
Exploring Dataset-Scale Indicators of Data Quality

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Modern computer vision foundation models are trained on massive amounts of data,
2 incurring large economic and environmental costs. Recent research has suggested
3 that improving data quality can significantly reduce the need for data quantity.
4 But what constitutes data quality in computer vision? We posit that the quality
5 of a given dataset can be decomposed into distinct sample-level and dataset-level
6 constituents, and that the former have been more extensively studied than the
7 latter. We ablate the effects of two important dataset-level constituents: label set
8 design, and class balance. By monitoring these constituents using key indicators
9 we provide, researchers and practitioners can better anticipate model performance,
10 measured in terms of its accuracy and robustness to distribution shifts.

11 1 Introduction

12 1.1 Motivation

13 In recent years, the field of computer vision has seen remarkable progress in a range of sub-
14 disciplines [Li et al., 2023; Mildenhall et al., 2021; Rombach et al., 2021]. Much of this progress has
15 derived from *foundation models* [Yuan et al., 2021; Radford et al., 2021]. Modern computer vision
16 foundation models are trained on massive amounts of data, and the general trend has been to achieved
17 improved performance by scaling up dataset size.

18 However, increasing the size of the dataset is not the only way to improve the downstream performance
19 of a given model. Recent work has introduced into the literature the importance of *data quality*, a
20 phrase which is meant to encompass all facets of a dataset which impact downstream performance,
21 aside from its size [Nguyen et al., 2022; Santurkar et al., 2022; Gadre et al., 2023].

22 The question of *which* facets of data quality are relevant for improved performance, however, is
23 under-explored in the literature. Even less frequently studied is whether there exist *quantifiable*,
24 *predictive indicators* of these facets. If reliable indicators of dataset quality can be found prior to
25 training, then researchers can better estimate the impact of modifications to their datasets (or even
26 guide dataset design), and ultimately lead to reduction of human, environmental, and economic costs
27 of large model training [Sharir et al., 2020; Changpinyo et al., 2021; Bommasani et al., 2022].

28 1.2 Our Contributions

29 In this (preliminary) work, we present several new findings on the constituents of dataset quality.
30 Unlike most previous works, our aim is to discover indicators that are properties at the dataset-scale,
31 rather than measures of image-level quality. Via a series of controlled ablation studies, we explore the
32 downstream performance impact of certain constituents, and show that these can be used as predictive
33 indicators of the quality of datasets prior to model training. In brief:

- 34 1. We conduct thorough ablation studies on two important dataset properties — *label set size*, and
35 *class imbalance* — and analyze the (sometimes) complex effects they have on downstream metrics
36 in image classification tasks.
37 2. For each of the above two constituents, we provide a list of *key indicators* which are predictive
38 of model performance. Our indicators are inexpensive to compute, scalable to very large dataset
39 sizes, and can be identified prior to model training.

40 Our results can be viewed as building blocks towards a systematic taxonomy of the notion of “quality”
41 at the dataset scale, which may enable improved design choices for datasets in computer vision.

42 2 Related Work

43 The question of how best to curate a dataset from a raw, large set of images is an important problem
44 in computer vision. The common approach is to use full supervision; classes are chosen in advance,
45 raw samples are algorithmically/manually filtered, and then manually labeled by annotators, and
46 (sometimes) class balance is enforced; the classic example is ImageNet [Deng et al., 2009]. Such
47 datasets tend to produce very strong baselines [Kornblith et al., 2019]. Other approaches take a
48 more relaxed approach to filtration, assigning unfiltered web-scraped images to human labelers and
49 applying a wide range of possible class labels [Kuznetsova et al., 2018].

50 As datasets scale from millions to billions of samples, human labeling becomes impractical. While
51 early approaches such as Thomee et al. [2016] initially curated datasets without any supervision,
52 *weak supervision* strategies have now become popular. With no human in the loop, proxy measures
53 of quality become essential. Sample-level quality indicators include encoding format, size, aspect
54 ratio, and offensive content [Sharma et al., 2018]. Labeling strategies for these images sometimes
55 rely on image tags from large social image-sharing sites [Sun et al., 2017; Mahajan et al., 2018].

56 One challenge in studying dataset curation in computer vision derives from the fact that many of the
57 largest image datasets cited in the literature are not publicly available [Sun et al., 2017; Mahajan
58 et al., 2018; Nguyen et al., 2022]. A notable recent exception to this is the work of Schuhmann et al.
59 [2021]. Subsequent works such as [Ilharco et al., 2021; Gadre et al., 2023; Feuer et al., 2023] have
60 taken advantage of these developments to train new models and design data-centric challenges.

