the computation is normalized. This implies that comparisons only hold within any given  $\mathcal{C}$ . In agreement with our prior findings, the third row with the highest diversity always outperforms the lower pretraining diversity models on in-domain evaluation, for both SimCLR and MoCoV3. However, only when we compare between different columns, *i.e.* different amounts of computation, we may observe that models pretrained with lower diversity may outperform higher diversity models. For example, for both SimCLR and MoCoV3, the models pretrained with  $N = 50 \times 10^3$  and  $\mathcal{C} = 10 \times 10^6$  outperform the models with higher data diversities  $N = 100 \times 10^3$  but less computation  $\mathcal{C} = 5 \times 10^6$ . As a result, we see that models with lower pretraining diversity can still outperform models with higher diversity, given that more computation is used. This highlights the importance of normalizing computation costs when evaluating the effects of diversity.

SimCLR						MoCoV3					
$N$ ( $\times 10^3$ )	$\mathcal{C}$ ( $\times 10^6$ )					$N$ ( $\times 10^3$ )	$\mathcal{C}$ ( $\times 10^6$ )				
	5	10	25	50	100		5	10	25	50	100
10	36.92	36.63	36.30	36.91	35.03	10	39.78	41.82	40.06	36.56	28.92
50	40.76	44.30	47.69	48.77	48.91	50	39.88	43.42	46.68	46.45	48.14
100	<b>41.43</b>	<b>44.76</b>	<b>49.32</b>	<b>49.83</b>	<b>51.62</b>	100	<b>40.35</b>	<b>44.03</b>	<b>47.63</b>	<b>48.58</b>	<b>50.71</b>

Table 6. **Pretraining Diversity With Increasing Computational Budget:** We show for both SimCLR (left) and MoCoV3 (right) that increasing pretraining diversity always leads to better in-domain downstream accuracies, given that computation is normalized, *i.e.*, comparisons hold within the columns of the tables. Comparing models between different columns may lead to inconsistencies, where lower diversity models with more computation obtain higher results than higher diversity models with less computation.

#### A.2. Epoch-based normalization

In Section 5.2, we adhered to a fixed computational budget of  $\mathcal{C} = 50 \times 10^6$ , pretraining models on  $\mathbb{D}_{\text{SSL}} = \mathbb{B}$  for 800 epochs, and on  $\mathbb{B} \cup \mathbb{A}$  for 400 epochs, considering that the latter dataset was twice the size. We further demonstrate the importance of a computationally-normalized evaluation by exposing the inconsistencies of an alternative epoch-based normalization, hence in which networks are trained for 400 epochs, regardless of the dataset size.

We propose this alternative scenario in Figure 5, where the compute-normalized baseline (the black dashed line in Figure 3) is replaced with an epoch-normalized baseline (indicated by the red dashed line), obtained by pretraining for 400 epochs



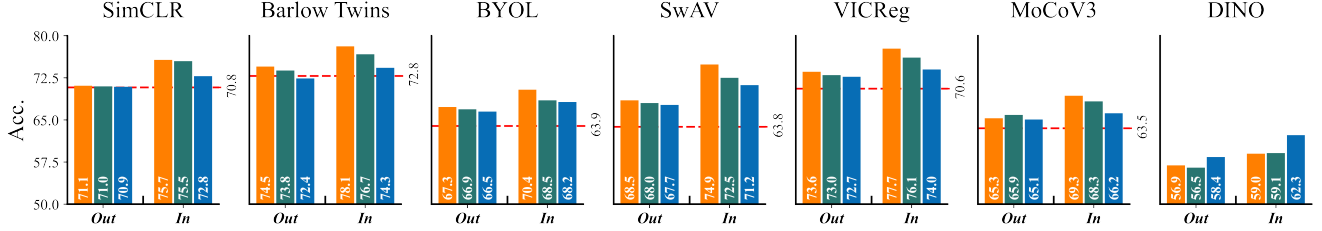


Figure 5. **Impact of Epoch Normalization on SSL Pretraining Performance:** This figure contrasts an epoch-normalized baseline (red line) with the trained methods in the main paper, Figure 3. Under epoch normalization, we notice contrasting findings, *i.e.* more diverse trainings, irrespective of their origin (source, web, or synthetic) and label distribution (in or out-of-distribution), consistently enhances performance. This is an unfair comparison due to the greater costs of each augmented pretraining if epochs are normalized. This illustrates how alternative normalization can lead to wrong conclusions compared to compute normalization. DINO’s epoch-normalized baseline is shown in text only (Acc. 41.14) for ease of visualization.

