

---

# *Few-Class Arena: A Benchmark for Efficient Selection of Vision Models and Dataset Difficulty Measurement*

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 A wide variety of benchmark datasets with many classes (80-1000) have been  
2 created to assist Computer Vision architectural evolution. An increasing number of  
3 vision models are evaluated with these many-class datasets. However, real-world  
4 applications often involve substantially fewer classes of interest (2-10). This gap  
5 between many and few classes makes it difficult to predict performance of the  
6 few-class applications using models trained on the available many-class datasets.  
7 To date, little has been offered to evaluate models in this *Few-Class Regime*. We  
8 propose *Few-Class Arena (FCA)*, as a unified benchmark with focus on testing  
9 efficient image classification models for few classes. We conduct a systematic  
10 evaluation of the ResNet family trained on ImageNet subsets from 2 to 1000 classes,  
11 and test a wide spectrum of Convolutional Neural Networks and Transformer  
12 architectures over ten datasets by using our newly proposed *FCA* tool. Furthermore,  
13 to aid an up-front assessment of dataset difficulty and a more efficient selection  
14 of models, we incorporate a difficulty measure as a function of class similarity.  
15 *FCA* offers a new tool for efficient machine learning in the *Few-Class Regime*,  
16 with goals ranging from a new efficient class similarity proposal, to lightweight  
17 model architecture design, to a new scaling law. *FCA* is user-friendly and can be  
18 easily extended to new models and datasets, facilitating future research work. Our  
19 benchmark is available at <https://github.com/fewclassarena/fca>.

## 20 1 Introduction

21 The de-facto benchmarks for evaluating efficient vision models are large scale with many classes  
22 (e.g. 1000 in ImageNet [1], 80 in COCO [2], etc.). Such benchmarks have expedited the advance of  
23 vision neural networks toward efficiency [3, 4, 5, 6, 7, 8, 9, 10] with the hope of reducing the financial  
24 and environmental cost of vision models [11, 12]. More efficient computation is facilitated by using  
25 quantization [13, 14, 15], pruning [16, 17, 18, 19], and data saliency [20]. Despite efficiency  
26 improvements such as these, many-class datasets are still the standard of model evaluation.

27 Real-world applications, however, typically comprise only a few number of classes (e.g, less than  
28 10) [21, 22, 23] which we termed *Few-Class Regime*. To deploy a vision model pre-trained on large  
29 datasets in a specific environment, it requires the re-evaluation of published models or even retraining  
30 to find an optimal model in an expensive architectural search space [24].

31 One major finding is that, apart from scaling down model and architectural design for efficiency,  
32 dataset difficulty also plays a vital role in model selection [25] (described in Section 4.3).

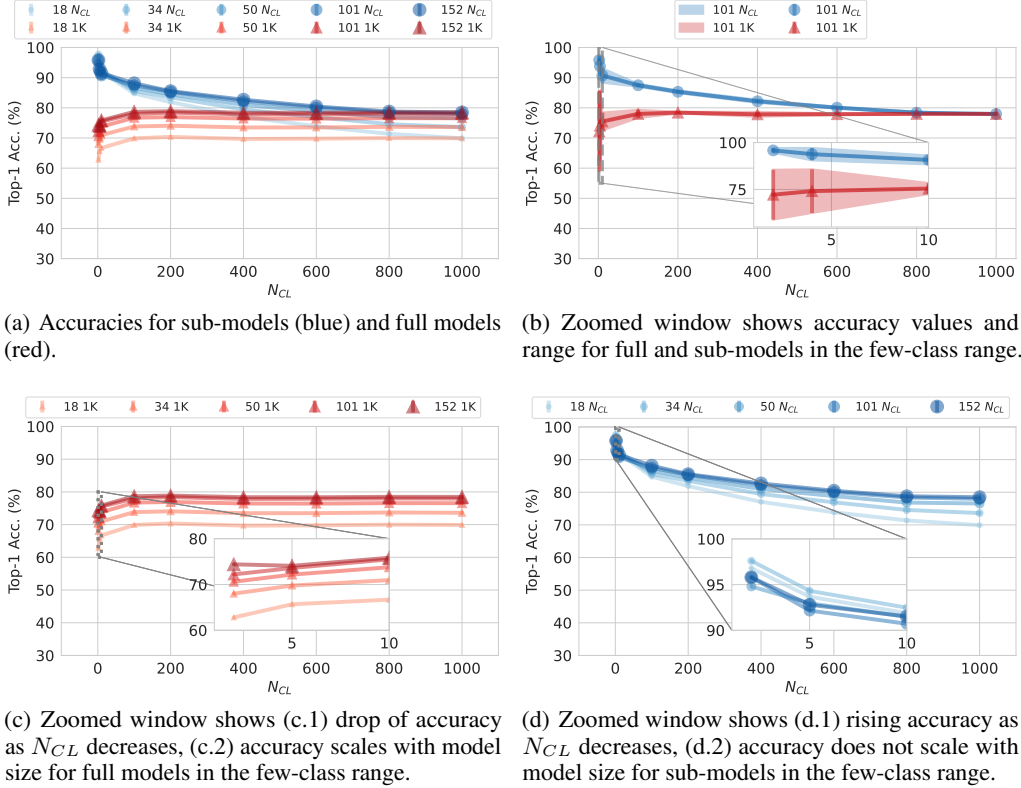


Figure 1: Top-1 accuracies of various scales of ResNet, whose model sizes are shown in the legend, and whose plots vary from dark to light by decreasing size. Plots range along number of classes  $N_{CL}$  from the full ImageNet size (1000) down to the *Few-Class Regime*. Each model is tested on 5 subsets whose  $N_{CL}$  classes are randomly sampled from the original 1000 classes. (a) Plots for sub-models trained on subsets of classes (blue) and full models trained on all 1000 classes (red). (b) Zoomed window shows the standard deviation of subset’s accuracies is much smaller than for the full model. (c.1) Full model accuracies drop when  $N_{CL}$  decreases. (c.2) Full model accuracies increase as model scales up in the *Few-Class Regime*. (d.1) Sub-model accuracies grow as  $N_{CL}$  decreases. (d.2) Sub-model accuracies do not increase when model scales up in the *Few-Class Regime*.

33 Figure 1 summarizes several key findings under the *Few-Class Regime*. On the left graph in red  
 34 are accuracy results for a range of number of classes  $N_{CL}$  for what we call the “full model”, that  
 35 is ResNet models pre-trained on the full 1000 classes of ImageNet (generally available from many  
 36 websites). On the right are accuracy results for what we call “sub-models”, each of which is trained  
 37 and tested on the same  $N_{CL}$ , where this number of classes is sampled from the full dataset down to  
 38 the *Few-Class Regime*. Findings include the following. (a) Sub-models attain higher upper-bound  
 39 accuracy than full models. (b) The range of accuracy widens for full models at few-classes, which  
 40 increases the uncertainty of a practitioner selecting a model for few classes. In contrast, sub-models  
 41 narrow the range. (c) Full models follow the scaling law [26] in the dimension of model size - larger  
 42 models (darker red) have higher accuracy from many to few classes. (4) Surprisingly, the scaling law  
 43 is violated for sub-models in the *Few-Class Regime* (see the zoomed-in subplot) where larger models  
 44 (darker blue) do not necessarily perform better than smaller ones (lighter blue). From these plots,  
 45 our key insight is that, instead of using full models, researchers and practitioners in the *Few-Class*  
 46 *Regime* should use sub-models for selection of more efficient models.

