

Beyond Supervised vs. Unsupervised: Representative Benchmarking and Analysis of Image Representation Learning

Matthew Gwilliam and Abhinav Shrivastava
University of Maryland, College Park

Abstract

By leveraging contrastive learning, clustering, and other pretext tasks, unsupervised methods for learning image representations have reached impressive results on standard benchmarks. The result has been a crowded field – many methods with substantially different implementations yield results that seem nearly identical on popular benchmarks, such as linear evaluation on ImageNet. However, a single result does not tell the whole story. In this paper, we compare methods using performance-based benchmarks such as linear evaluation, nearest neighbor classification, and clustering for several different datasets, demonstrating the lack of a clear front-runner within the current state-of-the-art. In contrast to prior work that performs only supervised vs. unsupervised comparison, we compare several different unsupervised methods against each other. To enrich this comparison, we analyze embeddings with measurements such as uniformity, tolerance, and centered kernel alignment (CKA), and propose two new metrics of our own: nearest neighbor graph similarity and linear prediction overlap. We reveal through our analysis that in isolation, single popular methods should not be treated as though they represent the field as a whole, and that future work ought to consider how to leverage the complimentary nature of these methods. We also leverage CKA to provide a framework to robustly quantify augmentation invariance, and provide a reminder that certain types of invariance will be undesirable for downstream tasks.

1. Introduction

Image features are critical components in many computer vision (CV) pipelines. In this paper, we define image features, also referred to as embeddings, encodings, or representations, as an n -dimensional vector that represents the content of an image. With the emergence of deep learning, classical approaches to computing image features have been supplanted by neural networks that use large amounts of data to generate powerful image representations. The most widespread method is straightforward: a neural net-

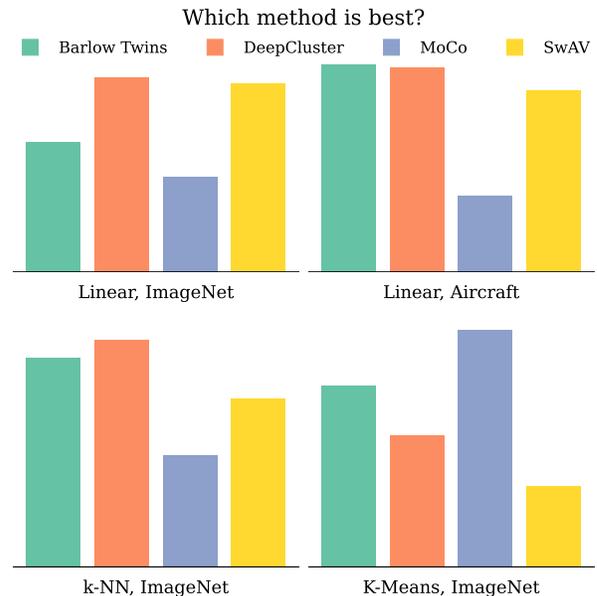


Figure 1. Results for a sample of classification benchmarks we perform in this paper. While these bar charts report real results, lack of axes is intentional – the exact numbers are in Section 4. Importantly, between the four tasks, there is no clear “best” method.

work (e.g., a ResNet50 [30]) is trained to classify the images in some large dataset, typically ImageNet. The portion of the network that performs the classification, usually just the final layer, is then removed, and the outputs of the penultimate layer for a given image are considered the features for that image. This process relies on image classification, a supervised learning task, and thus requires the availability of large amounts of annotated, high-quality data.

Recent successes make unsupervised learning a viable alternative paradigm where image features are learned without the need for class labels. Within unsupervised learning, methods can be considered either generative or discriminative. Generative methods are typically designed for reconstruction or similar tasks [5, 19, 20, 35, 51]. Since we are more concerned with a potential transfer to downstream tasks such as image classification and object detection, we

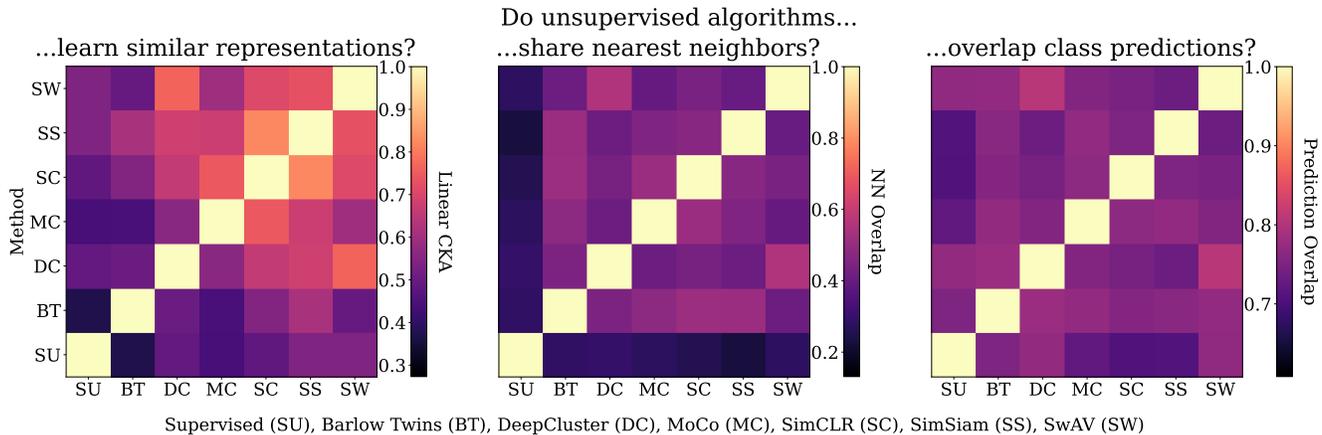


Figure 2. Similarity between learned representations, based on the outputs of ResNet50s on the validation images from ImageNet. For the metrics shown, which are described in more detail in Section 4, higher values indicate similarity. While the supervised model tends to be more dissimilar from the unsupervised models, there are many ways in which unsupervised methods differ substantially from each other.

choose to focus on discriminative methods.

There are many different ways to compare image representation learning algorithms. In this paper, we opt to focus on the role of the methods as feature extractors, where a model that is pre-trained for some task is expected to be able to generate useful features for unseen images. Thus, we only use benchmarks that keep the backbone (the portion of the neural network that generates the embedding) frozen. Prior works are often limited by their focus on a single benchmark, single method, or toy datasets. In contrast, we compare 6 SOTA unsupervised methods on ImageNet and 6 fine-grained visual categorization (FGVC) datasets using several different benchmarks. Figure 1 provides a sample of this angle of analysis.

Comparing these methods to each other is very important. Nevertheless, prior analysis works tend to lump unsupervised methods together, and often choose only a single representative such as MoCo or SimCLR for comparison against supervised representation learning [1, 14, 22, 26, 54]. This ignores the significant ways in which unsupervised methods differ from each other. In contrast to prior work, we extend existing methods and introduce novel methods to prove that unsupervised methods vary significantly in terms of how they learn to represent images, as shown in Figure 2.