Barlow Twins								
$\mathbb{D}_{SSL}$	$N$ ( $\times 10^6$ )	$\mathcal{D}$ ( $\times 10^{-3}$ )	Accuracy $\uparrow$					
			ImageNet	Cars	Flow.	Pets	Places	Food
ImageNet	0.128	1.31	57.17	51.51	85.84	75.81	45.87	63.72
ImageNet	1.281	13.0	<b>65.40</b>	<b>59.56</b>	89.16	<b>83.89</b>	49.19	<b>70.42</b>
YFCC	98.17	1000	57.85	44.76	83.46	67.21	<b>50.01</b>	65.20
$\mathcal{C} = 98 \times 10^6$								

BYOL								
$\mathbb{D}_{SSL}$	$N$ ( $\times 10^6$ )	$\mathcal{D}$ ( $\times 10^{-3}$ )	Accuracy $\uparrow$					
			ImageNet	Cars	Flow.	Pets	Places	Food
ImageNet	0.128	1.31	61.82	46.62	85.84	80.28	46.91	67.01
ImageNet	1.281	13.0	<b>68.39</b>	<b>52.51</b>	<b>88.77</b>	<b>84.75</b>	49.96	<b>73.52</b>
YFCC	98.17	1000	60.73	42.78	84.38	68.63	<b>50.92</b>	68.63
$\mathcal{C} = 98 \times 10^6$								

$\mathcal{C} = 98 \times 10^6$

$\mathcal{C} = 98 \times 10^6$

Table 7. **Non-contrastive pretraining:** We explore two more pretraining methods, namely, Barlow Twins and BYOL, for our large-scale pretraining experiments. Again, the budget is set to  $\mathcal{C} = 98 \times 10^6$ , we find that our earlier conclusions still hold here: (1) ImageNet pretraining outperforms YFCC100M (YFCC) pretraining except for Places365 due to the distribution shift (2) Increased pretraining diversity  $\mathcal{D}$  generally correlates with improved downstream performance with the exception. Those findings are consistent for both Barlow Twins and BYOL.

on  $\mathbb{B}$ . Here, we observe that augmenting with additional samples consistently enhances performance, irrespective of the data augmentation technique used and whether the sample labels are in or out of distribution. This finding does not align with the insights from Section 5.2, and we highlight that it does not take into account the difference in costs for training models for the same number of epochs, but on a dataset twice the size. Hence, this constitutes an unfair comparison that may lead to incorrect conclusions, advocating for the effectiveness of our computational-based normalization.

## B. Alternative Settings

### B.1. Non-Contrastive Methods

For large-scale experiments in section 5.3 we only considered SimCLR [10] and MoCoV3 [14], both of which are contrastive SSL methods. Here we show that the results are consistent for the non-contrastive methods Barlow Twins [62] and BYOL [27]. We highlight that these experiments are computationally intensive, hence we explore a reduced setting with a single backbone and lower  $\mathcal{D}$  variability. We pretrained a ResNet-50 backbone using Barlow Twins and BYOL on ImageNet and the same subsets from Section 5.3, as well as on the full YFCC100M dataset, ensuring that the total compute is equal to a single epoch on YFCC100M, *i.e.*  $\mathcal{C} = 98 \times 10^6$ . Again, we show linear evaluation on multiple downstream datasets including ImageNet [18], Stanford Cars [35], Flowers-102 [41], Oxford-IIIT Pets [45], Places365 [66], and Food-101 [5] in Table 7. In accordance with Section 5.3, we observe that pretraining on higher diversities leads to improved downstream accuracies when  $\mathbb{D}_{SSL} = \mathbb{D}_{task}$ , *i.e.* pretraining and evaluating on ImageNet. Also, the highest pretraining diversity in ImageNet leads to the best results for all downstream datasets, except for Places365, for which pretraining YFCC100M performs best, for which we refer again to distribution distances analysis in Section 6. Again for these methods, the maximum diversity of YFCC100M is not enough to diminish the effects of the domain gap between the pretraining data and the datasets other than Places365.