47 However, obtaining sub-models involves computationally expensive training and testing cycles since  
 48 they need to be converged on each of the few-class subsets. By carefully studying and comparing the  
 49 experiment and evaluation setup of these works in the literature, we observe that, how models scale  
 50 down to *Few-Class Regime* is rarely studied. The lack of comprehensive benchmarks for *few-class*  
 51 research impedes both researchers and practitioners from quickly finding models that are the most

52 efficient for their dataset size. To fill this need, we propose a new benchmark, *Few-Class Arena (FCA)*,  
53 with the goal of benchmarking vision models under few-class scenarios. To our best knowledge, *FCA*  
54 is the first benchmark for such a purpose.

55 We formally define *Few-Class Regime* as a scenario where the dataset has a limited number of classes.  
56 Real-world applications often comprise only a few number of classes (e.g.  $N_{CL} < 10$  or 10% classes  
57 of a dataset). Consequently, *Few-Class Arena* refers to a benchmark to conduct research experiments  
58 to compare models in the *Few-Class Regime*. This paper focuses on the image classification task,  
59 although *Few-Class Regime* can generalize to object detection and other visual tasks.

60 **Statement of Contributions.** Four contributions are listed below:

- 61 • To be best of our knowledge, we are the first to explore the problems in the *Few-Class*  
62 *Regime* and develop a benchmark tool *Few-Class Arena (FCA)* to facilitate scientific research,  
63 analysis, and discovery for this range of classes.
- 64 • We introduce a scalable few-class data loading approach to automatically load images and  
65 labels in the *Few-Class Regime* from the full dataset, avoiding the need to duplicate data  
66 points for every additional few-class subset.
- 67 • We incorporate dataset similarity as an inverse difficulty measurement in *Few-Class Arena*  
68 and propose a novel Silhouette-based similarity score named *SimSS*. By leveraging the visual  
69 feature extraction power of CLIP and DINOv2, we show that *SimSS* is highly correlated  
70 with ResNet performance in the *Few-Class Regime* with Pearson coefficient scores  $\geq 0.88$ .
- 71 • We conduct extensive experiments that comprise ten models on ten datasets and 2-1000  
72 numbers of classes on ImageNet, totalling 1591 training and testing runs. In-depth analyses  
73 on this large body of testing reveal new insights in the *Few-Class Regime*.

## 74 2 Related Work

75 **Visual Datasets and Benchmarks.** To advance deep neural network research, a wealth of large-scale  
76 many-class datasets has been developed for benchmarking visual neural networks over a variety of  
77 tasks. Typical examples <sup>1</sup> include 1000 classes in ImageNet [1] for image classification, and 80 object  
78 categories in COCO [2] for object detection. Previous benchmarks also extend vision to multimodal  
79 research such as image-text [27, 28, 29, 30]. While prior works often scale up the number of object  
80 categories for general purpose comparison, studies [31, 32] raise a concern on whether models trained  
81 on datasets with such a large number of classes (e.g. ImageNet) can be reliably transferred to real  
82 world applications often with far fewer classes. A close work to ours is vision backbone comparison  
83 [33] whose focus is on model architectures. Our perspective differs in a focus on cases with fewer  
84 number of classes, which often better aligns with real-world scenarios.

85 **Dataset Difficulty Measurement.** Research has shown the existence of inherent dataset difficulty  
86 [32] for classification and other analytic tasks. Efficient measurement methods are proposed to  
87 characterize dataset difficulty using Silhouette Score [34], K-means Fréchet inception distance  
88 [35, 36, 37], and Probe nets [25]. Prior studies have proposed image quality metrics using statistical  
89 heuristics, including peak signal-to-noise ratio (PSNR) [38], structural similarity (SSIM) Index  
90 [39], and visual information fidelity VIF [40]. A neuroscience-based image difficulty metric [32]  
91 is defined as the minimum viewing time related to object solution time (OST) [41]. Another type  
92 of difficulty measure method consists of additional procedures such as c-score [42], prediction  
93 depth [43], and adversarial robustness [44]. Our work aligns with the line of research [45, 46, 47]  
94 involving similarity-based difficulty measurements: similar images are harder to distinguish from  
95 each other while dissimilar images are easier. Previous studies are mainly in the image retrieval  
96 context [48, 49, 50]. Similarity score is used in [51] with the limitation that a model serving similarity  
97 measurement has to be trained for one dataset. We push beyond this limit by leveraging large vision  
98 models that learn general visual features using CLIP [52] and DINOv2 [53]. The study [32] shows  
99 that CLIP generalizes well to both easy and hard images, making it a good candidate for measuring

---

<sup>1</sup>A detailed list of many-class datasets used in this paper can be found in the Appendix.

100 image difficulty. Supported by the evidence that better classifiers can act as better perceptual feature  
101 extractors [54], in later sections we show how CLIP and DINOv2 will be used as our similarity base  
102 function.

103 Despite the innovation of difficulty measure algorithms on many-class datasets, little attention has  
104 been paid to leveraging these methods in the *Few-Class Regime*. We show that, as the number of  
105 classes decreases, sub-dataset difficulty in the *Few-Class Regime* plays a more critical role in efficient  
106 model selection. To summarize, unlike previous work on many-class benchmarks and difficulty  
107 measurements, our work takes few-class and similarity-based dataset difficulty into consideration,  
108 and in doing so we believe the work pioneers the development of visual benchmark dedicated to  
109 research in the *Few-Class Regime*.

### 110 3 Few-Class Arena (FCA)

111 We introduce the *Few-Class Arena (FCA)* benchmark in this section. In practice, we have integrated  
112 *FCA* into the MMPreTrain framework [55], implemented in Python3 and Pytorch<sup>2</sup>.

#### 113 3.1 Goals

114 **1. Generality.** All vision models and existing datasets for classification should be compatible in this  
115 framework. In addition, users can extend to custom models and datasets for their needs.

116 **2. Efficiency.** The benchmark should be time- and space-efficient for users. The experimental setup  
117 for the few-class benchmark should be easily specified by a few hyper-parameters (e.g. number of  
118 classes). Since the few-class regime usually includes sub-datasets extracted from the full dataset, the  
119 benchmark should be able to locate those sub-datasets without generating redundant duplicates for  
120 reasons of storage efficiency. For time-efficiency, it should conduct training and testing automatically  
121 through use of user-specified configuration files, without users' manual execution.