State-of-the-art convolution-based unsupervised algorithms, whether they use contrastive learning, clustering, or some pretext task such as colorization, all attempt to learn invariance to some class of augmentations. In other words, they seek to learn a function f , such that $f(I) = f(I_A)$ for some image I and some set of augmentations A that are applied to that image. Xiao *et al.* speculate on the negative effect this may have on learned representations and performance on downstream tasks [56]. However, a careful reading reveals they don't provide evidence for the existence of

transform invariance in unsupervised models, only that their method seems to perform better than MoCo on tasks related to transform invariance. This inspires Section 4.4 – we take a closer look at the presence of augmentation invariance in representations learned by different unsupervised methods.

Prior work is constrained by some combination of limited metrics, use of toy datasets, and a tendency to consider a single unsupervised method as though it is representative of the field. In contrast to this, we contribute the following:

- We utilize multiple methods for measuring properties of learned embeddings, including 3 performance-based benchmarks, and an extension of prior work on uniformity-tolerance analysis to more unsupervised methods across more realistic datasets.
- We perform novel comparison by extending Linear Centered Kernel Alignment (CKA) analysis beyond toy datasets, and by developing two new metrics for comparing embeddings: nearest neighbor graph similarity and linear overlap.
- We propose a framework for measuring augmentation invariance, and demonstrate its results across several methods, augmentations, and datasets.

We conclude in Section 5 with key insights for future unsupervised methods for representation learning:

- Currently, there is no clear “best” method.
- Unsupervised models share properties that are circumstantially undesirable, *e.g.*, color invariance.
- Unsupervised models have similar representations in most layers, but diverge substantially in the last layer.

2. Related Work

2.1. Unsupervised Learning Methods

Some of the first unsupervised methods of the deep learning era were fashioned after pretext tasks from natural language processing. A network would be trained to perform some auxiliary task before transfer for downstream tasks. These auxiliary pretexts task included solving jigsaw puzzles [44], colorization from grayscale [39, 60], inpainting [46], relative patch prediction [17], predicting rotation angle [23], or a combination of tasks [18]. However, the introduction of Noise Contrastive Estimation (NCE) [27] triggered a paradigm shift within unsupervised learning [55], and subsequent methods which utilized contrastive learning [10, 12, 42] would surpass all of these.

Contrastive Learning, which implicitly performs instance discrimination, involves training a model to attract positive pairs (typically augmented views of a given image) and repel negative pairs (augmented views of two different images) [27, 28]. Many papers have proposed successful methods using contrastive learning [4, 11, 12, 29, 31, 32, 42, 45, 48, 55, 58]. In this study, we include Barlow Twins [59], MoCov2 [12], and SimCLR [11].

Clustering has emerged as another important class of unsupervised methods [3, 6–8]. Popular methods such as DeepCluster [6] and SwAV [8], and even methods that don't explicitly attempt representation learning, such as SCAN [49], share many traits with contrastive learning. Among these are the tendency to rely on a large batch size, the use of strong augmentations, which all of these have in common, and implementation details such as the use of projection heads, which are used by SwAV and DeepClusterv2 [6, 8]. We take as our representatives from this category DeepClusterv2 and SwAV.

Other methods, such as SimSiam [13], use neither negative pairs nor clustering objectives. Other methods highlight the potential of vision transformers within unsupervised regimes [9, 40, 57]. However, in this paper, we choose to compare methods that use ResNet-50 backbones. Furthermore, we believe the sample of methods we select are sufficient to support our main points.

2.2. Analysis of Unsupervised Learning

While each paper proposing a new algorithm uses some tasks to attempt to demonstrate their success compared to the prior art, other popular papers have entirely focused on benchmarking, evaluation, and comparison of specific methods. [47] studies augmentation invariance from the perspective of accuracy by using natural images that attempt to vary certain conditions, such as illumination. [54] and [53] address properties of the learned embeddings such as alignment, uniformity, and tolerance. Other works benchmark unsupervised performance on various tasks and con-

ditions [14, 22, 33, 37]. [26] uses the centered kernel alignment (CKA) framework from [36] to compare supervised and unsupervised representations.

This prior work operates under certain constraints. Many papers consider only a single unsupervised method, either MoCo [1, 22, 54] or SimCLR [14, 26], as though it represents the entire field; of those that consider other methods, none consider more than three [33, 37, 47]. We extend uniformity-tolerance analysis [53] to multiple unsupervised methods. We extend CKA analysis [15, 26, 36] beyond tiny datasets, and to multiple methods. We perform FGVC benchmarking [14, 33] for additional metrics and algorithms. We perform CKA analysis [36] to examine augmentation invariance without relying on linear classification as a confounding intermediate step [47]. We develop additional methods for comparing pairs or groups of methods for our first-of-its-kind comprehensive analysis of the similarities and complementary attributes of pretrained unsupervised methods for image representation learning. Our work follows in the spirit of Ericsson et al. in that we perform a comprehensive analysis using a sample of many unsupervised methods [21]. However, by using an entirely different set of tasks and datasets, we are able to both provide further evidence for one of their main conclusions (that no one method is *the best*), and uncover novel insights as well.

3. Methods

3.1. Performance-based Comparison

To measure the quality of the learned representations, we perform three performance-based measurements. For linear evaluation, we use the VISSL repository [24] to train a linear classifier on frozen features. For k -nearest neighbor classification, we also use the settings in VISSL, varying the number of neighbors for fine-grained datasets as described in the appendix. For k -means clustering, we try 10 initializations with the k -means++ method [2], and for ImageNet we use mini-batches of 16,384 images.

To obtain the accuracy for k -means cluster assignments when the number of classes is equal to the number of clusters, we use hungarian matching, mapping clusters to classes one-to-one in such a way that maximizes correspondence to ground truth classes. For overclustering, we greedily map each cluster to the ground truth class which has the most images in the cluster. Accuracy is then computed normally, using the cluster-to-class mapping result as the prediction. For each of linear evaluation, k -NN classification, and k -means clustering, models are trained on training data, and results come from the evaluation on test data.

3.2. Uniformity-Tolerance Tradeoff

To analyze how embeddings are distributed on the hypersphere, we borrow two key properties from prior work:

uniformity [54] and tolerance [53]. Uniformity, U , which describes how closely the embeddings match a uniform distribution on a hypersphere, is defined in Equation 1, where t is a scaling hyperparameter that we set to 2, f represents the model, and x and y are any pair of images.

$$U = \log \mathbb{E}_{x, y \sim p_{\text{data}}} \left[e^{-t \|f(x) - f(y)\|_2^2} \right] \quad (1)$$

Tolerance, T , is given in Equation 2, where f , x , and y are the same as in Equation 1, and I is the indicator function for the ground truth labels, returning 1 when x and y belong to the same class, and 0 otherwise.

$$T = \mathbb{E}_{x, y \sim p_{\text{data}}} \left[(\|f(x)\|_2^T \|f(y)\|_2) \cdot I_{l(x)=l(y)} \right] \quad (2)$$

Whereas uniformity measures how equally spread out the features are, tolerance leverages ground truth labels to indicate how well the embeddings reflect the semantic relationships between the images.