MoCoV3							
$\mathbb{D}_{\text{SSL}}$	$\mathcal{C}$ ( $\times 10^6$ )	$N$ ( $\times 10^6$ )	$\mathcal{D}$ ( $\times 10^{-3}$ )	Network			
				ResNet-50	ResNet-101	ViT-S/16	ViT-B/16
ImageNet	98	0.128	1.31	58.1	58.9	56.3	57.9
ImageNet	98	0.640	6.54	65.4	67.2	64.7	67.2
ImageNet	98	1.281	13.0	<b>65.9</b>	<b>67.7</b>	<b>65.4</b>	<b>68.8</b>
ImageNet	294	0.128	0.43	57.5	59.0	52.9	56.9
ImageNet	294	0.640	2.17	69.8	71.4	68.9	71.9
ImageNet	294	1.281	4.35	<b>71.4</b>	<b>73.3</b>	<b>71.4</b>	<b>74.9</b>

Table 8. **Pretraining Diversity With Different Architectures Sizes:** We investigate how pretraining diversity, total computation budget and model architecture size interact for MoCoV3 when pretraining and evaluating on ImageNet. Regardless of  $\mathcal{C}$  and the architecture choice, increasing pretraining diversity remains a reliable method to improve downstream results. Further, increasing model size also seems to consistently lead to better learned representations. Again, comparing pretraining diversity values only holds when the model architecture and  $\mathcal{C}$  are fixed.

## B.2. Different Architectures

We investigate how pretraining diversity interacts with varying backbone architecture sizes, as well as the total computational budget. With this, we aim to highlight how different models react to pretraining diversity. To benchmark the interaction of these factors, we focus on MoCoV3 and pretrain and evaluate on ImageNet using  $\mathcal{C} = 98 \times 10^6$  and the tripled amount  $\mathcal{C} = 294 \times 10^6$ . We use two different architecture sizes for ViT backbones as well as ResNet backbones: ViT-Small/16 paired with ViT-Base/16 and ResNet-50 paired with ResNet-101.

Results are shown in Table 8, and the first observation we make, is that for any combination of architecture size and total amount of computation, the model pretrained with largest amount of pretraining diversity  $\mathcal{D} = 13.0 \times 10^{-3}$  always has the highest in-domain downstream performance. Increasing pretraining diversity thus remains a reliable method to improve the quality of learned representations, regardless of the architecture size. Secondly, we see that for every diversity value, regardless of the backbone type or the amount of computation, an increase in backbone size, *i.e.* ResNet-50 to ResNet-101 or ViT-S/16 to ViT-B/16, leads to an increased performance. It is thus again of importance to only compare models with different pretraining diversities for fixed model size, as we did with fixed computational budget.

Finally, keep in mind that the larger architectures require more computation, which is not incorporated in  $\mathcal{C}$  as this term only describes the number of images that are seen during pretraining.

**A note on MAE** Masked Autoencoders (MAE) [30] is a Transformer-specific pretraining method based on an autoregressive loss. This differs considerably from what has been presented in Section 3, and it has significant impact on the components needed for our normalized analysis. Indeed, for a  $\mathcal{C} = 98 \times 10^6$  budget, MAE is far from providing optimal performance [30], making comparisons unfair without incurring in unsustainable costs. Also, the reconstruction task used for supervision extracts features requiring a deep decoder for best performance in linear probing [30], and it results in considerably better performance with full finetuning exclusively. We will analyze pretraining diversity effects for MAE in a future work.