122 **3. Large-Scale Benchmark.** The tool should allow for large-scale benchmarking, including training  
123 and testing of different vision models on various datasets when the number of classes varies.

#### 124 3.2 Few-Class Dataset Preparation

125 *Few-Class Arena* provides an easy way to prepare datasets in the *Few-Class Regime*. By leveraging  
126 the MMPreTrain framework, users only need to specify the parameters of few-class subsets in the  
127 configuration files, which includes the list of models, datasets, number of classes ( $N_{CL}$ ), and the  
128 number of seeds ( $N_S$ ). *Few-Class Arena* generates the specific model and dataset configuration  
129 files for each subset, where subset classes are randomly extracted from the full set of classes, as  
130 specified by the seed number. Note that only one copy of the full, original dataset is maintained during  
131 the whole benchmarking life cycle because few-class subsets are created through the lightweight  
132 configurations, thus maximizing storage efficiency. We refer readers to the Appendix and the publicly  
133 released link for detailed implementations and use instructions.

#### 134 3.3 Many-Class Full Dataset Trained Benchmark

135 We conducted large-scale experiments spanning ten popular vision models (including CNN and  
136 ViT architectures) and ten common datasets<sup>3</sup>. Except for ImageNet1K, where pre-trained model  
137 weights are available, we train models in other datasets from scratch. While different models'

---

<sup>2</sup>Code is available at <https://github.com/fewclassarena/fca>, including detailed documentation and long-term plans of maintenance.

<sup>3</sup>Models include: ResNet50 (RN50), VGG16, ConvNeXt V2 (CNv2), Inception V3 (INCv3), EfficientNet V2 (EFv2), ShuffleNet V2 (SNv2), MobileNet V3 (MNv2), Vision Transformer base (ViTb), Swin Transformer V2 base (SWv2b) and MobileViT small (MViTs). Datasets include CalTech101 (CT101), CalTech256 (CT256), CIFAR100 (CF100), CUB200 (CB200), Food101 (FD101), GTSRB43, (GT43), ImageNet1K (IN1K), Indoor67 (ID67), Quickdraw345 (QD345) and Textures47 (TT47).

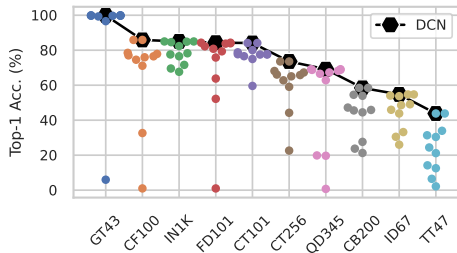
138 training procedures may incur various levels of complexity (particularly in our case for MobileNet  
 139 V3 and Swin Transformer V2 base), we have endeavored to minimize changes in the existing training  
 140 pipelines from MMPreTrain. The rationale is that if a model exhibits challenges in adapting it to a  
 141 dataset, then it is often not a helpful choice for a practitioner to select for deployment.

142 Results are summarized in Table 1. We make several key observations: (1) models in different datasets  
 143 (in rows) yield highly variable levels of performance by Top-1 accuracy; (2) no single best model  
 144 (bold, in columns) exists across all datasets; and (3) model rankings vary across various datasets.

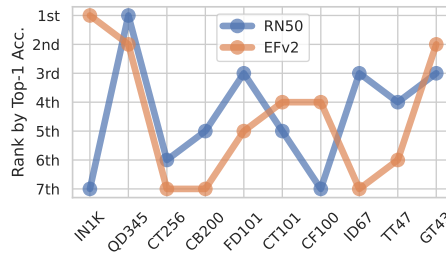
145 The first two observations are consistent with the findings in [25, 31]. For (1), it suggests there exists  
 146 underlying dataset-specific difficulty. To capture this characteristic, we adopt the reference dataset  
 147 classification difficulty number (DCN) [25] to refer to the empirically highest accuracy achieved in  
 148 a dataset from a finite number of models shown in Table 1 and Figure 2 (a). For observation (3),  
 149 we can examine the rankings among the ten models of ResNet50 and EfficientNet V2 in Figure 2  
 150 (b). ResNet50’s ranking varies dramatically for the different datasets, for instance ranking 7th on  
 151 ImageNet1K and 1st on Quickdraw345. This ranking variability is also observed in other models  
 152 (see all models in the Appendix). However, a common practice is to benchmark models – even for  
 153 efficiency – on large datasets, especially ImageNet1K. The varied dataset rankings in our experiments  
 154 expose the limitations of such a practice, further supporting our new benchmark paradigm, especially  
 155 in the *Few-Class Regime*. In later sections, we leverage DCN and image similarity for further analysis.

Dataset	RN50 [56]	VGG16 [57]	CNv2 [58]	INCv3 [59]	EFv2 [4]	SNv2 [9]	MNv3 [7]	ViTb [60]	SWv2b [61]	MViTs [10]	DCN [25]
GT43 [62]	99.85	96.60	99.83	99.78	<u>99.86</u>	<b>99.87</b>	5.98	99.31	99.78	99.69	99.87
CF100 [63]	74.56	71.12	<b>85.89</b>	75.97	<u>77.05</u>	77.89	1.00	32.65	<u>78.49</u>	76.51	85.89
IN1K [1]	76.55	71.62	<u>84.87</u>	77.57	<b>85.01</b>	69.55	67.66	82.37	84.6	78.25	85.01
FD101 [64]	83.76	75.82	63.80	<u>83.96</u>	80.82	79.36	0.99	52.21	<b>84.30</b>	82.23	84.30
CT101 [65]	77.70	74.99	77.52	<u>77.52</u>	77.82	<b>84.13</b>	76.58	59.59	78.82	80.06	84.13
CT256 [66]	65.07	59.08	<b>73.57</b>	66.09	62.80	<u>68.13</u>	22.63	44.23	67.28	<u>65.80</u>	73.57
QD345 [67]	<b>69.14</b>	19.86	62.86	68.25	<u>68.81</u>	67.32	0.72	19.67	66.54	68.76	69.14
CB200 [68]	45.86	21.26	27.61	45.58	44.48	53.95	47.22	23.73	<u>54.52</u>	<b>58.46</b>	58.46
ID67 [69]	53.75	26.01	33.21	45.95	43.85	<b>54.72</b>	49.10	30.51	48.58	<u>54.05</u>	54.72
TT47 [70]	30.43	12.55	6.49	14.20	21.17	<b>43.83</b>	2.18	31.38	<u>33.94</u>	24.41	43.83

Table 1: Top-1 accuracy across ten models in ten datasets. Models are trained and tested on full datasets with their original number of classes (e.g. 1K from ImageNet1K); this is denoted in the last few digits of the abbreviation of the dataset name. The best score is highlighted in bold while the second best is underlined for each dataset.



(a) Top-1 accuracy and DCN in ten full datasets.



(b) Ranking of ResNet50 (RN50) and EfficientNet V2 (EFv2) across 10 datasets by Top-1 acc.

Figure 2: Many-Class Full Dataset Benchmark.

