
Reduced, Reused and Recycled: The Life of a Dataset in Machine Learning Research

Bernard Koch

University of California, Los Angeles
bernardkoch@ucla.edu

Emily Denton

Google Research, New York
dentone@google.com

Alex Hanna

Google Research, Mountain View
alexhanna@google.com

Jacob Foster

University of California, Los Angeles
foster@soc.ucla.edu

Abstract

1 Benchmark datasets play a central role in the organization of machine learning
2 research. They coordinate researchers around shared research problems and serve
3 as a measure of progress towards shared goals. Despite the foundational role
4 benchmarking practices play in the field, relatively little attention has been paid
5 to the dynamics of benchmark dataset use and reuse within and across machine
6 learning subcommunities. In this work we dig into these dynamics, by studying
7 how dataset usage patterns differ across different machine learning subcommunities
8 and across time from 2015-2020. We find increasing concentration on fewer and
9 fewer datasets within task communities, significant adoption of datasets from other
10 tasks, and concentration across the field on datasets that have been introduced by
11 researchers situated within a small number of elite institutions. Our results have
12 implications for scientific evaluation, AI ethics, and equity/access within the field.

13 1 Introduction

14 Datasets form the backbone of machine learning research (MLR). They are deeply integrated into
15 work practices of machine learning researchers, operating as resources for training and testing
16 machine learning models. Moreover, datasets serve a central role in the organization of MLR as a
17 scientific field. Benchmark datasets establish stable points of comparison and coordinate scientists
18 around shared research problems. Improved performance on these benchmarks is considered a key
19 signal for collective progress; it is thus also an important form of scientific capital, sought after by
20 individual researchers and used to evaluate and rank their contributions.

21 Datasets also exemplify machine learning tasks, typically through a collection of input and output
22 pairs [1]. By institutionalizing benchmark datasets, task communities implicitly endorse these data
23 as meaningful abstractions of a task or problem domain. The institutionalization of benchmarks
24 influences the behavior of both researchers and end-users [2]. Because advancement on institutional
25 benchmarks is viewed as an indicator of progress, researchers are encouraged to make design choices
26 to maximize performance to gain credibility for their work. Institutionalization also signals to industry
27 adopters that models can be expected to perform in the real world as they do on the benchmark
28 datasets. The close alignment of datasets with “real world” tasks is thus critical not just to accurate
29 measurement of collective scientific progress, but safe, ethical, and effective deployment of models
30 in the wild.

31 Given their central role in the social and scientific organization of MLR, benchmark datasets have
32 also become a central object of critical inquiry in recent years [3]. Dataset audits have revealed
33 concerning biases that have direct implications for algorithmic bias and harms [4, 5, 6, 7]. Problematic
Submitted to the 35th Conference on Neural Information Processing Systems (NeurIPS 2021) Track on Datasets
and Benchmarks. Do not distribute.

34 categorical schemas have been identified in popular image datasets, including poorly formulated
35 categories and the inclusion of derogatory and offensive labels [8, 9]. Growing research into the
36 disciplinary norms of dataset development have revealed concerning practices relating to dataset
37 development and dissemination, such as unstandardized documentation and maintenance practices
38 [10, 11, 12]. There is also growing concern regarding the limitations of existing datasets and standard
39 metrics of evaluation for evaluating model behaviour in real-world settings and evaluating scientific
40 progress in a problem domain [13, 14].

41 Despite the increase in critical attention to benchmark datasets, surprisingly little empirical attention
42 has been paid to patterns of dataset use and reuse across the field as a whole. In this work we dig
43 into these dynamics, by studying how dataset usage patterns differ across different machine learning
44 subcommunities and across time from 2015-2020 in the Papers With Code (PWC) corpus.¹ More
45 specifically, we study machine learning subcommunities that have formed around different machine
46 learning tasks (e.g. *Sentiment Analysis* and *Face Recognition*) and examine: (i) the extent to which
47 research within task communities is concentrated or distributed across different benchmark datasets,
48 and (ii) patterns of dataset creation and movement between different task communities.

49 Overall, we find that the majority of papers within most tasks prefer datasets that were originally
50 created for other tasks over ones created explicitly for their own task, even though most tasks have
51 created more datasets than they have imported. Consistent with this finding, we see increasing
52 concentration on fewer and fewer datasets within task communities. Lastly, we find that these
53 dominant datasets have been introduced by researchers at just a handful of elite institutions.