3.3. Linear CKA for Comparing Representations

We follow procedures from prior work to compute Centered Kernel Alignment (CKA) values [26, 36], including using only a linear kernel [26]. To compute this, we first obtain the matrices containing the embeddings for two different methods, such as SimCLR and MoCo, which we represent X and Y . We then compute the Gram matrices of the embedding matrices: $K = XX^T$, $L = YY^T$. The CKA value is given by the normalized Hilbert-Schmidt Independence Criterion (HSIC) [25] as follows:

$$\text{CKA}(K, L) = \frac{\text{HSIC}(K, L)}{\sqrt{\text{HSIC}(K, K)\text{HSIC}(L, L)}} \quad (3)$$

Prior work [26] performs these computations for tiny datasets, consisting of 32×32 images. Since we extend this to 224×224 images, we compensate for the increased memory requirements by taking a set random sample of 10,000 test images when working with ImageNet.

3.4. Proposed Metric: NN Graph Similarity

Even for pure contrastive learning, images that belong to the same ground truth class tend to be semantically similar. In contrast to this tolerance, which relies on ground truth labels and describes the semantic structure in terms of a single model, we propose an unsupervised way to compare the structure of the semantic relationships between two or more learned representations. Specifically, we consider models in terms of their nearest neighbor graphs for a given dataset. Each image is a node, and an image’s top- k neighbors are represented by directed edges. We can thus compute the similarity between two representations by comparing their

nearest neighbor graphs. We choose to perform this computation in terms of neighbor overlap, where neighbor overlap refers to the average number of shared edges (neighbors) per node (image) for the graphs (unsupervised algorithms) considered. This neighbor overlap conveys the similarity in semantic structure learned by different algorithms. A score of 1.0 would indicate the structures are identical – the images have the same nearest k neighbors for both unsupervised algorithms. A score of 0.0 would indicate there are no shared neighbors. We use this method to compare pairs of models, as in Figure 2.

3.5. Proposed Metric: Linear Prediction Overlap

Unlike nearest neighbor graph similarity, this metric takes an indirect approach to compare two or more representations. For each image in a dataset, we get the predictions from the linear classifiers trained on each unsupervised backbone, described in Section 3.1. We then perform a few different calculations. For analysis that relies on ground truth labels, we calculate the portion of the dataset that all models classify correctly, or that no models classify correctly, or that only some subset of models from a group classifies correctly. For analysis that ignores the labels, we compute the percentage of the dataset for which some set of classifiers has the same prediction, regardless of the correctness of the prediction. For this measure, a score of 1.0 would indicate the set of classifiers make the same predictions for all images, while 0.0 would mean they do not have identical predictions for any image. We thus use linear overlap both to compare pairs of methods, as in Figure 2, as well as to compare sets of multiple models, as in Table 5.

3.6. Proposed Analysis: Augmentation Invariance

Popular unsupervised algorithms train models to be augmentation invariant. We develop a method for analyzing the prevalence of augmentation invariance in the representations learned by unsupervised models. Unlike prior work [47], we consider a broader set of unsupervised algorithms, and perform measurements on learned representations directly rather than relying on the performance of learned linear classifiers as a proxy. Instead, we use CKA to compare the similarity between embeddings of augmented and non-augmented images. We take $\text{CKA}(K, L)$ on $K = XX^T$, $L = X_A X_A^T$ where X is the embedding matrix for the images for a given dataset, and X_A is the embedding matrix for the same images after some augmentation A . We thus extend Linear CKA for unsupervised algorithms to overcome the limitations of previous methods and look directly at augmentation invariance for image representations.

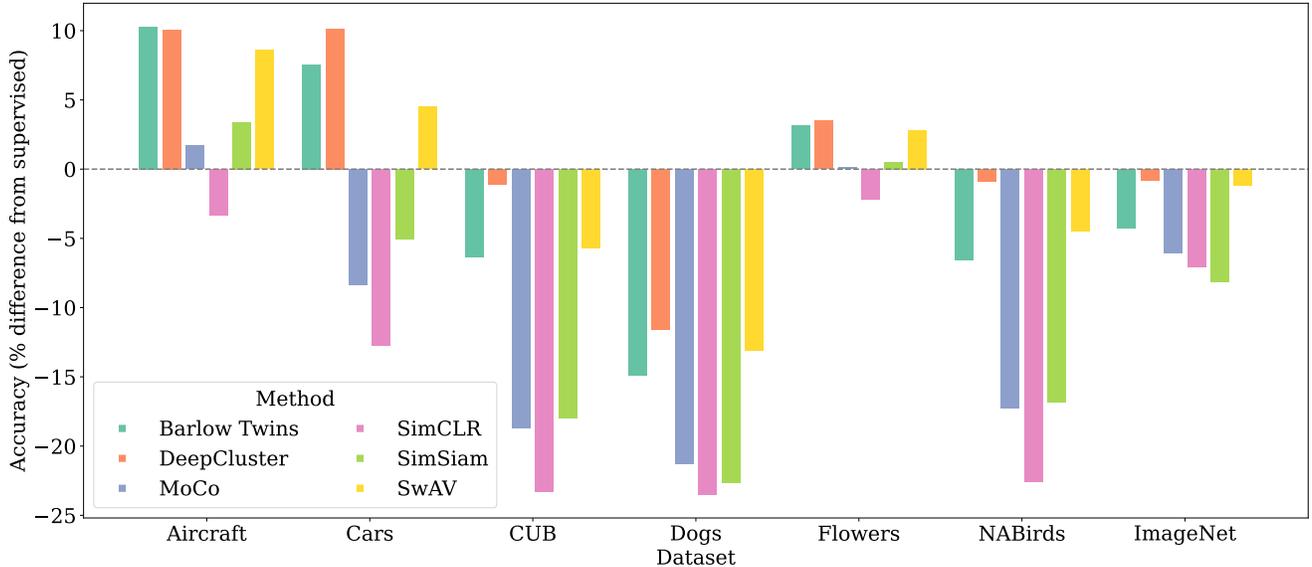


Figure 3. **Linear Classifier Results** on ImageNet and 6 FGVC datasets. Barlow Twins, DeepCluster, and SwAV tend to outperform Moco, SimCLR, and SimSiam, but there is no obvious winner.

4. Analysis

For our representative set of discriminative unsupervised methods, we consider contrastive methods SimCLR [11] and MoCo [12], clustering methods DeepCluster2 [6] and SwAV [8], as well as Barlow Twins [59], which attempts redundancy reduction, and SimSiam [13], which uses neither negative pairs nor clustering. For SimCLR, DeepCluster, and SwAV, we use the 800 epoch checkpoints from the VISSL model zoo [24], unless otherwise specified. For Barlow Twins, we use VISSL’s 1000 epoch checkpoint. For MoCo, we use the 800 epoch checkpoint from the authors. For SimSiam, we use the author-provided 100 epoch checkpoint. While we could have opted to retrain the models, we believe that discrepancy in training time is not a legitimate confounding factor in our analysis, and any attempt to create some “fair” setting would inevitably favor whichever models perform best under that setting. These methods can be tuned on a variety of hyperparameters, and any given setting that is “fair” to one inevitably favors others, so we opt to use the settings for which models are available. Additionally, our ablation experiments (Figure 6) indicate training time would not impact any of our findings.

We perform the experiments here on some subset of the datasets in Table 1, which contain realistic images.