156 In the next subsections, we introduce three new types of benchmarks: (1) Few-Class, Full Dataset  
 157 Trained Benchmark (FC-Full), which benchmarks vision models trained on the full dataset with the  
 158 original number of classes; (2) Few-Class, Subset Trained Benchmark (FC-Sub), which benchmarks  
 159 vision models trained on subsets of a fewer number of classes than the full dataset, and (3) Few-Class  
 160 Similarity Benchmark (FC-Sim), which benchmarks image similarity methods and their correlation  
 161 with model performance.

162 **3.4 Few-Class Full Dataset Trained Benchmark (FC-Full)**

163 Traditionally, a large number of models are trained and compared on many-class datasets. However,  
164 results for such benchmarks are not directly useful to the *Few-Class Regime* and many real-world  
165 scenarios. Therefore, we introduce the Few-Class Full Dataset Trained Benchmark (FC-Full), with the  
166 objective of effortlessly conducting large-scale experiments and analyses in the *Few-Class Regime*.

167 The procedure of FC-Full consists of two main stages. In the first stage, users select the models  
168 and datasets upon which they would like to conduct experiments. They can choose to download  
169 pre-trained model weights, which are usually available on popular model hubs (PyTorch Hub [71],  
170 TensorFlow Hub [72], Hugging Face [73], MMPreTrain [55] etc.). In case of no pre-trained weights  
171 available from public websites, users can resort to the option of training from scratch. To that end,  
172 our tool is designed and implemented to generate bash scripts for easily configurable and modifiable  
173 training through the use of configuration files.

174 In the second stage, users conduct benchmarking in the *Few-Class Regime*. By specifying the list of  
175 classes, *Few-Class Arena* automatically loads pre-trained weights of the chosen models and evaluates  
176 performance of the models on the selected datasets. Note that this process is accomplished through  
177 configuration files created by the user’s specifications, thus enabling hundreds of experiments to be  
178 launched by a single command. This dramatically reduces human effort that would otherwise be  
179 expended to run these experiments without *Few-Class Arena*.

180 **3.5 Few-Class Subset Trained Benchmark (FC-Sub)**

181 Our study in Figure 1 (red lines) reveals the limits of existing pre-trained models in the *Few-Class*  
182 *Regime*. To facilitate further research and analyze the upper bound performance in the *Few-Class*  
183 *Regime*, we introduce the Few-Class Subset Trained Benchmark (FC-Sub).

184 FC-Sub follows a similar procedure to FC-Full, except that, when evaluating a model in a subset with  
185 a specific number of classes, that model should have been trained on that same subset. Specifically, in  
186 Stage One (described for FC-Full), users specify models, datasets and the list of number of classes in  
187 configuration files. Then *Few-Class Arena* generates bash scripts for model training on each subset.  
188 In Stage two, *Few-Class Arena* tests each model in the same subset that it was trained on.

189 **3.6 Few-Class Similarity Benchmark (FC-Sim)**

190 One objective of our tool is to provide the Similarity Benchmark as a platform for researchers to  
191 design custom similarity scores for efficient comparison of models and datasets.

192 The intrinsic image difficulty of a dataset affects a model’s classification performance (and human)  
193 [74, 75, 32]. We show – as is intuitive – that the more similar two images are, the more difficult it is  
194 for a vision classifier to make a correct prediction. This suggests that the level of similarity of images  
195 in a dataset can be used as a proxy for a dataset difficulty measure. In this section, we first adopt and  
196 provide the basic formulation of similarity, the baseline of a similarity metric. Then we propose a  
197 Similarity-Based Silhouette Score to capture the characteristic of image similarity in a dataset.

198 We first adopt the basic similarity formulation from [51]. **Intra-Class Similarity**  $S_\alpha^{(C)}$  is defined as a  
199 scalar describing the similarity of images within a class by taking the average of all the distinct class  
200 pairs in  $C$ , while **Inter-Class Similarity** denotes a scalar describing the similarity among images in  
201 two different classes  $C_1$  and  $C_2$ . For a dataset  $D$ , these are defined as the mean of their similarity  
202 scores over all classes, respectively:

$$S_\alpha^{(D)} = \frac{1}{|L|} \sum_{l \in L} S_\alpha^{(C_l)} = \frac{1}{|L| \times |P^{(C_l)}|} \sum_{l \in L} \sum_{i, j \in C_l; i \neq j} \cos(\mathbf{Z}_i, \mathbf{Z}_j), \quad (1)$$

203

$$S_\beta^{(D)} = \frac{1}{|P^{(D)}|} \sum_{a, b \in L; a \neq b} S_\beta^{(C_a, C_b)} = \frac{1}{|P^{(D)}| \times |P^{(C_1, C_2)}|} \sum_{a, b \in L; a \neq b} \sum_{i \in C_1, j \in C_2} \cos(\mathbf{Z}_i, \mathbf{Z}_j). \quad (2)$$

204 where  $|L|$  is the number of classes in a dataset,  $Z_i$  is the visual feature of an image  $i$ ,  $|P^{(C)}|$  is the  
 205 total number of distinct image pairs in class  $C$ ,  $|P^{(D)}|$  is the total number of distinct class pairs, and  
 206  $|P^{(C_1, C_2)}|$  is the total number of distinct image pairs excluding same-class pairs.

207 Averaging these similarities provides a single scalar score at the class or dataset level. However,  
 208 this simplicity neglects other cluster-related information that can better reveal the underlying dataset  
 209 difficulty property of a dataset. In particular, the **(1) tightness of a class cluster** and **(2) distance to**  
 210 **other classes** of class clusters, are features that characterize the inherent class difficulty, but are not  
 211 captured by  $S_\alpha$  or  $S_\beta$  alone.

212 To compensate the aforementioned drawback, we adopt the Silhouette Score (SS) [34, 76]:  $SS(i) =$   
 213  $\frac{b(i) - a(i)}{\max(a(i), b(i))}$ , where  $SS(i)$  is the Silhouette Score of the data point  $i$ ,  $a(i)$  is the average dissimilarity  
 214 between  $i$  and other instances in the same class, and  $b(i)$  is the average dissimilarity between  $i$  and  
 215 other data points in the closest different class.