61 There have been many proposed approaches to predict the behavior of models for a given test set,
62 coalescing in a NeurIPS competition in December 2020 [Jiang et al., 2021]. Broadly, submitted
63 methods fell into three meta-categories: (i) generalization measures derived from theoretically
64 motivated generalization bounds; (ii) data augmentation methods, which estimate the generalization
65 error of a trained model by computing its accuracy on synthetic data, and (iii) geometry and statistics
66 of intermediate representations. The most successful approach was that of [Natekar & Sharma, 2020],
67 which used a combination of (ii) and (iii).

68 An important dimension in which foundation models have been found to outperform smaller computer
69 vision models — and which we use as a key metric in our experiments — is *distributional robustness*,
70 a test-time paradigm which aims to estimate model robustness to distribution shifts [Recht et al.,
71 2019]. A *distribution shift* is defined as evaluation data which differs from the data on which a
72 model was trained due to natural factors. Real world image classifiers require predictable model
73 behavior under such shifts. Models trained on large, heterogeneous datasets tend to provide greater
74 distributional robustness than their counterparts trained on less data [Feuer et al., 2023].

75 3 Experimental Setup

76 **Measurable indicators of dataset quality.** As discussed above in Sec. 2, most existing works which
77 attempt to evaluate data quality do so at the *sample level*; they focus on properties pertaining to
78 individual samples, such as image resolution and image fidelity. Our focus instead is on exploring
79 holistic *dataset-level* properties, which are typically determined by the dataset’s creators and (if that
80 dataset is used as a benchmark) are typically treated as immutable.

81 In this short paper, we focus on understanding the predictive power of two such properties: (i) number
82 of classes (label set size), and (ii) number of samples per class (individual class size). We chose these
83 properties since they are ubiquitous, concretely measurable, and tend to have large impacts on the
84 performance of the corresponding trained models.

85 **Data scaling.** We define *horizontal scaling* (H-scaling) as scaling up the label set size in the dataset
86 while holding individual class size constant. We define *vertical scaling* (V-scaling) as scaling up
87 individual class size while holding the label set size constant.

88 **Architecture.** Our baseline architecture against which all variations are compared is a modified
89 ResNet-50 with a 1000-class linear classification head [He et al., 2015]. The specifics of the
90 modifications are described in [Ilharco et al., 2021]. Whenever we use an architecture other than our
91 standard baseline, we refer to it by its name in the timm library from Wightman [2019].

92 **Pretraining datasets.** The experiments described in this paper were primarily conducted using
93 the JANuS dataset, introduced by [Feuer et al., 2023]. JANuS is a composite image-caption dataset
94 sourced from four different datasets of origin. Every sample in JANuS contains one or more labels (for
95 training conventional computer vision models) and one or more captions (for training vision-language
96 models) with varying degrees of supervision. These properties make it an ideal dataset for conducting
97 controlled experiments for H-scaling and V-scaling. In addition to JANuS, we make use of LAION,
98 ImageNet and OpenImages, from Deng et al. [2009]; Kuznetsova et al. [2018]; Schuhmann et al.
99 [2021], respectively; further details are provided in the experiments below.

100 **Training details.** In our experiments, we train with mixed precision, at a batch size of 256, and do
101 not use gradient clipping. We use the AMP library to implement the training process. Learning rate is
102 chosen once, via grid search, for each new architecture / dataset pair. Models are typically distributed
103 across a single node with 4 NVIDIA A100 GPUs. All models are trained for 256 epochs. Following
104 Santurkar et al. [2022], we use SimCLR augmentations (resize, crop, flip, jitter, blur, grayscale) rather
105 than CLIP augmentations (resize and crop) for model training. A few of our models are not trained
106 from scratch, but are instead evaluated zero-shot using weights sourced from Wightman [2019]; we
107 note this whenever it is the case.

108 **Labels.** We evaluate on a single, broad-scope label set of 100 classes corresponding to ImageNet-100
109 (IN100), which is the first constituent dataset of JANuS. In our tables, we refer to the validation
110 set for IN100 as IN100-Val, and the average shift accuracy as IN100-Avg. Rob. OI100 refers to
111 OpenImages-100, the second constituent dataset of JANuS.