54 The remainder of this paper is organized as follows. First, we motivate our research questions
55 by underscoring the critical importance of benchmarks in coordinating machine learning research.
56 Second, we describe our analyses on the PWC corpus, a catalog of datasets and their usage jointly
57 curated manually by the machine learning community and algorithmically by Facebook AI Research.
58 We then present our findings and discuss their implications for scientific validity, the ethical usage
59 of MLR, and inequity within the field. We close by offering recommendations for possible reform
60 efforts for the field.

61 **2 The scientific, social, and ethical, importance of benchmark datasets**

62 Following [1], we understand machine learning benchmarks as community resources against which
63 models are evaluated and compared. Benchmarks typically formalize a particular task through a
64 dataset and associated quantitative metric of evaluation. Benchmarking is the dominant paradigm
65 for evaluation in MLR, and the field collectively views upward trends on benchmarks as noisy but
66 meaningful indicators of scientific progress [2, 1, 15]. Over time, MLR has evolved strong norms
67 to facilitate widespread benchmarking including the development of open-access datasets, formal
68 competitions and challenges, and accompanying “black-box” software that allows researchers to test
69 their algorithms on benchmark datasets with minimal effort.

70 The establishment of benchmark datasets as shared resources for evaluation across the MLR com-
71 munity has unique advantages for coordinating scientists around common goals. First, barriers to
72 participation in MLR are reduced since well resourced institutions can shoulder the costs of dataset
73 curation and annotation². Second, by reducing otherwise complex comparisons to a single agreed
74 upon measure, the scientific community can easily align on the value of research contributions and
75 assess whether progress is being made on a particular task. Finally, a complete commitment to
76 benchmarking has allowed MLR to relax reliance on slower institutions for evaluating progress
77 like peer-review or theoretical integration. Together, these advantages have contributed to MLR’s
78 unprecedented transformation into a “rapid discovery science” in the past decade [16].

79 While there are clear advantages to benchmarking as a methodology of comparing algorithms and
80 measuring progress in a problem domain, there are growing concerns regarding benchmarking cultures
81 in MLR which tend to valorize state-of-the-art (SOTA) results on established benchmark datasets over
82 other forms of quantitative or qualitative analysis. The necessity of SOTA results on well established
83 benchmarks for publication acceptance has been identified as a barrier to the development of new
84 ideas [17] and there have been growing calls for more rigorous and comprehensive empirical analysis

¹paperswithcode.com

²However, machine learning model development still remains a resource intensive activity.

85 of models beyond standard top-line metrics, including reporting model size, energy consumption,
86 fairness metrics, and more [18, 19, 20, 21]. The standard benchmarking paradigm also contributes
87 to underspecification challenges in ML pipelines since a given level of performance on a held out
88 benchmark test set doesn't guarantee a model has learned the appropriate causal structure of a problem
89 [14]. In short, while community alignment on benchmarks and metrics can enable rapid algorithmic
90 advancement, hyper focus on singular metrics at the expense of other more comprehensive forms of
91 rigorous evaluation can lead the community astray and risk the development of models that poorly
92 generalize to the real world.

93 The MLR community has begun to reflect on the utility of established benchmarks in the field and their
94 appropriateness for evaluative purposes. For example, the Fashion-MNIST dataset was introduced
95 because MNIST is perceived to be over-utilized and too easy [22], and the utility of ImageNet — one
96 of the most influential ML benchmark in existence — as a meaningful measure of progress has been
97 a focus of critical examination in the past couple years [23, 24]. SOTA chasing concerns are also
98 compounded by the great capacity ML algorithms have to be “right for the wrong reason” [25],
99 enabling SOTA results that rely on “shortcuts” rather than learning the causal structure dictated by
100 the task [13]. [26] suggests the NLP community may have been “led down the garden path” by
101 over-focusing on “beating” benchmark tasks with models that can easily manipulate linguistic form
102 without any real capacity for language understanding. Recent dataset audits have also revealed that
103 established benchmark datasets tend to reflect very narrow — typically white, male, Western — slices
104 of the world [4, 5, 6, 7]. Thus, over-concentration of research on a small number of datasets and
105 metrics can distort perceptions of progress within the field and have serious ethical implications for
106 communities impacted by deployed models. Despite these discussions, little empirical work has
107 considered whether over-concentration of research on a small number of datasets is a systemic issue.
108 This prompts our first research question:

109 **RQ1: How concentrated are machine learning task communities on specific datasets and has**
110 **this changed over time?**