216 Observe that the above Intra-Class Similarity  $S_\alpha^{(C)}$  already represents the tightness of the class ( $C$ ),  
 217 therefore  $a(i)$  can be replaced with the inverse of Intra-Class Similarity  $a(i) = -S_\alpha(i)$ . For the  
 218 second term  $b(i)$ , we adopt the previously defined Inter-Class Similarity  $S_\beta^{(C_1, C_2)}$  and introduce a new  
 219 similarity score as **Nearest Inter-Class Similarity**  $S'_\beta^{(C)}$ , which is a scalar describing the similarity  
 220 among instances between class  $C$  and the closest class of each instance in  $C$ . The dataset-level  
 221 Nearest Inter-Class Similarity  $S'^{(D)}_\beta$  is expressed as:

$$S'^{(D)}_\beta = \frac{1}{|L|} \sum_{l \in L} S'^{(C_l, \hat{C}_l)}_\beta = \frac{1}{|L| \times |P^{(C_l, \hat{C}_l)}|} \sum_{l \in L} \sum_{i \in C_l, j \in \hat{C}_l} \cos(\mathbf{Z}_i, \mathbf{Z}_j). \quad (3)$$

222 where  $\hat{C}$  is the set of the nearest class to  $C$  ( $\hat{C} \neq C$ ). To summarize, we introduce our novel  
 223 **Similarity-Based Silhouette Score**  $SimSS^4$ :

$$SimSS^{(D)} = \frac{1}{|L| \times |C_l|} \sum_{i \in C_l} \frac{S_\alpha(i) - S'_\beta(i)}{\max(S_\alpha(i), S'_\beta(i))}. \quad (4)$$

## 224 4 Experimental Results

### 225 4.1 Results on FC-Full

226 In this section, we present the results of FC-Full. A model trained on the dataset with its original  
 227 number of classes (e.g. 1000 in ImageNet1K) is referred to as a *full-class model*. These experiments  
 228 are designed to understand how full-class model performance changes when the number of classes  
 229  $N_{Cl}$  decreases from many to few classes. We analyze the results of DCN-Full, shown in Figure 3  
 230 (details of all models are presented in the Appendix), and we make two key observations when  $N_{Cl}$   
 231 reduces to the *Few-Class Regime* (from right to left). (1) The best performing models do not always  
 232 increase its accuracy for fewer classes, as shown by the solid red lines that represent the average of  
 233 DCN for each  $N_{Cl}$ . (2) The variance, depicted by the light red areas, of the best models broaden  
 234 dramatically for low  $N_{Cl}$ , especially for  $N_{Cl} < 10$ .

235 Both observations support evidence of the limitations of using the common many-class benchmark  
 236 for application model selection in the *Few-Class Regime*, since it is not consistent between datasets  
 237 that a model can be made smaller with higher accuracy. Furthermore, the large variance in accuracy  
 238 means that prediction of performance for few classes is unreliable for this approach.

### 239 4.2 Results on FC-Sub

240 In this section, we show how using *Few-Class Arena* can help reveal more insights in the *Few-Class*  
 241 *Regime* to mitigate the issues of Section 4.1.

<sup>4</sup>The extended derivation is detailed in the Appendix.

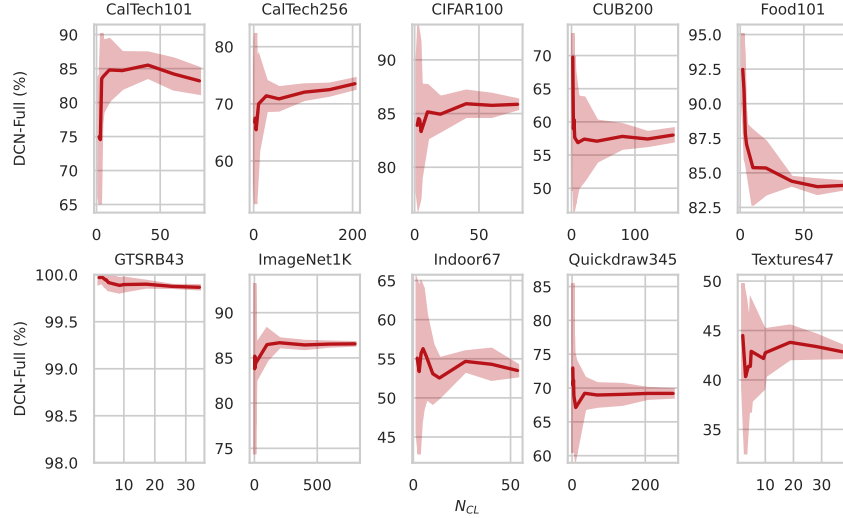


Figure 3: DCN-Full by Top-1 Accuracy (%).  $N_{CL}$  ranges from many to 2.

242 FC-Sub results are displayed in Figure 4. Recall that a *sub-class* model is a model trained on a subset  
 243 of the dataset where  $N_{CL}$  is smaller than the original number of classes in the full dataset. Observe  
 244 that in the *Few-Class Regime* (when  $N_{CL}$  decreases from 4 to 2) that: (1) DCN increases as shown by  
 245 the solid blue lines, and (2) variance reduces as displayed by the light blue areas.

246 The preceding observation for FC-Full 4.1 seems to contradict the common belief that, the fewer the  
 247 classes, the higher is the accuracy that a model can achieve. Conversely, the FC-Sub results do align  
 248 with this belief. We argue that a full-class model needs to accommodate many parameters to learn  
 249 features that will enable high performance across all classes in a many-class, full dataset. With the  
 250 same parameters, however, a sub-class model can adapt to finer and more discriminative features that  
 251 improve its performance when the number of target classes are much smaller.

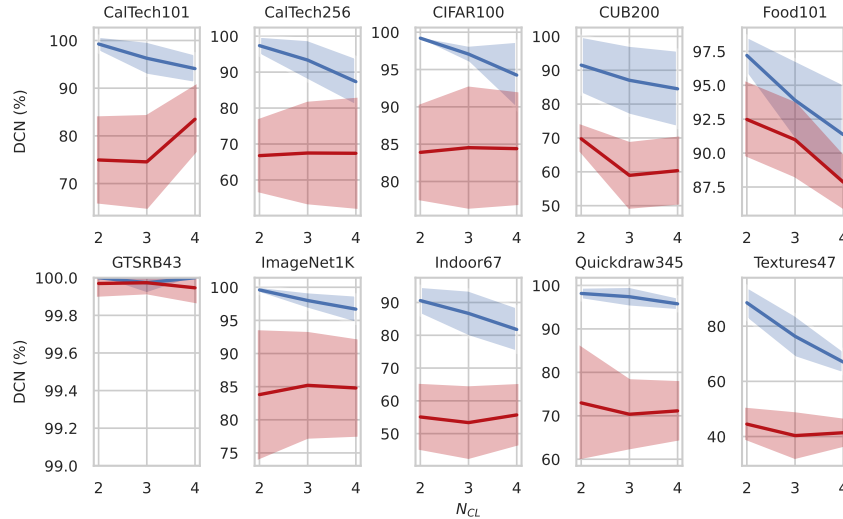


Figure 4: DCN-Sub (red) and DCN-Full (blue) by Top-1 Accuracy (%).  $N_{CL}$  ranges from 2 to 4.

### 252 4.3 Results on FC-Sim

253 In this section, we analyze the use of SimSS (Equation 4) as proxy for few-class dataset difficulty.  
 254 Experiments are conducted on ImageNet1K using the ResNet family for the lower  $N_{CL} \leq 10\%$  range  
 255 of the original 1000 classes,  $N_{CL} \in \{2, 3, 4, 5, 10, 100\}$ , and the results are shown in Figure 5. Each  
 256 datapoint of DCN-Full (diamond in red) or DCN-Sub (square in blue) represents an experiment in a



257 subset of a specific  $N_{CL}$ , where classes are sampled from the full dataset. For reproducible results,  
 258 we use seed numbers from 0 to 4 to generate 5 subsets for one  $N_{CL}$  by default. A similarity base  
 259 function ( $sim()$ ) is defined as the atomic function that takes a pair of images as input and outputs a  
 260 scalar that represents their image similarity.