112 **Distribution Shifts.** Following the literature, whenever we measure robustness, we report the average
113 of four shifts on ImageNet. *ImageNet-V2* was designed to duplicate, as closely as possible, the
114 original ImageNet test set [Recht et al., 2019]. *Imagenet-Sketch* is a distribution shift covering
115 sketches, paintings, drawings and illustrations [Wang et al., 2019]. *Imagenet-R* is a 200-class subset
116 of ImageNet-2012 focused on renditions of everyday objects [Hendrycks et al., 2021a]. *Imagenet-A*
117 is a 200-class subset of ImageNet-2012 which was algorithmically selected [Hendrycks et al., 2021b].

118 **Data filtration.** We define *data filtration* as any strategy which sub-selects from a larger pool of
119 possible samples. A simple example of a filtration strategy would be conducting a web search for the
120 target classes in a dataset, and selecting the first k samples in the search.

121 **Metrics for distributional robustness.** Our primary metric is *average robustness* (abbrev: Avg. Rob.),
122 which is the average test-set accuracy of a model on a set of distribution shifts. Although this measure
123 is easy to interpret, it can conceal substantial performance differences between shifts.

124 4 Results

125 The fact that both label set size and class (im)balance impact image classification models should be
126 folklore to computer vision practitioners. However, to the best of our knowledge, these two properties
127 have not been directly contrasted, given an overall data budget. We address the following questions:

- 128 1. **For a given overall data budget, is it better to scale up individual class sizes (V-scaling), or to**
129 **scale the number of classes (H-scaling)?** See Sec. 4.1.
- 130 2. **How does class imbalance impact accuracy and robustness?** See Sec. 4.2.

131 4.1 Label Set Size

132 Scaling up dataset size has become the *de facto* driving force for improving the accuracy (and
133 robustness) of image classification models. But what is the *right* way to scale up datasets: should we
134 just scale up samples per class? Or are there benefits if the model is trained on a larger set of classes?

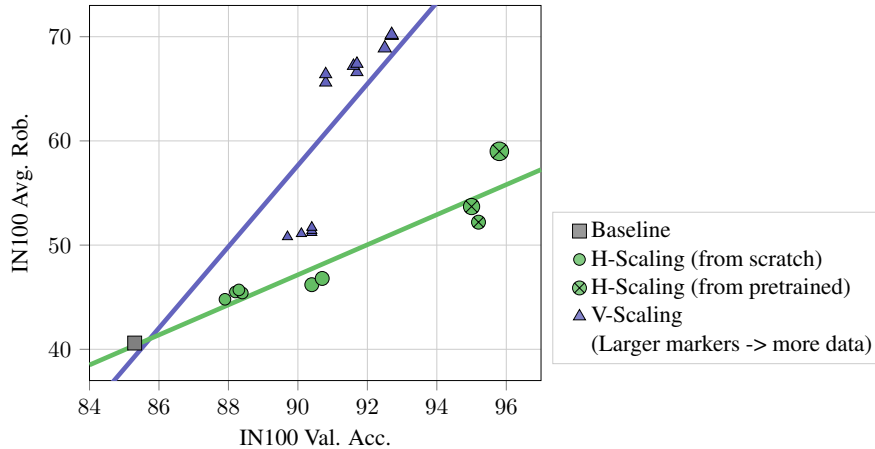


Figure 1: **Models benefit from scaling data horizontally, as well as vertically.** Vertical scaling (increasing the number of samples for in-distribution classes) is generally understood to improve accuracy and robustness. We show that horizontal scaling, increasing the number of out-of-distribution samples, also helps. Dot size represents size of training dataset (vertical scaling is ID, horizontal is OOD). Darker points represent models with more parameters. Image best viewed in color.

135 To fairly compare these two choices, we keep the test label set constant. The net effect of the latter
 136 case is that the model sees image examples that are OOD with respect to the test set, and we zero out
 137 the logits of the OOD classes at test time.

138 We enumerate results on two distinct sets of large scale classification models. All vertically scaled
 139 models are trained by us. For horizontal scaling, we employ a mix of pretrained and from-scratch
 140 (with the larger models generally evaluated from pretrained checkpoints). We distinguish between
 141 them in Fig. 1 with distinct labels.

142 4.1.1 Key Indicators

143 A natural indicator for vertical scaling is per-class dataset size (though this approach demands separate
 144 attention to class balance; see Sec. 4.2). Training label set size is the indicator for horizontal scaling.