4.1. Performance Benchmarks

As explained in Section 3, we perform VISSL’s linear evaluation, where we train linear classifiers on frozen features from the first convolutional layer and for each of the 4 bottleneck blocks of the ResNet-50 network. We show results for the linear classifier trained on the outputs of the

final block in Figure 3. We show results for k-NN classification in Table 2 and for k-means clustering in Table 3. As an additional benchmark for DeepCluster and SwAV, we also compare their clustering heads, which partition ImageNet into 3000 clusters, to k-means clustering in Table 4. We don’t consider other performance benchmarks such as full finetuning because we are trying to evaluate only the learned embeddings, and not the network initializations.

Table 1. **Datasets** used for experiments in this paper.

Dataset	#Cls	#Train	#Test
FGVC Aircraft [41] (Aircraft)	100	6,667	3,333
Stanford Cars [38] (Cars)	196	8,144	8,041
Caltech Birds [52] (CUB)	200	5,994	5,794
Stanford Dogs [34] (Dogs)	120	12,000	8,580
Oxford Flowers [43] (Flowers)	102	2,040	6,149
NABirds V1 [50] (NABirds)	555	23,929	24,633
ImageNet [16]	1000	1.3mil	50,000

Figure 3 shows that contrary to claims in some prior work, unsupervised methods do not necessarily struggle on FGVC datasets [33]. This is possibly because that work used linear SVMs that were perhaps less accommodating to the ways the unsupervised embeddings tend to be distributed; since our linear evaluation protocol uses batch norm, it is able to better account for this. Nevertheless, we demonstrate while performance is worse for the two birds datasets and the dogs dataset, Barlow Twins, DeepCluster, and SwAV all outperform supervised pre-training for aircraft, cars, and flowers. We suggest that significant over-

Table 2. **k-NN Results.** Again, there is no obvious frontrunner.

Method	Dataset			
	ImageNet	Aircraft	Flowers	NABirds
Supervised	73.41	31.59	77.96	43.25
BTwins	62.90	31.83	86.18	22.29
DCv2	63.70	32.70	84.76	21.05
MoCo	58.59	21.39	74.53	15.40
SimCLR	54.57	21.21	74.78	14.03
SimSiam	53.66	27.39	80.01	15.18
SwAV	61.14	28.77	82.24	15.72

Table 3. **K-Means Results.** Supervised is best for most datasets.

Method	Dataset			
	ImageNet	Aircraft	Flowers	NABirds
Supervised	58.92	15.69	54.97	25.95
BTwins	34.88	13.20	63.70	11.87
DCv2	31.79	13.92	60.20	10.86
MoCo	38.30	9.84	43.34	10.75
SimCLR	29.78	11.16	43.99	9.08
SimSiam	26.20	12.66	54.51	9.53
SwAV	28.69	12.60	56.04	9.26

Table 4. **K-Means Overclustering Results.** The clustering heads of DeepCluster v2 and SwAV outperform k-means on the learned embeddings of DeepCluster v2 and SwAV.

Method	ImageNet		
	$k = 1000$	$k = 3000$	Δ
Supervised K-Means	58.92	65.66	+6.74
DCV2 K-Means	31.79	43.02	+11.23
DCV2 Clustering Head	n/a	54.35	n/a
SwAV K-Means	28.69	37.94	+9.25
SwAV Clustering Head	n/a	48.9	n/a

lap with ImageNet contributes to part of the gap for performance on CUB and Dogs, which likely confounds those results; nevertheless, NABirds results confirm that unsupervised methods have substantial struggles on that dataset.

Table 2 echoes the results in Figure 3. Table 3 and Table 4, however, show that supervised pre-training dominates the k-means metric, except on Flowers. It seems clear that pre-training with labels gives supervised learning a strong advantage for k-NN classification and k-means clustering on ImageNet, to the extent that supervised representations even outperform the clustering heads of DeepCluster and SwAV in the overclustering regime on ImageNet.

We distill 3 key findings from this. First, from each of our benchmarks, unsupervised methods are comparable with supervised for generating embeddings for FGVC, and

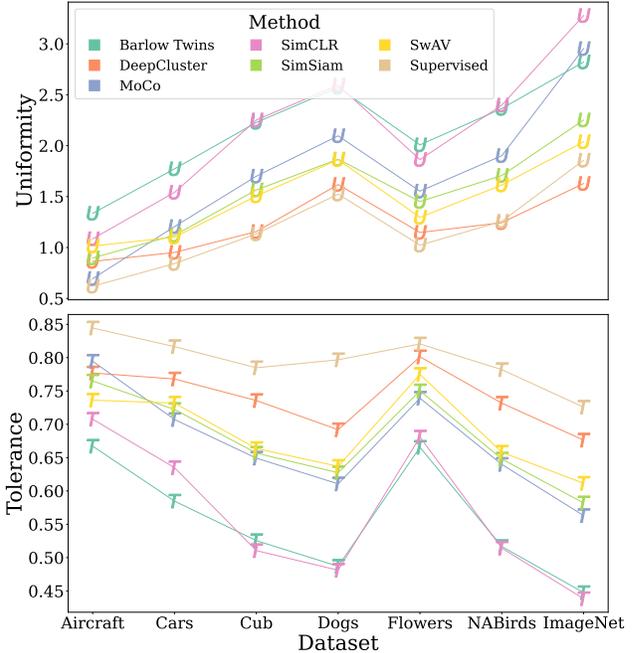


Figure 4. **Uniformity (U) and Tolerance (T)** on ImageNet and 6 FGVC datasets, with datasets sorted by ascending uniformity. Methods with similar objectives (such as the contrastive methods: Barlow Twins, MoCo, and SimCLR) tend to have similar scores.

methods like Barlow Twins, DeepCluster, and SwAV seem particularly competitive. Second, setup matters – architectural decisions such as the design of the classification head can create subtle biases that favor certain methods, such as the SVM analysis from [33] favoring supervised representations. Finally, breadth is helpful; each of our benchmarking methods relies on some assumptions, and indirectly evaluates the robustness of the learned embeddings. Taken together, linear evaluation, k-NN classification, and k-means clustering give a more holistic view of how the representations compare.

4.2. Uniformity-Tolerance Tradeoff

High values for uniformity and tolerance are simultaneously desirable, as they indicate favorable distribution of the embeddings on the hypersphere. Nevertheless, our results in Figure 4 reinforce that in practice, these values have an inverse correlation. This is because as embeddings are more spread out on the hypersphere in general, they tend to also be more spread out with respect to each ground truth class. It is perhaps for this reason that supervised pre-training has, in general, the most tolerant and least uniform embeddings.