111 There are also growing concerns regarding the gap between benchmark datasets and the problem
112 domains that they are being used to evaluate progress in. For example, [12] found that computer
113 vision datasets tend to be developed in a manner that is decontextualized from a particular task or
114 application area. Supposedly “general purpose” benchmarks are often valued within the field, though
115 the precise bounds of what makes a dataset suitable for general evaluative purposes remains unclear
116 [15]. These observations prompt our second research question:

117 **RQ2: How frequently do machine learning researchers borrow datasets from other tasks in-**
118 **stead of using one created explicitly for that task?**

119 Despite widespread recognition that datasets are critical to the advancement of the field, slow careful
120 dataset development is often undervalued and disincentivized, especially relative to algorithmic
121 contributions [27, 12]. Given the high value the MLR community places on SOTA performance
122 on established benchmarks, researchers are also incentivized to reuse recognizable benchmarks to
123 legitimize their contributions. Moreover, dataset development is time and labor intensive, making
124 large scale dataset development potentially inaccessible to lower-resourced institutions. These
125 observations prompt our final research question:

126 **RQ3: What institutions are responsible for the major ML benchmarks in circulation?**

127 **3 Data**

128 Our primary data source for this work is Papers With Code³ (PWC), an open source repository
129 for machine learning papers, datasets, and evaluation tables created by researchers at Facebook AI
130 Research. PWC is largely community contributed — anyone can add a benchmarking result or a task,
131 provided the benchmarking result is published in a paper as pre-print, in a conference or a journal.
132 Once tasks and datasets are introduced by humans, PWC scrapes ArXiv using keyword searches to
133 find other examples of the task or uses of the dataset.

134 We downloaded the complete PWC dataset on 06/16/2021 (licensed under CC BY-SA 4.0). In this
135 study, we focus primarily on the “Datasets” archive, as well as papers utilizing those datasets. Each

³www.paperswithcode.com

136 dataset in the archive is associated with metadata such as the modality of the dataset (e.g., texts,
 137 images, video, graphs), the date the dataset was introduced, and the paper title that introduced the
 138 dataset (if relevant). At the time we found 4,384 datasets on the site and scraped 60,647 papers that
 139 PWC associates with those datasets using a PWC internal API.

140 In PWC papers, benchmarks, and (by transitivity) datasets are associated with tasks. For this analysis
 141 we were constrained to the 46,668 papers that use a dataset and are labeled with a task (see Figure
 142 6 for a truncated histogram of usage across datasets). These papers collectively use 3,511 datasets.
 143 Studying the transfer of datasets between tasks imposes an additional constraint that we must know
 144 both the task of the paper that introduced the dataset (“the origin task”) and the task of the paper that
 145 used the dataset later in time (“the destination tasks”). For example, ImageNet [28], was introduced
 146 as a benchmark for *Object Recognition* and *Object Localization* (origin tasks), but is now regularly
 147 utilized as a benchmark for *Image Generation* (destination task) among many others.

148 2,583 datasets on PWC were formally introduced in a paper affiliated with a task, utilized by 39,465
 149 unique papers. An additional 640 datasets were introduced in a paper, but not labeled with tasks.
 150 Two authors manually labeled 50 of these dataset papers with tasks (see supplemental spreadsheet
 151 tab “Manually Tasked Datasets” for justifications) allowing us to include another 17,219 utilizing
 152 papers. We do not utilize the remaining 590 datasets and 2790 utilizing papers (14%).

153 PWC includes a taxonomy of tasks and subtasks but the graph is cyclic, making it hard to disen-
 154 tangle dataset transfer between broad tasks and finer-grained tasks (see data supplement tab “Task
 155 Relations”). For each transfer, we annotate both the transfer between the origin and destination,
 156 and the transfer between the origin’s parents and the destination’s parents. This approach allows
 157 us to accurately capture both dynamics between larger tasks (e.g., *Image Classification* and *Image*
 158 *Generation*), and between finer-grained tasks (e.g., *Image-to-Image Translation* and *Image Inpainting*
 159 who are both children of *Image Generation*). Because we found dataset usages to be noisy (i.e., a
 160 paper would be associated with a dataset if the keyword appeared multiple times in the paper), we
 161 restricted each transfer to destination tasks that PWC had already associated with that dataset.

162 **Datasets for Analyses 1 and 2 (RQ1, RQ2):** Because our annotation system double counts transfers
 163 of a single dataset across different levels of organization, we chose to focus exclusively on high-level
 164 transfers between 313 parent tasks. Because the metrics we use in the analyses (particularly Gini and
 165 Creation Ratio) are biased in small samples, we chose to focus only on parent tasks with more than
 166 the median number of 31 papers. This resulted in a final sample of 133 tasks with 47,607 collective
 167 uses and 924 unique datasets (see supplemental spreadsheet tab “List of Tasks”).