145 4.1.2 Findings

146 Somewhat surprisingly, we find that models achieve more robust and accurate overall representations
 147 via horizontal scaling. It is well known that models benefit from seeing more ID samples during

Model	Dataset Size	Training Stages	IN100-Val / Avg. Rob.
resnetv2_50x1_bit	15.4 Mn	PT 21800 FT 1000 ZS 100	95.8% / 59.0%
resnetv2_50	12.6 Mn	PT 1000 FT 12000 ZS 100	95.0% / 53.7%
resnetv2_50	1.2 Mn	PT 1000 ZS 100	95.2% / 52.2%
swin_base_patch4_window7	14.2 Mn	PT 21800 ZS 100	86.7% / 43.4%
flexivit_base_1000ep	14.2 Mn	PT 21800 ZS 100	68.9% / 18.8%
convnext_base	14.2 Mn	PT 21800 ZS 100	66.4% / 36.2%
resnetv2_50x1_bit	14.2 Mn	PT 21800 ZS 100	26.7% / 5.2%

Table 1: **Horizontal scaling with large label sets.** We find that horizontal scaling with large label sets works much better when training occurs in multiple stages. Surprisingly, the process works almost as well whether the larger label set comes first or second; label set size is not an inherent obstacle to the success of horizontal scaling. Model names correspond to those found in the timm library [Wightman, 2019]. PT (pretraining) refers to the first stage of model training. FT (fine-tuning) refers to subsequent stages of model training. ZS (zero-shot) refers to the evaluation process described in Sec. 4.1.2. Parameter count is rounded to the nearest million; label set size to the nearest hundred.

148 training, and that contrastive models require positive as well as large batch negative samples to
 149 learn [Shalev-Shwartz & Ben-David, 2014; Radford et al., 2021]. Less well-studied is whether
 150 non-contrastive models also benefit from scaling batch negatives (OOD classes, in this case). We find
 151 that increasing the number of OOD classes from 0 to 900 leads to reliable gains in both accuracy and
 152 robustness.

OOD Data Source	Dataset Size	IN100-Val / Avg. Rob.
OpenImages	1.2 Mn	91.5% / 50.2%
ImageNet	1.3 Mn	90.7% / 46.8%
FractalDB	1 Mn	85.4% / 40.2%

Table 2: **Horizontal scaling with different out-of-distribution approaches.** Horizontal scaling works as well or better when the out-of-distribution classes are not from ImageNet. Horizontal scaling on synthetic out-of-distribution classes underperforms scaling on natural images. Dataset size is rounded to the nearest 100,000.

153 *When we pretrain and fine-tune models, even very large label sets can benefit from horizontal scaling.*
 154 As the label set grows extremely large, our naive approach to horizontal scaling fails; we examine
 155 a range of timm model checkpoints pretrained on the entirety of ImageNet (which contains almost
 156 22,000 classes), even larger and better-performing architectures fail to match much smaller models
 157 fine-tuned on IN1000 (see Tab. 1).

158 There are many potential explanations for this phenomenon. Candidates include the lack of a
 159 validation set against which to optimize, or the extreme class imbalance. Another possible explanation
 160 is that horizontal scaling cannot succeed when the class space is very large. We show that this is not
 161 the case – multi-stage training, even on very large label sets, (see Tab. 1, lines 1, 2) outperforms naive
 162 horizontal scaling on 1000 classes (Tab. 1, line 3).

163 The fact that pretraining and fine-tuning seems to work in *either direction* is particularly surprising
 164 to us. Starting with fewer classes and fine-tuning on many more works almost as well as the more
 165 intuitive approach of starting with a model pretrained on many classes and fine-tuning on a much
 166 smaller number of classes. We leave further exploration of these dynamics to future work.

167 *Horizontal scaling works with non-ImageNet images.* In (Tab. 2), we describe the results of horizontal
 168 scaling on both ImageNet and non-ImageNet out of distribution classes. Specifically, we experiment
 169 with adding classes from OpenImages which do not overlap with ImageNet-1000 classes, and
 170 synthetic classes from the FractalDB dataset [Kataoka et al., 2020]. We find that the OpenImages
 171 classes actually *out-perform* the ImageNet classes substantially, despite the dataset being slightly
 172 smaller. This finding suggests that dataset blending, combined with horizontal scaling, is a promising
 173 approach for training more performant computer vision models.