261 In our experiments, we leverage the general visual feature extraction ability of CLIP (image + text)  
 262 [52] and DINOv2 (image) [53] by self-supervised learning. Specifically, a pair of images are fed into  
 263 its latent space from which the cosine score is calculated and normalized to 0 to 1. Note that we  
 264 only use the Image Encoder in CLIP.

265 **Comparing Accuracy and Similarity** To evaluate SimSS, we compute the Pearson correlation  
 266 coefficient (PCC) ( $r$ ) between model accuracy and SimSS. Results in Figure 5 (a) (b) show that  
 267 SimSS is poorly correlated with DCN-Full ( $r = 0.18$  and  $r = 0.26$  for CLIP and DINOv2) due to the  
 268 large variance shown in Section 4.1. In contrast, SimSS is highly correlated with DCN-Sub (shown  
 269 in blue squares), with  $r = 0.90$  and  $r = 0.88$  using CLIP (dashed) and DINOv2 (solid), respectively.  
 270 The high PCC [77, 78] demonstrates that SimSS is a reliable metric to estimate few-class dataset  
 271 difficulty, and this can help predict the empirical upper-bound accuracy of a model in the *Few-Class*  
 272 *Regime*. Comparison between SimSS and all models can be found in the Appendix. Such a high  
 273 correlation suggests this offers a reliable scaling relationship to estimate model accuracy by similarity  
 274 for other values of  $N_{CL}$  without an exhaustive search. Due to the dataset specificity of the dataset  
 275 difficulty property, this score is computed once and used for all times the same dataset is used. We  
 276 have made available difficulty scores for many datasets at the *Few-Class Arena* site.

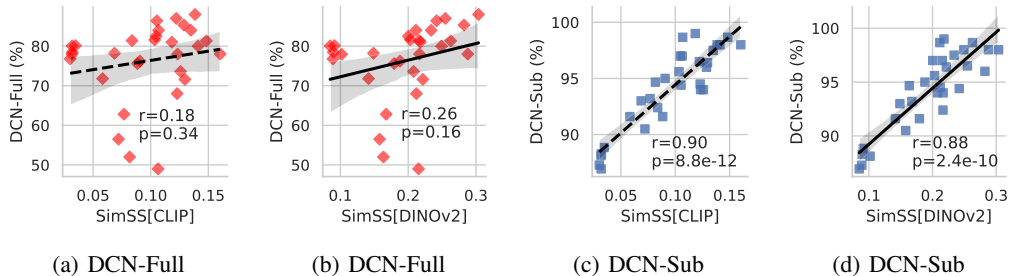


Figure 5: Pearson correlation coefficient ( $r$ ) between DCN and SimSS when  $N_{Cl} \in \{2, 3, 4, 5, 10, 100\}$ . DCN-Sub (blue squares) is more highly correlated than DCN-Full (red diamonds) with SimSS using both similarity base functions of CLIP (dashed line) and DINOv2 (solid line) with  $r \geq 0.88$ .

## 277 5 Conclusion

278 We have proposed *Few-Class Arena* and a dataset difficulty measurement, which together form  
 279 a benchmark tool to compare and select efficient models in the *Few-Class Regime*. Extensive  
 280 experiments and analyses over 1500 experiments with 10 models on 10 datasets have helped identify  
 281 new behavior that is specific to the *Few-Class Regime* as compared to for many-classes. One finding  
 282 reveals a new  $n_{Cl}$ -scaling law whereby dataset difficulty must be taken into consideration for accuracy  
 283 prediction. Such a benchmark will be valuable to the community by providing both researchers and  
 284 practitioners with a unified framework for future research and real applications.

285 **Limitations and Future Work.** We note that the convergence of sub-models is contingent on various  
 286 factors in a training scheduler, such as learning rate. A careful tuning of training procedure may  
 287 increase a model’s performance, but it shouldn’t change the classification difficulty number drastically  
 288 since this represents a dataset’s intrinsic difficulty property. The current difficulty benchmark supports  
 289 image similarity while in the future it can be expanded to other difficulty measurements [25]. As  
 290 CLIP and DINOv2 are trained toward general visual features, it is unclear if they will be appropriate  
 291 for other types of images such as sketches without textures in Quickdraw [67]. For this reason, a  
 292 universal similarity foundation model would be appealing that applies to any image type. In summary,  
 293 *Few-Class Arena* identifies a promising new path to achieve efficiencies that are focused on the  
 294 important and practical *Few-Class Regime*, establishing this as a baseline for future work.