DeepCluster is similar to supervised for both these metrics. This is unsurprising when considering the DeepCluster objective: cross-entropy loss for pseudolabels. SwAV, with its own clustering objective, exhibits some of these same

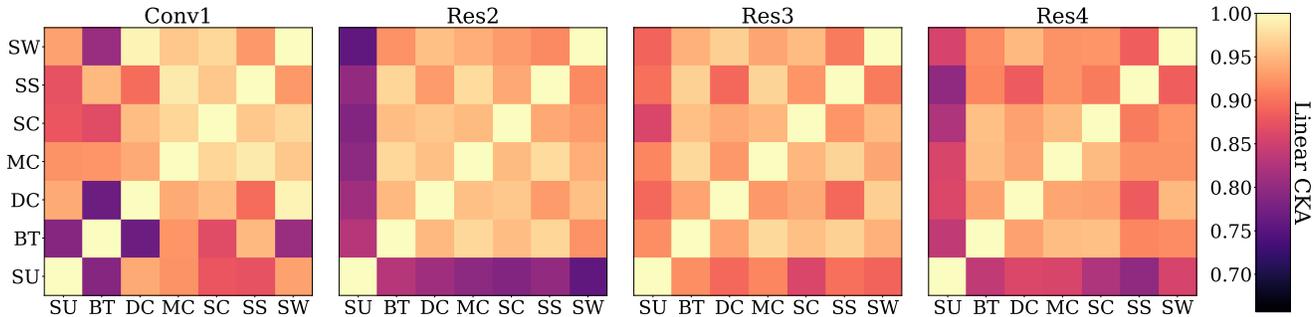


Figure 5. **Linear CKA for initial convolutional layer and first 3 bottleneck blocks** for validation set of Imagenet. Linear CKA for the last block can be found in Figure 2. In contrast to the final layer, representations are fairly similar in the initial and intermediate blocks.

tendencies, whereas the contrastive methods are the opposite. This strengthens support for our hypothesis – unsupervised methods are quite diverse in terms of the distribution of learned embeddings on the hypersphere.

4.3. Measuring Similarity of Representations

We consider three main metrics for measuring the similarity between representations: linear CKA, which compares embeddings for any given pair of models, nearest neighbor graph similarity, which compares nearest neighbors for a set of models, and overlap in linear predictions, which compares the predictions of linear classifiers trained on a set of models. See Figure 2 for linear CKA, neighbor similarity, and overlap in linear predictions for the outputs of each ResNet50’s final block for Imagenet. For linear CKA for the other layers we measure, see Figure 5. In Table 5, we perform an analysis of linear overlap for groups of models, and leverage ground truth labels to evaluate the uniqueness of the different linear classifiers.

Results in Figure 2 indicate that similar representations tend to have similar neighbors, and classifiers trained on more similar representations tend to make more similar predictions. The metrics reveal similarities between algorithms with related objectives, such as DeepCluster and SwAV. More surprisingly, we find more similarity between supervised and unsupervised representations such as DeepCluster than we do between some unsupervised representations, such as MoCo and Barlow Twins. Figure 5 simultaneously confirms these findings and those of [26], who found that supervised and unsupervised representations diverged the most in the final layer. We thus extend their hypothesis from SimCLR to additional unsupervised methods, and provide evidence that unsupervised algorithms differ enough from each other to have very different final representations. This highlights the idea that it is unreasonable to make the “supervised” vs. “unsupervised” comparisons that are so common in the literature where the “unsupervised” is represented by only a couple of algorithms.

We also compute CKA for unsupervised models under

Table 5. **Results for linear overlap.** We examine overlap in predictions for linear classifiers trained on frozen features. On the top, we report how many images were predicted correctly by classifiers for both supervised and any unsupervised method, by either, and by neither. On the bottom, we compare within unsupervised, considering how many images were uniquely classified correctly by a single linear classifier. We find substantial uniqueness for each algorithm, which attests to their complementary nature.

Method	Dataset		
	ImageNet	Aircraft	NABirds
Sup. and Unsup.	73.64	81.40	55.05
Sup. Only	2.40	0.03	6.05
Unsup. Only	10.34	18.45	17.78
Neither	13.62	0.12	21.12
All Unsup.	58.19	80.08	30.49
BTwins Only	0.97	0.24	2.46
DCv2 Only	1.74	0.18	4.16
MoCo Only	0.69	0.09	0.87
SimSiam Only	0.64	0.00	0.86
SwAV Only	1.74	0.21	2.80
No Unsup.	16.02	0.15	27.17

non-default settings, to validate our other findings. Figure 6, when compared to Figure 2, shows that training time, at least in the case of SimCLR, induces a comparatively small difference in learned representations, despite the large gaps in performance on benchmarks such as linear evaluation. While linear evaluation accuracies are different by several percentage points, SimCLR models trained 800 epochs are more similar to the other SimCLR checkpoints than to any other unsupervised algorithm. Our finding on training time stands in contrast to the results from other settings. We find details such as cropping strategy and batch size can have a massive impact on the similarity of neural representations, to the extent that SwAV with its full batch size and cropping strategy is more related to DeepCluster with the same settings than to SwAV methods that utilize small batch sizes. We suggest that future work should probe these effects fur-

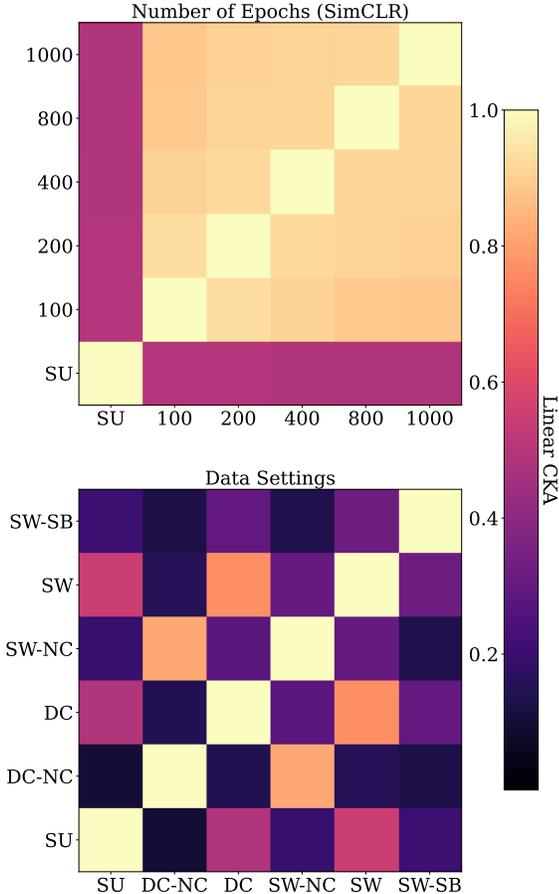


Figure 6. **Linear CKA for ablation settings** for final residual outputs on the validation set of ImageNet. On the top, we compare supervised (SU) with SimCLR trained for 100, 200, 400, 800, and 1000 epochs. On the bottom, we compare SU with DeepCluster (DC) and SwAV (SW) trained for 400 epochs, as well as with SwAV and DeepCluster trained for 400 epochs with no crops (NC) and SwAV for 400 epochs with a smaller batch size (256).

ther, and examine the extent to which the number of crops, batch size, and data augmentations can affect learned representations and downstream applications for various unsupervised algorithms.