174 4.2 Class (Im)Balance

175 It has long been observed that class-imbalanced datasets underperform compared to balanced
 176 ones [Goodfellow et al., 2016]. In line with those observations, we confirm that class imbalance has a
 177 substantial negative effect on both validation accuracy and average robustness for IN100 classification.
 178 Less frequently explored is the question of *why* imbalanced datasets underperform.

179 Inspired by recent work investigating Zipfian distributions Chan et al. [2022], we posit two hypotheses.
 180 Zipf’s law stipulates that the frequency of an event is inversely proportional to its rank in a frequency
 181 table. Zipfian distributions occur in natural language, where a small number of words (like “the” and
 182 “and”) occur very frequently, while the majority of words occur rarely. We interrogate each of these
 183 hypotheses: first, that the underperformance is due to the existence of a few overrepresented classes,
 184 a property which we call **left-skewedness**; second, that the underperformance is due to the existence
 185 of long-tail classes with very few samples in them, which we call **long-tailedness**.

186 4.2.1 Key indicators

187 We propose two potential key indicators to better compare these hypotheses. Our proposed indicator
 188 for left-skewedness is the percentage of samples in the dataset which are members of the most

189 common $k\%$ of classes. For the problems described in this paper, we heuristically set $k = 5$. In a
 190 perfectly balanced dataset, then, left-skewedness will be 5%. In OI100 without rebalancing, it is 64%.
 191 Our proposed indicator of long-tailedness is the percentage of classes with fewer than k samples in
 192 them ($\%<k$ Classes). We heuristically explore two choices: $k = 500$ and $k = 100$.

193 4.2.2 Findings

Data Source	Dataset Size	Left-skew	Long-tail @ 500 / Long-tail @ 100	IN100-Val / Avg. Rob.
in100 (100%)	125,000	5%	0% / 0%	85.3% / 40.6%
in100 (62%)	130,000	5%	0% / 0%	82.5% / 43.1%
oi100 (71%)	190,000	45%	0% / 0%	82.2% / 44.3%
oi100 (60%)	101,000	13%	0% / 0%	79.3% / 41.3%
oi100 (88%)	90,000	25%	0% / 0%	76.6% / 38.8%
oi100 (67%)	135,000	31%	9% / 9%	73.9% / 40.7%
oi100 (57%)	105,000	12%	9% / 9%	73.4% / 39.1%
oi100 (100%)	135,000	64%	64% / 9%	67.7% / 37.2%
oi100 (100%)	53,000	18%	67% / 9%	58.2% / 31.1%

Table 3: **Under-represented classes trigger performance declines.** We ablate class imbalance by blending samples from ImageNet and OpenImages, which share a data source. Surprisingly, class imbalance alone does not cause model performance to degrade (Left-skew, the percentage of samples which are in one of the 5 most common classes, is not predictive of performance). Rather, it is the existence of long-tail classes containing very few samples (Long-tail @ k). The 9 longest tailed classes in OI100 ($k = 100$) account for the majority of performance decline. Sample sizes are rounded to the nearest 1,000 samples. Percentiles are rounded to the nearest percentile.

194 Our main results on class imbalance can be found in Tab. 3, where we report data source, dataset size,
 195 our key indicators, validation accuracy, and average accuracy under shift.

196 *Surprisingly, we find that class imbalance alone does not trigger a performance decline.* When we
 197 preserve the imbalance in the largest classes of OpenImages, but rebalance small classes, performance
 198 *improves* slightly (see Tab. 3 line 3) compared to the baseline (Tab. 3 line 8). Furthermore, when
 199 we decrease the degree of imbalance in OpenImages by truncating the largest classes, performance
 200 degrades (see Tab. 3 line 9). We conclude that models improve when training on more samples of a
 201 given class, even when classes are imbalanced.

202 *Rather, the performance declines are attributable the presence of underrepresented classes.* In Tab. 3,
 203 our long-tailedness indicator is in perfect rank-order agreement with IN100 validation accuracy.
 204 Dataset size and left-skewedness exhibit only weak agreement.

205 *Furthermore, the very smallest classes account for most of the decline.* We ablate this by rebalancing
 206 all classes with more than $k = 100$ samples (using this approach, 91 of 100 classes are balanced).
 207 We find that so doing accounts for only 39% of the overall accuracy gains. The majority of accuracy
 208 loss from class imbalance is attributable to the very small classes.