## 295 References

- 296 [1] Deng, J., W. Dong, R. Socher, et al. Imagenet: A large-scale hierarchical image database. In  
297 *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- 298 [2] Lin, T.-Y., M. Maire, S. Belongie, et al. Microsoft coco: Common objects in context. In  
299 *European conference on computer vision*, pages 740–755. Springer, 2014.
- 300 [3] Tan, M., Q. Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In  
301 *International conference on machine learning*, pages 6105–6114. PMLR, 2019.
- 302 [4] Tan, M., Q. L. Efficientnetv2: Smaller models and faster training. In *International conference*  
303 *on machine learning*, pages 10096–10106. PMLR, 2021.
- 304 [5] Sinha, D., M. El-Sharkawy. Thin mobilenet: An enhanced mobilenet architecture. In *2019*  
305 *IEEE 10th annual ubiquitous computing, electronics & mobile communication conference*  
306 *(UEMCON)*, pages 0280–0285. IEEE, 2019.
- 307 [6] Sandler, M., A. Howard, M. Zhu, et al. Mobilenetv2: Inverted residuals and linear bottlenecks.  
308 In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages  
309 4510–4520. 2018.
- 310 [7] Howard, A., M. Sandler, G. Chu, et al. Searching for mobilenetv3. In *Proceedings of the*  
311 *IEEE/CVF international conference on computer vision*, pages 1314–1324. 2019.
- 312 [8] Iandola, F. N., S. Han, M. W. Moskewicz, et al. Squeezenet: Alexnet-level accuracy with 50x  
313 fewer parameters and < 0.5 mb model size. *arXiv preprint arXiv:1602.07360*, 2016.
- 314 [9] Ma, N., X. Zhang, H.-T. Zheng, et al. Shufflenet v2: Practical guidelines for efficient cnn  
315 architecture design. In *Proceedings of the European conference on computer vision (ECCV)*,  
316 pages 116–131. 2018.
- 317 [10] Mehta, S., M. Rastegari. Mobilevit: Light-weight, general-purpose, and mobile-friendly vision  
318 transformer. arxiv 2021. *arXiv preprint arXiv:2110.02178*.
- 319 [11] Patterson, D., J. Gonzalez, Q. Le, et al. Carbon emissions and large neural network training.  
320 *arXiv preprint arXiv:2104.10350*, 2021.
- 321 [12] Rae, J. W., S. Borgeaud, T. Cai, et al. Scaling language models: Methods, analysis & insights  
322 from training gopher. *arXiv preprint arXiv:2112.11446*, 2021.
- 323 [13] Gysel, P., J. Pimentel, M. Motamedi, et al. Ristretto: A framework for empirical study of  
324 resource-efficient inference in convolutional neural networks. *IEEE transactions on neural*  
325 *networks and learning systems*, 29(11):5784–5789, 2018.
- 326 [14] Han, S., H. Mao, W. J. Dally. Deep compression: Compressing deep neural networks with  
327 pruning, trained quantization and huffman coding. *arXiv preprint arXiv:1510.00149*, 2015.
- 328 [15] Leng, C., Z. Dou, H. Li, et al. Extremely low bit neural network: Squeeze the last bit out with  
329 adm. In *Proceedings of the AAAI conference on artificial intelligence*, vol. 32. 2018.
- 330 [16] Cheng, Y., D. Wang, P. Zhou, et al. A survey of model compression and acceleration for deep  
331 neural networks. *arXiv preprint arXiv:1710.09282*, 2017.
- 332 [17] Blalock, D., J. J. Gonzalez Ortiz, J. Frankle, et al. What is the state of neural network pruning?  
333 *Proceedings of machine learning and systems*, 2:129–146, 2020.
- 334 [18] Li, H., A. Kadav, I. Durdanovic, et al. Pruning filters for efficient convnets. *arXiv preprint*  
335 *arXiv:1608.08710*, 2016.
- 336 [19] Shen, M., H. Yin, P. Molchanov, et al. Structural pruning via latency-saliency knapsack. *arXiv*  
337 *preprint arXiv:2210.06659*, 2022.
- 338 [20] Yeung, S., O. Russakovsky, G. Mori, et al. End-to-end learning of action detection from frame  
339 glimpses in videos. In *Proceedings of the IEEE conference on computer vision and pattern*  
340 *recognition*, pages 2678–2687. 2016.

- 341 [21] Shao, Z., L. Wang, Z. Wang, et al. Saliency-aware convolution neural network for ship detection  
342 in surveillance video. *IEEE Transactions on Circuits and Systems for Video Technology*,  
343 30(3):781–794, 2020.
- 344 [22] A. Delplanque, P. L. J. L. J. T., S. Foucher. Multispecies detection and identification of african  
345 mammals in aerial imagery using convolutional neural networks. *Remote Sensing in Ecology  
346 and Conservation*, 8(April):166–179, 2022.
- 347 [23] Cai, Y., T. Luan, H. Gao, et al. Yolov4-5d: An effective and efficient object detector for  
348 autonomous driving. *IEEE Transactions on Instrumentation and Measurement*, 70:1–13, 2021.
- 349 [24] Scheidegger, F., L. Benini, C. Bekas, et al. Constrained deep neural network architecture search  
350 for iot devices accounting for hardware calibration. *Advances in Neural Information Processing  
351 Systems*, 32, 2019.
- 352 [25] Scheidegger, F., R. Istrate, G. Mariani, et al. Efficient image dataset classification difficulty  
353 estimation for predicting deep-learning accuracy. *The Visual Computer*, 37(6):1593–1610, 2021.
- 354 [26] Kaplan, J., S. McCandlish, T. Henighan, et al. Scaling laws for neural language models. *arXiv  
355 preprint arXiv:2001.08361*, 2020.
- 356 [27] Lee, T., M. Yasunaga, C. Meng, et al. Holistic evaluation of text-to-image models. *Advances in  
357 Neural Information Processing Systems*, 36, 2024.
- 358 [28] Le, T., V. Lal, P. Howard. Coco-counterfactuals: Automatically constructed counterfactual  
359 examples for image-text pairs. *Advances in Neural Information Processing Systems*, 36, 2024.
- 360 [29] Laurençon, H., L. Saulnier, L. Tronchon, et al. Obelics: An open web-scale filtered dataset of  
361 interleaved image-text documents. *Advances in Neural Information Processing Systems*, 36,  
362 2024.
- 363 [30] Bitton, Y., N. Bitton Guetta, R. Yosef, et al. Winogavil: Gamified association benchmark to  
364 challenge vision-and-language models. *Advances in Neural Information Processing Systems*,  
365 35:26549–26564, 2022.
- 366 [31] Fang, A., S. Kornblith, L. Schmidt. Does progress on imagenet transfer to real-world datasets?  
367 *Advances in Neural Information Processing Systems*, 36, 2024.
- 368 [32] Mayo, D., J. Cummings, X. Lin, et al. How hard are computer vision datasets? calibrating dataset  
369 difficulty to viewing time. *Advances in Neural Information Processing Systems*, 36:11008–  
370 11036, 2023.
- 371 [33] Goldblum, M., H. Souri, R. Ni, et al. Battle of the backbones: A large-scale comparison of  
372 pretrained models across computer vision tasks. *Advances in Neural Information Processing  
373 Systems*, 36, 2024.
- 374 [34] Rousseeuw, P. J. Silhouettes: a graphical aid to the interpretation and validation of cluster  
375 analysis. *Journal of computational and applied mathematics*, 20:53–65, 1987.
- 376 [35] Dowson, D., B. Landau. The fréchet distance between multivariate normal distributions. *Journal  
377 of multivariate analysis*, 12(3):450–455, 1982.
- 378 [36] Heusel, M., H. Ramsauer, T. Unterthiner, et al. Gans trained by a two time-scale update rule  
379 converge to a local nash equilibrium. *Advances in neural information processing systems*, 30,  
380 2017.
- 381 [37] Lucic, M., K. Kurach, M. Michalski, et al. Are gans created equal? a large-scale study. *Advances  
382 in neural information processing systems*, 31, 2018.
- 383 [38] Hore, A., D. Ziou. Image quality metrics: Psnr vs. ssim. In *2010 20th international conference  
384 on pattern recognition*, pages 2366–2369. IEEE, 2010.
- 385 [39] Wang, Z., A. C. Bovik, H. R. Sheikh, et al. Image quality assessment: from error visibility to  
386 structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004.
- 387 [40] Sheikh, H. R., A. C. Bovik. Image information and visual quality. *IEEE Transactions on image  
388 processing*, 15(2):430–444, 2006.