## C. Additional Insights on Label Quantity

In Section 6 we considered how pretraining diversity affects the number of labels necessary for the best downstream ImageNet accuracies when pretrained on ImageNet ( $\mathbb{D}_{\text{SSL}} = \mathbb{D}_{\text{task}}$ ) or on YFCC100M ( $\mathbb{D}_{\text{SSL}} \neq \mathbb{D}_{\text{task}}$ ). Here explore the setting where upstream and downstream data are the same, *i.e.*,  $\mathbb{D}_{\text{SSL}} = \mathbb{D}_{\text{task}}$ , and we repeat the experiment with models pretrained on various diversities on ImageNet. Table 9 shows the in-domain results on ImageNet for SimCLR and MoCoV3 pretrained with ResNet-50 and ViT-B/16 backbones. It is clear that for in-domain evaluation, the models pretrained with largest pretraining diversity always perform the best, regardless of the label quantity used for linear evaluation. More interestingly, it is possible to achieve better performance with less labels if a model is pretrained with higher  $\mathcal{D}$ . For example, for every combination of backbone and SSL method, the models pretrained with maximum diversity  $\mathcal{D} = 13.0 \times 10^{-3}$  using 50% of the labels outperform the models pretrained with  $\mathcal{D} = 2.61 \times 10^{-3}$  with 100% of the labels. Thus, if models are evaluated or deployed in few-shot downstream tasks, it may be desirable to use models pretrained with the highest pretraining diversity available.

SimCLR								MoCoV3							
Network	$\mathbb{D}_{SSL}$	$N$	$\mathcal{D}$	Accuracy $\uparrow$				Network	$\mathbb{D}_{SSL}$	$N$	$\mathcal{D}$	Accuracy $\uparrow$			
		$(\times 10^6)$	$(\times 10^{-3})$	5%	10%	50%	100%			$(\times 10^6)$	$(\times 10^{-3})$	5%	10%	50%	100%
ResNet-50	ImageNet	0.128	1.31	42.1	45.4	53.2	56.9	ResNet-50	ImageNet	0.128	1.31	43.9	47.4	55.2	58.1
	ImageNet	0.256	2.61	46.8	50.3	57.8	61.1		ImageNet	0.256	2.61	49.0	52.7	60.4	62.9
	ImageNet	0.640	6.54	49.1	53.0	60.6	63.7		ImageNet	0.640	6.54	51.3	55.4	63.1	65.4
	ImageNet	1.281	13.0	<b>50.3</b>	<b>54.2</b>	<b>61.7</b>	<b>64.5</b>		ImageNet	1.281	13.0	<b>52.1</b>	<b>56.2</b>	<b>63.8</b>	<b>65.9</b>
ViT-B/16	ImageNet	0.128	1.31	40.5	44.8	53.0	54.2	ViT-B/16	ImageNet	0.128	1.31	47.2	51.3	58.0	57.9
	ImageNet	0.256	2.61	48.0	52.2	59.6	61.3		ImageNet	0.256	2.61	53.5	57.6	63.1	63.7
	ImageNet	0.640	6.54	53.4	57.6	64.3	65.5		ImageNet	0.640	6.54	57.9	61.7	66.7	67.2
	ImageNet	1.281	13.0	<b>55.2</b>	<b>59.4</b>	<b>65.6</b>	<b>66.7</b>		ImageNet	1.281	13.0	<b>59.3</b>	<b>63.0</b>	<b>68.3</b>	<b>68.8</b>

Table 9. **Evaluating Network Accuracy With Varied Label Quantity:** We evaluate the accuracy of networks pretrained on ImageNet with various pretraining diversities and evaluate with different labeling percentages of  $\mathbb{D}_{task} = \text{ImageNet}$ . For in-distribution data, one can get away with fewer labels using more diverse pretraining data.