168 **Dataset for Analysis 3 (RQ3):** To study the distribution of successful datasets across institutions,
 169 we linked dataset-introducing papers to the Microsoft Academic Graph (MAG) [29]. Affiliation
 170 concentration analyses were performed on the 2,461 datasets with papers that had the last author
 171 affiliation annotated in MAG.

172 4 Methods and Findings

173 4.1 Analysis 1 (RQ1): Concentration in Task Communities on Datasets

174 4.1.1 Methods

175 To measure how concentrated task communities are on certain datasets (RQ1), we calculated the Gini
 176 coefficient across the distribution of observed dataset usages within each task. Gini is a continuous
 177 measure of dispersion in frequency distributions. The metric is frequently used in social science to
 178 study inequality [30]. The Gini score varies between 0 and 1, with 0 indicating that the papers within
 179 a task use all datasets in equal proportions, and 1 indicating that only a single dataset is used across
 180 all dataset-using papers. Gini is calculated as the average absolute difference in the usage of all pairs
 181 of datasets used in the task, divided by the average usage of datasets. Formally, if x_i is the number of
 182 usages of dataset i out of all n datasets used in the task, then the Gini coefficient of dataset usage is,

$$G = \frac{\sum_{i=1}^n \sum_{j=1}^n |x_i - x_j|}{2 \sum_{i=1}^n \sum_{j=1}^n x_j} = \frac{\sum_{i=1}^n \sum_{j=1}^n |x_i - x_j|}{2n \sum_{j=1}^n x_j} = \frac{\sum_{i=1}^n \sum_{j=1}^n |x_i - x_j|}{2n^2 \bar{x}} \quad (1)$$

183 Because Gini can be biased in small samples [31], we use the the sample corrected Gini, $G_s = \frac{n}{n-1}G$,
 184 and excluded tasks (or task-years when disaggregating by time) with fewer than 10 papers.

185 **Regression Model 1:** In addition to descriptive statistics, we built a regression model to assess the
 186 extent to which observed trends in Gini from year-to-year could be attributable to confounds like
 187 task size, task age, or other task-specific traits at that time. Our outcome is G_s in each task year
 188 from 2015-2020 (Figure 7 shows PWC coverage is limited for papers published before 2015). Our
 189 predictors of interest are:

- 190 1. **Year** (since we are interested in trends in concentration over time)
- 191 2. **CV, NLP, Methods** (three dummy variables indicating whether the task belongs to the
 192 Computer Vision, Natural Language Processing, or Methodology categories in PWC).

193 To absorb additional heterogeneity, we also included the following control covariates:

- 194 1. **Task size** in number of dataset-using/introducing papers for that task in that year
- 195 2. **Task age** (because younger tasks may have higher Gini coefficients)
- 196 3. **Random intercepts for each task** (because we have repeated observations over time)

197 Gini is bounded between 0 and 1 so we use beta regression [30], but apply the smoothing transforma-
 198 tion in [32] to deal with the occasional task-year where the Gini is 0. We use a fully restricted model
 199 with the following interactions:

$$Beta(G_s) = \alpha + \beta_1 Year + \beta_2 Task Size + \beta_3 Task Age + \tag{2}$$

$$\beta_4 CV + \beta_5 NLP + \beta_6 Methods + \beta_7 Full size + \tag{3}$$

$$\beta_8 CV * Year + \beta_9 NLP * Year + \beta_{10} Methods * Year +$$

$$\beta_{11} Year * Task Age * Task Size \tag{4}$$

200 This model was favored over all simpler models on AICc.

201 4.1.2 Findings

202 Controlling for task age, task size, and task-specific effects, Model 1 finds significant evidence for
 203 increasing concentration in task communities for the full dataset over time, predicting a marginal
 204 increase in Gini of .065 from 2015-2020 (Figure 1 top green; Table 1). This trend is also visible in
 205 the overall distributions of Gini coefficients over this period (Figure 1 bottom). By 2020, the Gini
 206 coefficient for a task was .60. There are no statistically significant differences between Computer
 207 Vision and Methodology tasks compared to the full sample (Figure 1 top, Figure 5), but Model 1
 208 suggests increases in concentration are attenuated in 2019-2020 for Natural Language Processing.
 209 During these two years, the model predicts NLP concentrations to decrease by .013 while the full
 210 sample increases by .012.