- 389 [41] Kar, K., J. Kubilius, K. Schmidt, et al. Evidence that recurrent circuits are critical to the ventral  
390 stream’s execution of core object recognition behavior. *Nature neuroscience*, 22(6):974–983,  
391 2019.
- 392 [42] Jiang, Z., C. Zhang, K. Talwar, et al. Characterizing structural regularities of labeled data in  
393 overparameterized models. *arXiv preprint arXiv:2002.03206*, 2020.
- 394 [43] Baldock, R., H. Maennel, B. Neyshabur. Deep learning through the lens of example difficulty.  
395 *Advances in Neural Information Processing Systems*, 34:10876–10889, 2021.
- 396 [44] Goodfellow, I. J., J. Shlens, C. Szegedy. Explaining and harnessing adversarial examples. *arXiv*  
397 *preprint arXiv:1412.6572*, 2014.
- 398 [45] Arun, S. Turning visual search time on its head. *Vision Research*, 74:86–92, 2012.
- 399 [46] Trick, L. M., J. T. Enns. Lifespan changes in attention: The visual search task. *Cognitive*  
400 *Development*, 13(3):369–386, 1998.
- 401 [47] Wolfe, J. M., E. M. Palmer, T. S. Horowitz. Reaction time distributions constrain models of  
402 visual search. *Vision research*, 50(14):1304–1311, 2010.
- 403 [48] Zhang, D., G. Lu. Evaluation of similarity measurement for image retrieval. In *International*  
404 *conference on neural networks and signal processing, 2003. proceedings of the 2003*, vol. 2,  
405 pages 928–931. IEEE, 2003.
- 406 [49] Wang, J., Y. Song, T. Leung, et al. Learning fine-grained image similarity with deep ranking.  
407 In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages  
408 1386–1393. 2014.
- 409 [50] Tudor Ionescu, R., B. Alexe, M. Leordeanu, et al. How hard can it be? estimating the difficulty  
410 of visual search in an image. In *Proceedings of the IEEE Conference on Computer Vision and*  
411 *Pattern Recognition*, pages 2157–2166. 2016.
- 412 [51] Cao, B. B., L. O’Gorman, M. Coss, et al. Data-side efficiencies for lightweight convolutional  
413 neural networks. *arXiv preprint arXiv:2308.13057*, 2023.
- 414 [52] Radford, A., J. W. Kim, C. Hallacy, et al. Learning transferable visual models from natural  
415 language supervision. In *International conference on machine learning*, pages 8748–8763.  
416 PMLR, 2021.
- 417 [53] Oquab, M., T. Darcet, T. Moutakanni, et al. Dinov2: Learning robust visual features without  
418 supervision. *arXiv preprint arXiv:2304.07193*, 2023.
- 419 [54] Kumar, M., N. Houlsby, N. Kalchbrenner, et al. Do better imagenet classifiers assess perceptual  
420 similarity better? *arXiv preprint arXiv:2203.04946*, 2022.
- 421 [55] Contributors, M. Openmmlab’s pre-training toolbox and benchmark. [https://github.com/  
422 open-mmlab/mmpretrain](https://github.com/open-mmlab/mmpretrain), 2023.
- 423 [56] He, K., X. Zhang, S. Ren, et al. Deep residual learning for image recognition. In *Proceedings*  
424 *of the IEEE conference on computer vision and pattern recognition*, pages 770–778. 2016.
- 425 [57] Simonyan, K., A. Zisserman. Very deep convolutional networks for large-scale image recogni-  
426 tion. *arXiv preprint arXiv:1409.1556*, 2014.
- 427 [58] Woo, S., S. Debnath, R. Hu, et al. Convnext v2: Co-designing and scaling convnets with masked  
428 autoencoders. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*  
429 *Recognition*, pages 16133–16142. 2023.
- 430 [59] Szegedy, C., V. Vanhoucke, S. Ioffe, et al. Rethinking the inception architecture for computer  
431 vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,  
432 pages 2818–2826. 2016.
- 433 [60] Dosovitskiy, A., L. Beyer, A. Kolesnikov, et al. An image is worth 16x16 words: Transformers  
434 for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.

- 435 [61] Liu, Z., H. Hu, Y. Lin, et al. Swin transformer v2: Scaling up capacity and resolution. In  
436 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages  
437 12009–12019. 2022.
- 438 [62] Stallkamp, J., M. Schlipsing, J. Salmen, et al. Man vs. computer: Benchmarking machine  
439 learning algorithms for traffic sign recognition. *Neural Networks*, (0):–, 2012.
- 440 [63] Krizhevsky, A., G. Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- 441 [64] Bossard, L., M. Guillaumin, L. Van Gool. Food-101 – mining discriminative components with  
442 random forests. In *European Conference on Computer Vision*. 2014.
- 443 [65] Li, F.-F., M. Andreeto, M. Ranzato, et al. Caltech 101, 2022.
- 444 [66] Griffin, G., A. Holub, P. Perona. Caltech 256, 2022.
- 445 [67] Ha, D., D. Eck. A neural representation of sketch drawings. *arXiv preprint arXiv:1704.03477*,  
446 2017.
- 447 [68] Wah, C., S. Branson, P. Welinder, et al. *The Caltech-UCSD Birds-200-2011 Dataset*. 2011.
- 448 [69] Quattoni, A., A. Torralba. Recognizing indoor scenes. In *2009 IEEE conference on computer  
449 vision and pattern recognition*, pages 413–420. IEEE, 2009.
- 450 [70] Cimpoi, M., S. Maji, I. Kokkinos, et al. Describing textures in the wild. In *Proceedings of the  
451 IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 2014.
- 452 [71] Foundation, T. L. Pytorch hub. <https://pytorch.org/hub>, 2024. Accessed on 2024-06-04.
- 453 [72] Inc., G. Tensorflow hub. <https://www.tensorflow.org/hub>, 2024. Accessed on 2024-06-  
454 04.
- 455 [73] Face, H. Hugging face models. <https://huggingface.co/models>, 2024. Accessed on  
456 2024-06-04.
- 457 [74] Geirhos, R., D. H. Janssen, H. H. Sch"utt, et al. Comparing deep neural networks against  
458 humans: object recognition when the signal gets weaker. *arXiv preprint arXiv:1706.06969*,  
459 2017.
- 460 [75] Rajalingham, R., E. B. Issa, P. Bashivan, et al. Large-scale, high-resolution comparison of the  
461 core visual object recognition behavior of humans, monkeys, and state-of-the-art deep artificial  
462 neural networks. *Journal of Neuroscience*, 38(33):7255–7269, 2018.
- 463 [76] Shahapure, K. R., C. Nicholas. Cluster quality analysis using silhouette score. In *2020 IEEE  
464 7th international conference on data science and advanced analytics (DSAA)*, pages 747–748.  
465 IEEE, 2020.
- 466 [77] Wicklin, R. Weak or strong? how to interpret a spearman or kendall correlation. [https:  
467 //blogs.sas.com/content/iml/2023/04/05/interpret-spearman-kendall-corr.  
468 html](https://blogs.sas.com/content/iml/2023/04/05/interpret-spearman-kendall-corr.html), 2024. Accessed on 2024-06-04.
- 469 [78] Schober, P., C. Boer, L. A. Schwarte. Correlation coefficients: appropriate use and interpretation.  
470 *Anesthesia & analgesia*, 126(5):1763–1768, 2018.