## D. Additional details

### D.1. SSL Methods

In Section 3 we described a general framework for self-supervised pretraining that is common to many state-of-the-art SSL methods. Although all the methods we use in our experiments mostly follow this procedure, they do differ in loss functions as well as in certain architectural choices. For each of the methods we use for our experiments, we describe in depth the key aspects that specifically define the SSL method and make them different from the introduced framework in Section 3. Further details on the methods can be found in their respective papers and repositories.

SimCLR [10], Barlow Twins [62] and VICReg [4] closely follow the general framework, and mainly differ in the loss function  $\mathcal{L}_{SSL}$  used during optimization. SimCLR makes use of the InfoNCE loss [43], which is applied to the representations of each positive pair of samples in the batch, and also incorporates negatives from the current batch. Barlow Twins uses a loss function that makes the cross-correlation matrix between the representations of the distorted samples as close to the identity matrix as possible. As a result, representations of correlated views of the same sample are forced to become similar, while redundancy between the components of these vectors is minimized. The loss function used for VICReg is a combination of the mean-squared euclidean distance between representations with an additional variance and covariance term for feature decorrelation and avoiding representation collapse.

BYOL [27], DINO [9] and MoCoV3 [14] have slightly more evident differences from the proposed framework, as they do not use shared parameters  $\theta_f$  and  $\theta_g$  for the feature extractor and projection head between the two augmented views. Instead the two augmented views pass through two different networks: a student network with feature extractor  $f_{\theta_f}$  and prediction head  $g_{\theta_g}$ , parameterised with  $\theta_f$  and  $\theta_g$ , and a teacher network with its own respective components  $f'_{\theta_{f'}}$  and  $g'_{\theta_{g'}}$ , with unique parameters  $\theta_{f'}$  and  $\theta_{g'}$ . The weights  $\theta_{f'}$  and  $\theta_{g'}$  in the teacher network are an exponential moving average of the weights  $\theta_f$  and  $\theta_g$  in the student network. The three methods differ in how they compute  $\mathcal{L}_{SSL}$  from the representations  $\mathbf{z}_A$  from the student and  $\mathbf{z}_B$  from the teacher. For BYOL, after the correlated views are passed through the two networks, an additional predictor network  $q_{\theta_q}$ , parameterised with  $\theta_q$ , tries to predict the representation of the teacher network  $\mathbf{z}_B$  from the output of the student network as  $q_{\theta_q}(\mathbf{z}_A)$ , and the mean squared error between the teacher representation and the student prediction is minimised. DINO performs knowledge distillation in the student by minimising the cross-entropy loss between the direct outputs  $\mathbf{z}_A$  and  $\mathbf{z}_B$ . MoCoV3 uses the student and teacher network to generate representations from the augmented views called *queries*  $\mathbf{z}_B$  and *keys*  $\mathbf{z}_B$ , and stores the keys in a queue. The contrastive InfoNCE loss is again used as SSL objective, and uses matching queries and keys as positive samples and the recent keys from the dictionary as negative samples. For all three methods, a stop-gradient operator is applied after the teacher network, to avoid any weight updates in the teacher network during backpropagation.

SwAV [8] does share weights for  $f_{\theta_f}$  and  $g_{\theta_g}$  between correlated views, but instead relies on additional components. First, the representations of different views  $\mathbf{z}_A$  and  $\mathbf{z}_B$  are assigned to prototype vectors, resulting in *codes*  $\mathbf{q}_A$  and  $\mathbf{q}_B$ . The prototype vectors are trainable vectors and represent the notion of clusters. A swapped prediction problem is solved afterwards, where the code of one augmented view is predicted using the other view. The swapped prediction problem is used to define the loss as  $\mathcal{L}_{SSL}(\mathbf{z}_A, \mathbf{z}_B) = \ell(\mathbf{z}_A, \mathbf{q}_B) + \ell(\mathbf{z}_B, \mathbf{q}_A)$ , where  $\ell$  measures the fit between features and codes.