209 5 Conclusions and Future Work

210 In this paper, we outline useful indicators of data quality which are inexpensive to compute and can
 211 be used during pretraining. In Sec. 4.1, we showed that horizontal scaling benefits models trained
 212 from scratch. Based on this, for small datasets, as an alternative to pretraining and fine-tuning, we
 213 suggest co-training on target classes and classes drawn from existing image datasets in order to ensure
 214 that the overall label set size is large. In Sec. 4.2, we showed that underrepresented classes, rather
 215 than class imbalance, harms model performance. In a setting with limited resources, we recommend
 216 rebalancing the classes with the fewest samples first.

217 In future work, we hope to introduce additional important dataset-level factors that influence pretrain-
 218 ing and provide indicators which make use of the image space as well as the caption space.

219 **References**

- 220 Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx,
221 Michael S. Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, Erik Brynjolfsson,
222 Shyamal Buch, Dallas Card, Rodrigo Castellon, Niladri Chatterji, Annie Chen, Kathleen Creel,
223 Jared Quincy Davis, Dora Demszky, Chris Donahue, Moussa Doumbouya, Esin Durmus, Stefano
224 Ermon, John Etchemendy, Kawin Ethayarajh, Li Fei-Fei, Chelsea Finn, Trevor Gale, Lauren
225 Gillespie, Karan Goel, Noah Goodman, Shelby Grossman, Neel Guha, Tatsunori Hashimoto, Peter
226 Henderson, John Hewitt, Daniel E. Ho, Jenny Hong, Kyle Hsu, Jing Huang, Thomas Icard, Saahil
227 Jain, Dan Jurafsky, Pratyusha Kalluri, Siddharth Karamcheti, Geoff Keeling, Fereshte Khani, Omar
228 Khattab, Pang Wei Koh, Mark Krass, Ranjay Krishna, Rohith Kuditipudi, Ananya Kumar, Faisal
229 Ladhak, Mina Lee, Tony Lee, Jure Leskovec, Isabelle Levent, Xiang Lisa Li, Xuechen Li, Tengyu
230 Ma, Ali Malik, Christopher D. Manning, Suvir Mirchandani, Eric Mitchell, Zanele Munyikwa,
231 Suraj Nair, Avaniika Narayan, Deepak Narayanan, Ben Newman, Allen Nie, Juan Carlos Nieves,
232 Hamed Nilforoshan, Julian Nyarko, Giray Ogut, Laurel Orr, Isabel Papadimitriou, Joon Sung
233 Park, Chris Piech, Eva Portelance, Christopher Potts, Aditi Raghunathan, Rob Reich, Hongyu
234 Ren, Frieda Rong, Yusuf Roohani, Camilo Ruiz, Jack Ryan, Christopher Ré, Dorsa Sadigh,
235 Shiori Sagawa, Keshav Santhanam, Andy Shih, Krishnan Srinivasan, Alex Tamkin, Rohan Taori,
236 Armin W. Thomas, Florian Tramèr, Rose E. Wang, William Wang, Bohan Wu, Jiajun Wu, Yuhuai
237 Wu, Sang Michael Xie, Michihiro Yasunaga, Jiaxuan You, Matei Zaharia, Michael Zhang, Tianyi
238 Zhang, Xikun Zhang, Yuhui Zhang, Lucia Zheng, Kaitlyn Zhou, and Percy Liang. On the
239 Opportunities and Risks of Foundation Models, July 2022. URL <http://arxiv.org/abs/2108.07258>. arXiv:2108.07258 [cs].
- 241 Stephanie C. Y. Chan, Adam Santoro, Andrew K. Lampinen, Jane X. Wang, Aaditya Singh, Pierre H.
242 Richemond, Jay McClelland, and Felix Hill. Data Distributional Properties Drive Emergent
243 In-Context Learning in Transformers, November 2022. URL <http://arxiv.org/abs/2205.05055>. arXiv:2205.05055 [cs].
- 245 Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12M: Pushing
246 web-scale image-text pre-training to recognize long-tail visual concepts. In *CVPR*, 2021.
- 247 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hier-
248 archical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*,
249 2009.
- 250 Benjamin Feuer, Ameya Joshi, Minh Pham, and Chinmay Hegde. Distributionally Robust Classifi-
251 cation on a Data Budget. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856.
252 URL <https://openreview.net/forum?id=D5Z2E8CNsD>.
- 253 Samir Yitzhak Gadre, Gabriel Ilharco, Alex Fang, Jonathan Hayase, Georgios Smyrnis, Thao
254 Nguyen, Ryan Marten, Mitchell Wortsman, Dhruva Ghosh, Jieyu Zhang, Eyal Orgad, Rahim
255 Entezari, Giannis Daras, Sarah Pratt, Vivek Ramanujan, Yonatan Bitton, Kalyani Marathe, Stephen
256 Mussmann, Richard Vencu, Mehdi Cherti, Ranjay Krishna, Pang Wei Koh, Olga Saukh, Alexander
257 Ratner, Shuran Song, Hannaneh Hajishirzi, Ali Farhadi, Romain Beaumont, Sewoong Oh, Alex
258 Dimakis, Jenia Jitsev, Yair Carmon, Vaishal Shankar, and Ludwig Schmidt. DataComp: In search
259 of the next generation of multimodal datasets, July 2023. URL <http://arxiv.org/abs/2304.14108>. arXiv:2304.14108 [cs].
- 261 Ian J. Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, Cambridge,
262 MA, USA, 2016. <http://www.deeplearningbook.org>.
- 263 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
264 recognition. *CoRR*, abs/1512.03385, 2015. URL <http://arxiv.org/abs/1512.03385>.
- 265 Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul
266 Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, Dawn Song, Jacob Steinhardt, and Justin Gilmer.
267 The many faces of robustness: A critical analysis of out-of-distribution generalization. *ICCV*,
268 2021a.
- 269 Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adversarial
270 examples. *CVPR*, 2021b.