211 4.2 Analysis 2 (RQ2): Changes in Rates of Adoption and Creation of Datasets Over Time

212 4.2.1 Methods

213 We created two proportions to better understand patterns of dataset usage and creation within tasks as
 214 outcomes:

$$\text{Adoption Ratio} = \frac{\# \text{ of Papers Using Datasets from Other Tasks}}{\# \text{ of Papers Using Datasets from Other Tasks} + \# \text{ of Papers Using Datasets from This Task}}$$

$$\text{Creation Ratio} = \frac{\# \text{ of Datasets Created Within This Task}}{\# \text{ of Datasets Created within this Task} + \# \text{ of Datasets Imported from Other Tasks}}$$

215
 216 **Aggregated Descriptive Analyses:** We first computed these proportions for each of the 133 parent
 217 tasks aggregated across all years, and subsetted these by the “Computer Vision,” “Natural Language
 218 Processing,” and “Methodology” categories.

219 **Regression Models 2A & 2B:** Because our outcomes are now ratios of “successful” counts out of
 220 “all” counts, they naturally follow a binomial distribution. We used a mixed effects logistic regression
 221 to model these outcomes with the same predictors as Model 1.

222 **4.2.2 Findings**

223 The top row of Figure 2 shows a wide variance in adoption ratios in both the full sample and the
 224 subcategories. Within the full sample, more than half of task communities use adopted datasets 57.1%
 225 of the time. However, this number varies dramatically across the three PWC subcategories. In more
 226 than half of Computer Vision communities, authors adopt 71.7% of their datasets from a different
 227 task, while half of Natural Language Processing communities adopt datasets less than 28.3% of the
 228 time. Methodology tasks adopt datasets from other tasks at very high rates as well (76.0%).

229 In the bottom row of Figure 2, we see a largely inverted trend. Of all unique datasets used in a task
 230 community, 66.7% are created specifically for that task in more than half of tasks. Within Computer
 231 Vision and Methods tasks, the median is lower at 58.9% and 63.3% with similar distributions across
 232 tasks. Most strikingly, 78% of datasets are created specifically for the task in more than half of NLP
 233 communities with a much tighter variance. We do note that there is a significant correlation between
 234 creation ratio and task size (Spearman’s $\rho = .26 p = 0$).

235 Regression Models 2A and 2B do not find any trends in adoption or creation ratios over time (data
 236 not shown).

237

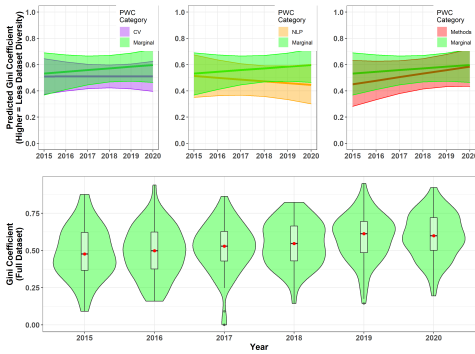


Figure 1: **Top: Predicted concentration on datasets across task communities over time.** Gini predicted by Model 1 holding task size/age to means. Green plots show the estimated effects of the full dataset, other colors are fixed effects for categories. 95% confidence intervals shown. **Bottom: Distributions of concentrations over time.** Higher Gini indicates greater concentration on fewer datasets. We observe significant spread of Gini across different task communities, with the median trending upwards over time.

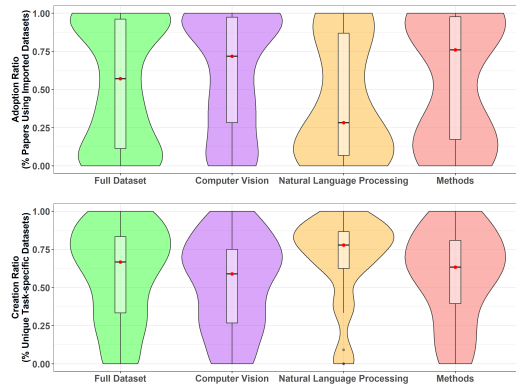


Figure 2: **Adoption (top) and Creation (bottom) Ratios for PWC parent Tasks** . Full dataset in green, tasks in the Computer Vision category in purple, Methods tasks in red, and Natural Language Processing tasks in orange. Red dot and line in boxplot indicate median. Width of violins indicates distribution of tasks.