4.4. Augmentation Invariance

We use linear CKA to test for augmentation invariance with respect to color jitter, blurring, jitter with blurring, horizontal and vertical flipping, and rotation. Figure 7 provides evidence that, contrary to the conclusions of [47], and confirming most other prior work, unsupervised algorithms learn representations that are invariant to their training augmentations. We note that the invariance is at least somewhat weaker for the clustering algorithms, SwAV and DeepCluster. Also, the unsupervised methods tend to be somewhat more invariant to augmentations not used at training time,

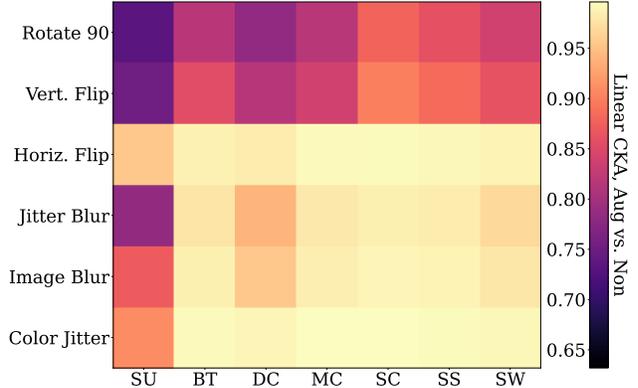


Figure 7. **Augmentation Invariance**, using Linear CKA for 7 algorithms on ImageNet, for 1 augmentation used in all pretrainings (horizontal flip), 3 used in unsupervised pretraining (color jitter, blurring, and both simultaneously), as well as 2 not used. Unsupervised models exhibit invariance to their training augmentations.

rotation, and vertical flips. Nevertheless, these experiments suggest that for applications where color is critical, such as bird classification, methods that rely on learned color invariance are destined to underperform. Future works which seek to mitigate color invariance can leverage our framework as a way to evaluate success.

5. Key Takeaways and Conclusion

We highlight the following as key takeaways from our findings. First, there is no clear “best” method. Therefore, it is essential to avoid over-indexing on a given metric, and in the context of applications, representations should be selected to optimize for both downstream data and task. Second, unsupervised methods share properties that are situationally undesirable, such as robust color invariance. Thus, it is important for future work to develop methods that mitigate certain invariances when necessary, and our CKA-based framework can be utilized to validate their success. Finally, representations for all algorithms we considered are fairly similar until the last layer, where specialized loss functions or even training settings such as augmentation strategy and batch size can induce learning of substantially different representations. Therefore, it is critical to not assume that one method, such as MoCo or SimCLR, can act as a representative of the field. Additionally, taken in the context of our other findings, we suggest researchers continue to pursue meta-learning, distillation, ensembles, and other approaches that effectively combine different unsupervised algorithms to leverage their complementary nature.

Acknowledgements. This project was partially funded by the DARPA SAIL-ON (W911NF2020009) program, an independent grant from Facebook AI, and Amazon Research Award to AS.

References

- [1] Srikar Appalaraju, Yi Zhu, Yusheng Xie, and István Fehérvári. Towards good practices in self-supervised representation learning. *CoRR*, abs/2012.00868, 2020. 2, 3
- [2] David Arthur and Sergei Vassilvitskii. k-means++: The advantages of careful seeding. Technical report, Stanford, 2006. 3
- [3] Yuki Markus Asano, Christian Rupprecht, and Andrea Vedaldi. Self-labelling via simultaneous clustering and representation learning. *arXiv preprint arXiv:1911.05371*, 2019. 3
- [4] Philip Bachman, R Devon Hjelm, and William Buchwalter. Learning representations by maximizing mutual information across views. *arXiv preprint arXiv:1906.00910*, 2019. 3
- [5] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural image synthesis. *arXiv preprint arXiv:1809.11096*, 2018. 1
- [6] Mathilde Caron, Piotr Bojanowski, Armand Joulin, and Matthijs Douze. Deep clustering for unsupervised learning of visual features. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 132–149, 2018. 3, 5
- [7] Mathilde Caron, Piotr Bojanowski, Julien Mairal, and Armand Joulin. Unsupervised pre-training of image features on non-curated data. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2959–2968, 2019. 3
- [8] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. *arXiv preprint arXiv:2006.09882*, 2020. 3, 5
- [9] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. *arXiv preprint arXiv:2104.14294*, 2021. 3
- [10] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 1597–1607. PMLR, 13–18 Jul 2020. 3
- [11] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020. 3, 5
- [12] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020. 3, 5
- [13] Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15750–15758, 2021. 3, 5
- [14] Elijah Cole, Xuan Yang, Kimberly Wilber, Oisín Mac Aodha, and Serge J. Belongie. When does contrastive visual representation learning work? *CoRR*, abs/2105.05837, 2021. 2, 3
- [15] Corinna Cortes, Mehryar Mohri, and Afshin Rostamizadeh. Algorithms for learning kernels based on centered alignment. *Journal of Machine Learning Research*, 13(28):795–828, 2012. 3
- [16] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009. 5
- [17] Carl Doersch, Abhinav Gupta, and Alexei A Efros. Unsupervised visual representation learning by context prediction. In *Proceedings of the IEEE international conference on computer vision*, pages 1422–1430, 2015. 3
- [18] Carl Doersch and Andrew Zisserman. Multi-task self-supervised visual learning. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2051–2060, 2017. 3
- [19] Jeff Donahue, Philipp Krähenbühl, and Trevor Darrell. Adversarial feature learning. *arXiv preprint arXiv:1605.09782*, 2016. 1
- [20] Jeff Donahue and Karen Simonyan. Large scale adversarial representation learning. *arXiv preprint arXiv:1907.02544*, 2019. 1
- [21] Linus Ericsson, Henry Gouk, and Timothy M. Hospedales. How well do self-supervised models transfer? *CoRR*, abs/2011.13377, 2020. 3
- [22] Wouter Van Gansbeke, Simon Vandenhende, Stamatios Georgoulis, and Luc Van Gool. Revisiting contrastive methods for unsupervised learning of visual representations. *CoRR*, abs/2106.05967, 2021. 2, 3
- [23] Spyros Gidaris, Praveer Singh, and Nikos Komodakis. Unsupervised representation learning by predicting image rotations. *arXiv preprint arXiv:1803.07728*, 2018. 3
- [24] Priya Goyal, Quentin Duval, Jeremy Reizenstein, Matthew Leavitt, Min Xu, Benjamin Lefauieux, Mannat Singh, Vinicius Reis, Mathilde Caron, Piotr Bojanowski, Armand Joulin, and Ishan Misra. *Vissl*. <https://github.com/facebookresearch/vissl>, 2021. 3, 5
- [25] Arthur Gretton, Kenji Fukumizu, Choon Hui Teo, Le Song, Bernhard Schölkopf, Alexander J Smola, et al. A kernel statistical test of independence. In *Nips*, volume 20, pages 585–592. Citeseer, 2007. 4
- [26] Tom George Grigg, Dan Busbridge, Jason Ramapuram, and Russ Webb. Do self-supervised and supervised methods learn similar visual representations?, 2021. 2, 3, 4, 7
- [27] Michael Gutmann and Aapo Hyvärinen. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 297–304. JMLR Workshop and Conference Proceedings, 2010. 3
- [28] Raia Hadsell, Sumit Chopra, and Yann LeCun. Dimensionality reduction by learning an invariant mapping. In *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, volume 2, pages 1735–1742. IEEE, 2006. 3
- [29] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross B. Girshick. Momentum contrast for unsupervised vi-