## D.2. Data Collection Strategies

This section outlines the data collection strategies for our three approaches detailed in Section 5.2: Source, Web, and Synthetic. We base these strategies on the Base dataset  $\mathbb{B}$  (introduced in 5.2), consisting of half of ImageNet100, totaling 65,000 samples.

**Source Dataset** We expand the dataset  $\mathbb{B}$  by integrating the remaining half of ImageNet100, forming  $\mathbb{A}_{\text{Source}}^{\text{In}}$ . For  $\mathbb{A}_{\text{Source}}^{\text{Out}}$ , we begin by selecting 100 random, non-overlapping classes from ImageNet. We then gather 65,000 corresponding samples from these classes and add them to  $\mathbb{B}$ .

**Web Dataset** We utilize three search engines—Flickr, Bing, and DuckDuckGo—to gather web samples while employing Safe-Search for content appropriateness. Our queries, based on class names, are carefully crafted to avoid ambiguity. For  $\mathbb{A}_{\text{Web}}^{\text{In}}$ , we collect approximately 100,000 samples from ImageNet100 classes, selecting the top 65,000 for inclusion in  $\mathbb{B}$ . Similarly, for  $\mathbb{A}_{\text{Web}}^{\text{Out}}$ , we follow the same process for the 100 randomly selected classes from the Source dataset.

**Synthetic Dataset** For synthetic sample generation, we employ Stable Diffusion V2.1 (SDV2.1). Using the prompt 'A photo of a `class_name`', we generate images for each class in ImageNet100 for  $\mathbb{A}_{\text{Synthetic}}^{\text{In}}$  and the 100 distinct classes from ImageNet for  $\mathbb{A}_{\text{Synthetic}}^{\text{Out}}$ . Each class contributes 650 images, totaling 65,000 samples. We utilize the DPMSolver++ scheduler with start and end  $\beta$  values of 0.00085 and 0.012, respectively. The guidance scale is set at  $w = 7.5$ , and each image is generated in 50 steps.

## D.3. Details on Distribution Distances

In Section 6 of our study, we explored the relationship between upstream datasets—specifically ImageNet and YFCC100M—and downstream datasets, which include ImageNet, Stanford Cars, Flowers102, OxfordIIITPets, Places365, and Food101. To quantitatively measure the distance between these datasets, we employed two distinct metrics: VisualDNA (VDNA) and the Fréchet Inception Distance (FID).

Our methodology for calculating distribution distances involved selecting a substantial number of samples from each dataset. We used 50,000 samples from each dataset for computing the distribution distances, if the dataset is composed of less than 50,000 samples, the whole dataset is used. This approach ensured a robust and comprehensive analysis of the dataset distributions. For the implementation of VisualDNA, we utilized two different architectures: MUGS ViT-B/16, as recommended by the original paper [47], and DinoV2 ViT-B/16. The FID scores were computed using a standard approach with an Inception network, as detailed in [31].

The results, as discussed in Section 6, revealed consistent findings across both VDNA and FID metrics. Our analysis showed that a greater distance between the upstream and downstream datasets correlated with a decrease in downstream classification accuracy.

## E. Implementation

In all our experiments, we have utilized the `solo-learn` library [16] as our main codebase. For the ImageNet100 and CIFAR100 experiments, we used the parameters provided by `solo-learn`. In the case of ImageNet, while we began with the parameters provided in the original papers, we made slight modifications to optimize performance. These modifications included changes to the number of warmup epochs and an adjustment of the learning rate. For the YFCC100M dataset, we found the parameters optimized for ImageNet to be the most effective, whereas for TinyImageNet, we used the CIFAR100 parameters provided by `solo-learn`.

To create different fractions of each dataset, our first step involved the creation of an H5 file containing all image paths. This file is then shuffled and saved. When a specific percentage of the data is required for our SSL pretraining, we simply select the first  $k\%$  of the image paths from this H5 file. Since we use a fixed computational budget, we scaled the number of epochs accordingly. This scaling is achieved by a factor of  $1/k \times 100$ . For example, if we utilized 10% of the data for pretraining, we would increase the base number of epochs by a factor of 10.