- 271 Gabriel Ilharco, Mitchell Wortsman, Nicholas Carlini, Rohan Taori, Achal Dave, Vaishaal Shankar,
272 Hongseok Namkoong, John Miller, Hannaneh Hajishirzi, Ali Farhadi, and Ludwig Schmidt.
273 Openclip, 2021. URL <https://doi.org/10.5281/zenodo.5143773>.
- 274 Yiding Jiang, Parth Natekar, Manik Sharma, Sumukh K. Aithal, Dhruva Kashyap, Natarajan Subra-
275 manyam, Carlos Lassance, Daniel M. Roy, Gintare Karolina Dziugaite, Suriya Gunasekar, Isabelle
276 Guyon, Pierre Foret, Scott Yak, Hossein Mobahi, Behnam Neyshabur, and Samy Bengio. Methods
277 and analysis of the first competition in predicting generalization of deep learning. In Hugo Jair
278 Escalante and Katja Hofmann (eds.), *Proceedings of the NeurIPS 2020 Competition and Demon-
279 stration Track*, volume 133 of *Proceedings of Machine Learning Research*, pp. 170–190. PMLR,
280 06–12 Dec 2021. URL <https://proceedings.mlr.press/v133/jiang21a.html>.
- 281 Hirokatsu Kataoka, Kazushige Okayasu, Asato Matsumoto, Eisuke Yamagata, Ryosuke Yamada,
282 Nakamasa Inoue, Akio Nakamura, and Yutaka Satoh. Pre-training without natural images. In
283 *Asian Conference on Computer Vision (ACCV)*, 2020.
- 284 Simon Kornblith, Jonathon Shlens, and Quoc V. Le. Do Better ImageNet Models Transfer Better?
285 In *CVPR*. arXiv, June 2019. doi: 10.48550/arXiv.1805.08974. URL [http://arxiv.org/abs/
286 1805.08974](http://arxiv.org/abs/1805.08974). arXiv:1805.08974 [cs, stat].
- 287 Alina Kuznetsova, Hassan Rom, Neil Gordon Aldrin, Jasper R. R. Uijlings, Ivan Krasin, Jordi
288 Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Mallocci, Alexander Kolesnikov, Tom Duerig,
289 and Vittorio Ferrari. The open images dataset v4. *International Journal of Computer Vision*, 2018.
- 290 Feng Li, Hao Zhang, Peize Sun, Xueyan Zou, Shilong Liu, Jianwei Yang, Chunyuan Li, Lei Zhang,
291 and Jianfeng Gao. Semantic-SAM: Segment and Recognize Anything at Any Granularity, July
292 2023. URL <http://arxiv.org/abs/2307.04767>. arXiv:2307.04767 [cs].
- 293 Dhruv Mahajan, Ross B. Girshick, Vignesh Ramanathan, Kaiming He, Manohar Paluri, Yixuan
294 Li, Ashwin Barambe, and Laurens van der Maaten. Exploring the limits of weakly supervised
295 pretraining. In Vittorio Ferrari, Martial Hebert, Cristian Sminchisescu, and Yair Weiss (eds.),
296 *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14,
297 2018, Proceedings, Part II*, volume 11206 of *Lecture Notes in Computer Science*, pp. 185–201.
298 Springer, 2018. doi: 10.1007/978-3-030-01216-8_12. URL [https://doi.org/10.1007/
978-3-030-01216-8_12](https://doi.org/10.1007/
299 978-3-030-01216-8_12).
- 300 Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and
301 Ren Ng. NeRF: representing scenes as neural radiance fields for view synthesis. *Communications
302 of the ACM*, 65(1):99–106, December 2021. ISSN 0001-0782. doi: 10.1145/3503250. URL
303 <https://dl.acm.org/doi/10.1145/3503250>.
- 304 Parth Natekar and Manik Sharma. Representation based complexity measures for predicting general-
305 ization in deep learning, 2020.
- 306 Thao Nguyen, Gabriel Ilharco, Mitchell Wortsman, Sewoong Oh, and Ludwig Schmidt. Quality not
307 quantity: On the interaction between dataset design and robustness of clip. *NeurIPS*, 2022.
- 308 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
309 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever.
310 Learning transferable visual models from natural language supervision. In *ICML*, 2021.
- 311 Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do imagenet classifiers
312 generalize to imagenet? In *ICML*, 2019.
- 313 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
314 resolution image synthesis with latent diffusion models, 2021.
- 315 Shibani Santurkar, Yann Dubois, Rohan Taori, Percy Liang, and Tatsunori Hashimoto. Is a caption
316 worth a thousand images? a controlled study for representation learning. *ICLR*, 2022.
- 317 Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis,
318 Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of
319 clip-filtered 400 million image-text pairs. *Data Centric AI NeurIPS Workshop*, 2021.