- sual representation learning. *CoRR*, abs/1911.05722, 2019. [3](#)
- [30] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition, 2015. [1](#)
- [31] Olivier Henaff. Data-efficient image recognition with contrastive predictive coding. In *International Conference on Machine Learning*, pages 4182–4192. PMLR, 2020. [3](#)
- [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. *arXiv preprint arXiv:1808.06670*, 2018. [3](#)
- [33] Grant Van Horn, Elijah Cole, Sara Beery, Kimberly Wilber, Serge J. Belongie, and Oisín Mac Aodha. Benchmarking representation learning for natural world image collections. *CoRR*, abs/2103.16483, 2021. [3](#), [5](#), [6](#)
- [34] Aditya Khosla, Nityananda Jayadevaprakash, Bangpeng Yao, and Li Fei-Fei. Novel dataset for fine-grained image categorization. In *First Workshop on Fine-Grained Visual Categorization, IEEE Conference on Computer Vision and Pattern Recognition*, Colorado Springs, CO, June 2011. [5](#)
- [35] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. [1](#)
- [36] Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey Hinton. Similarity of neural network representations revisited. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 3519–3529. PMLR, 09–15 Jun 2019. [3](#), [4](#)
- [37] Klemen Kotar, Gabriel Ilharco, Ludwig Schmidt, Kiana Ehsani, and Roozbeh Mottaghi. Contrasting contrastive self-supervised representation learning pipelines. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 9949–9959, October 2021. [3](#)
- [38] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *4th International IEEE Workshop on 3D Representation and Recognition (3dRR-13)*, Sydney, Australia, 2013. [5](#)
- [39] Gustav Larsson, Michael Maire, and Gregory Shakhnarovich. Learning representations for automatic colorization. In *European conference on computer vision*, pages 577–593. Springer, 2016. [3](#)
- [40] Chunyuan Li, Jianwei Yang, Pengchuan Zhang, Mei Gao, Bin Xiao, Xiyang Dai, Lu Yuan, and Jianfeng Gao. Efficient self-supervised vision transformers for representation learning. *arXiv preprint arXiv:2106.09785*, 2021. [3](#)
- [41] S. Maji, J. Kannala, E. Rahtu, M. Blaschko, and A. Vedaldi. Fine-grained visual classification of aircraft. Technical report, 2013. [5](#)
- [42] Ishan Misra and Laurens van der Maaten. Self-supervised learning of pretext-invariant representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6707–6717, 2020. [3](#)
- [43] Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number of classes. In *Indian Conference on Computer Vision, Graphics and Image Processing*, Dec 2008. [5](#)
- [44] Mehdi Noroozi and Paolo Favaro. Unsupervised learning of visual representations by solving jigsaw puzzles. In *European conference on computer vision*, pages 69–84. Springer, 2016. [3](#)
- [45] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. [3](#)
- [46] Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor Darrell, and Alexei A Efros. Context encoders: Feature learning by inpainting. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2536–2544, 2016. [3](#)
- [47] Senthil Purushwalkam and Abhinav Gupta. Demystifying contrastive self-supervised learning: Invariances, augmentations and dataset biases. *CoRR*, abs/2007.13916, 2020. [3](#), [4](#), [8](#)
- [48] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive multiview coding. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XI 16*, pages 776–794. Springer, 2020. [3](#)
- [49] Wouter Van Gansbeke, Simon Vandenhende, Stamatios Georgoulis, Marc Proesmans, and Luc Van Gool. Scan: Learning to classify images without labels. In *European Conference on Computer Vision*, pages 268–285. Springer, 2020. [3](#)
- [50] Grant Van Horn, Steve Branson, Ryan Farrell, Scott Haber, Jessie Barry, Panos Ipeirotis, Pietro Perona, and Serge Belongie. Building a bird recognition app and large scale dataset with citizen scientists: The fine print in fine-grained dataset collection. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 595–604, 2015. [5](#)
- [51] Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol. Extracting and composing robust features with denoising autoencoders. In *Proceedings of the 25th international conference on Machine learning*, pages 1096–1103, 2008. [1](#)
- [52] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. The Caltech-UCSD Birds-200-2011 Dataset. Technical Report CNS-TR-2011-001, California Institute of Technology, 2011. [5](#)
- [53] Feng Wang and Huaping Liu. Understanding the behaviour of contrastive loss. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2495–2504, June 2021. [3](#), [4](#)
- [54] Tongzhou Wang and Phillip Isola. Understanding contrastive representation learning through alignment and uniformity on the hypersphere. *CoRR*, abs/2005.10242, 2020. [2](#), [3](#), [4](#)
- [55] Zhirong Wu, Yuanjun Xiong, Stella X. Yu, and Dahua Lin. Unsupervised feature learning via non-parametric instance discrimination. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. [3](#)
- [56] Tete Xiao, Xiaolong Wang, Alexei A Efros, and Trevor Darrell. What should not be contrastive in contrastive learning. In *International Conference on Learning Representations*, 2021. [2](#)

- [57] Zhenda Xie, Yutong Lin, Zhuliang Yao, Zheng Zhang, Qi Dai, Yue Cao, and Han Hu. Self-supervised learning with swin transformers. *arXiv preprint arXiv:2105.04553*, 2021. [3](#)
- [58] Mang Ye, Xu Zhang, Pong C. Yuen, and Shih-Fu Chang. Un-supervised embedding learning via invariant and spreading instance feature. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. [3](#)
- [59] Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and Stéphane Deny. Barlow twins: Self-supervised learning via redundancy reduction. *CoRR*, abs/2103.03230, 2021. [3](#), [5](#)
- [60] Richard Zhang, Phillip Isola, and Alexei A Efros. Colorful image colorization. In *European conference on computer vision*, pages 649–666. Springer, 2016. [3](#)

Beyond Supervised vs. Unsupervised: Representative Benchmarking and Analysis of Image Representation Learning

Supplementary Material

A. Code and Assets

To reproduce our results, please visit our repository at <https://github.com/mgwillia/unsupervised-analysis>. Where specified (see our repository for details), we use code from VISSL (<https://github.com/facebookresearch/vissl/>) and SCAN (<https://github.com/wvangansbeke/Unsupervised-Classification>). These repositories have an MIT License and Creative Commons License, respectively.

B. k-NN Details

For k-NN classification, we use the VISSL defaults for ImageNet: 200 neighbors. For the FGVC datasets, there are too few images per class to use this approach. Instead, we try values in the set $\{0, 5, 10, 15, 20, 25, 30, 35, 40, 45, 50\}$ and choose whichever value maximizes accuracy.

C. More Benchmark Results

Here, we give an expanded look at our benchmarks. Table 6 complements Figure 3 by providing the same data, in tabular form. Figures 8 and 9 along with Tables 7 and 8 do the same for k-NN and k-means, offering an expansion of the results shown in Tables 2 and 3.

We verify claims we make about the SimCLR models from the main paper. Specifically, we say that training time has a significant impact on results, while not changing the representations substantially (see Figure 6). Table 9 offers evidence supporting our claim.