- 320 Shai Shalev-Shwartz and Shai Ben-David. *Understanding Machine Learning: From Theory to*
321 *Algorithms*. Cambridge University Press, USA, 2014. ISBN 1107057132.
- 322 Or Sharir, Barak Peleg, and Yoav Shoham. The Cost of Training NLP Models: A Concise Overview,
323 April 2020. URL <http://arxiv.org/abs/2004.08900>. arXiv:2004.08900 [cs].
- 324 Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned,
325 hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of the 56th*
326 *Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp.
327 2556–2565, Melbourne, Australia, July 2018. Association for Computational Linguistics. doi:
328 10.18653/v1/P18-1238. URL <https://aclanthology.org/P18-1238>.
- 329 Chen Sun, Abhinav Shrivastava, Saurabh Singh, and Abhinav Gupta. Revisiting Unreasonable
330 Effectiveness of Data in Deep Learning Era. In *2017 IEEE International Conference on Computer*
331 *Vision (ICCV)*, pp. 843–852, October 2017. doi: 10.1109/ICCV.2017.97. ISSN: 2380-7504.
- 332 Bart Thomee, David A. Shamma, Gerald Friedland, Benjamin Elizalde, Karl S. Ni, Douglas N.
333 Poland, Damian Borth, and Li-Jia Li. Yfcc100m: the new data in multimedia research. *Commun.*
334 *ACM*, 59:64–73, 2016.
- 335 Haohan Wang, Songwei Ge, Eric P. Xing, and Zachary Chase Lipton. Learning robust global
336 representations by penalizing local predictive power. In *NeurIPS*, 2019.
- 337 Ross Wightman. Pytorch image models. *GitHub repository*, 2019. doi: 10.5281/zenodo.4414861.
- 338 Lu Yuan, Dongdong Chen, Yi-Ling Chen, Noel Codella, Xiyang Dai, Jianfeng Gao, Houdong Hu,
339 Xuedong Huang, Boxin Li, Chunyuan Li, Ce Liu, Mengchen Liu, Zicheng Liu, Yumao Lu, Yu Shi,
340 Lijuan Wang, Jianfeng Wang, Bin Xiao, Zhen Xiao, Jianwei Yang, Michael Zeng, Luowei Zhou,
341 and Pengchuan Zhang. Florence: A New Foundation Model for Computer Vision, November 2021.
342 URL <http://arxiv.org/abs/2111.11432>. arXiv:2111.11432 [cs].