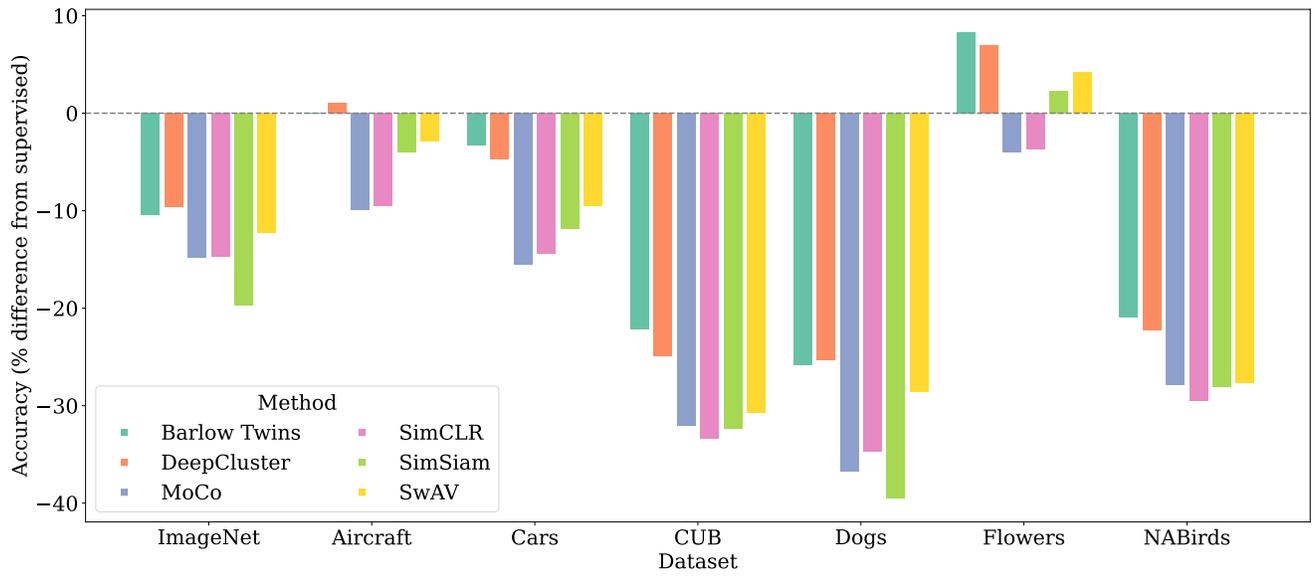


Figure 8. **k-NN** results.

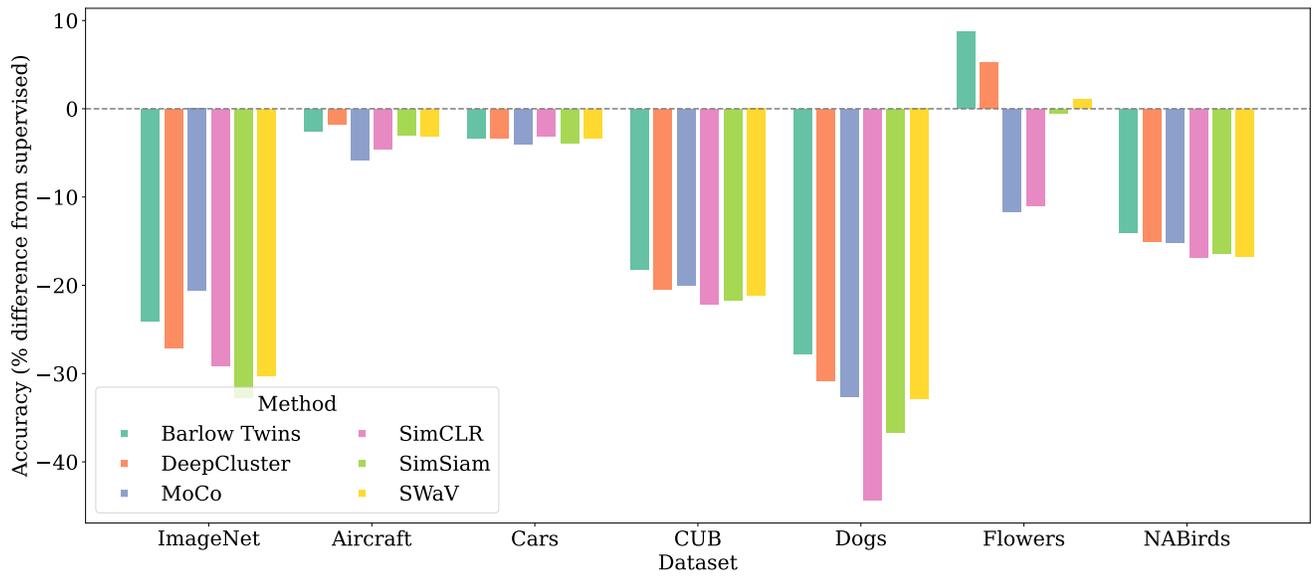


Figure 9. **K-Means** results.

Table 6. **Linear evaluation results.**

Method	Dataset						
	ImageNet	Aircraft	Cars	CUB	Dogs	Flowers	NABirds
Supervised	76.04	48.05	57.72	70.57	88.92	91.30	61.10
Barlow Twins	71.78	58.51	65.30	63.98	74.35	94.21	54.59
DeepCluster	75.19	58.36	67.76	69.82	77.80	94.46	59.89
MoCo	69.95	49.88	49.68	51.95	67.93	91.14	43.51
SimCLR	68.95	44.39	45.00	47.36	65.84	88.90	38.43
SimSiam	67.89	52.22	53.34	52.55	66.50	91.86	44.18
SwAV	74.87	55.73	61.95	65.10	75.99	93.97	56.52

Table 7. **k-NN results.**

Method	Dataset						
	ImageNet	Aircraft	Cars	CUB	Dogs	Flowers	NABirds
Supervised	73.41	31.59	30.16	56.63	88.38	77.96	43.25
Barlow Twins	62.90	31.83	26.94	34.41	62.53	86.18	22.29
DeepCluster	63.70	32.70	25.48	31.74	62.97	84.76	21.05
MoCo	58.59	21.39	14.64	24.35	51.60	74.53	15.40
SimCLR	54.57	21.21	14.74	23.21	49.63	74.78	14.03
SimSiam	53.66	27.39	18.41	24.20	48.97	80.01	15.18
SwAV	61.14	28.77	20.84	25.75	59.87	82.24	15.72

Table 8. **K-Means results.**

Method	Dataset						
	ImageNet	Aircraft	Cars	CUB	Dogs	Flowers	NABirds
Supervised	58.92	15.69	11.95	35.23	53.69	54.97	25.95
DeepCluster	31.79	13.92	8.66	14.81	22.84	60.20	10.86
MoCo	38.30	9.84	7.98	15.21	21.10	43.34	10.75
Barlow Twins	34.88	13.20	8.63	17.07	25.94	63.70	11.87
SimCLR	29.78	11.16	8.80	13.07	9.41	43.99	9.08
SimSiam	26.20	12.66	8.03	13.57	17.07	54.51	9.53
SwAV	28.69	12.60	8.66	14.05	20.79	56.04	9.26

Table 9. **Linear Evaluation for SimCLR with varying training time.**

Method	Dataset						
	ImageNet	Aircraft	Cars	CUB	Dogs	Flowers	NABirds
100 Epochs	64.76	44.81	44.67	43.17	60.44	88.72	34.34
200 Epochs	66.92	45.56	46.31	46.05	62.48	89.39	36.90
400 Epochs	67.93	44.84	46.36	46.00	64.35	89.08	37.23
800 Epochs	68.95	44.39	45.00	47.36	65.84	88.90	38.43
1000 Epochs	64.57	45.26	44.55	46.93	66.25	88.57	37.93