238 **4.3 Analysis 3 (RQ3): Concentration in Dataset-Introducing Institutions Over Time**

239 **4.3.1 Methods**

240 To look at trends in Gini inequality across institutions and datasets over time for the larger set of
 241 dataset-using papers, we calculated the Gini coefficient G_s in each year for dataset usages by both
 242 dataset and by institution. We regressed this Gini on year, as well as the total number of papers used
 243 to estimate G_s , using a standard beta regression. We also mapped dataset-introducing institutions
 244 using the longitude and latitude coordinates provided for the first author’s institution on Microsoft
 245 Academic.

246 **4.3.2 Findings**

247 Overall, we find that widely-used datasets are introduced by only a handful of elite institutions (Figure
 248 3A). In fact, over 50% of dataset usages in PWC as of June 2021 can be attributed to just thirteen
 249 institutions. Moreover, this concentration on elite institutions as measured through Gini has increased

250 to the mid .80s in recent years (Figure 3B red). This trend is also observed in Gini concentration on
251 datasets in PWC more generally (Figure 3B blue).

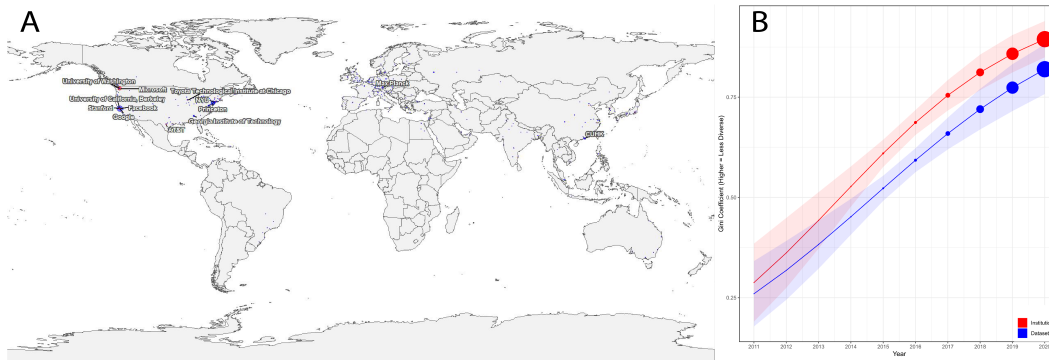


Figure 3: **Increases in concentration of dataset usages on institutions and datasets (non-task specific) over time.** **A:** Map of dataset usages per institution as of June 2021. Dot size indicates number of usages. Dot color indicates whether the institution is for-profit or not-for-profit. Institutions accounting for 50%+ usages labeled. **B:** Gini coefficient for concentration of dataset usages across the whole PWC dataset over time for both institutions and datasets. Ribbons indicate 95% confidence intervals.

252 5 Discussion

253 In this paper, we find that task communities are both heavily concentrated on a limited number
254 of datasets, and that this concentration has been increasing over time (see Figure 1). Moreover, a
255 significant portion of the datasets being used for benchmarking purposes within these communities
256 were originally developed for a different task (see Figure 2). This result is striking given the fact that
257 communities *are* creating new datasets — in most cases more than the unique number that have been
258 imported from other tasks — but the newly introduced datasets are being used at lower rates. When
259 examining PWC as a whole, we find that there is increasing inequality in dataset usage globally, and
260 that more than 50% of all dataset usages in our sample of 28,749 were for datasets introduced by
261 thirteen elite, primarily western, institutions.

262 While striking, there are valid reasons to expect widespread adoption and concentration on key
263 datasets. First, a certain degree of research focus on a particular benchmark is both necessary and
264 healthy to establish the validity and utility of the benchmark — or in some cases contest these
265 properties — and gain community alignment around the benchmark as a meaningful measure of
266 progress. Second, the curation of large-scale datasets is not just costly in terms of resources, but may
267 require unique or privileged data (e.g., anonymized medical records, self-driving car logs) accessible
268 to only a few elite academic and corporate institutions. Nevertheless, the extent of concentration we
269 observe poses questions relating to the scientific rigor and ecological validity of machine learning
270 research and underscores benchmarking as a vehicle for inequality in the field. In the remainder of this
271 section we discuss our findings in relation to these two broad themes and outline recommendations
272 that can be enacted at an individual and institutional level. We close by discussing limitations of this
273 analysis and outlining directions for future work.

274 5.1 Scientific rigor and ecological validity of MLR

275 Heavy concentration of research on a small number of datasets for each task community is a fairly
276 unsurprising result given the value placed on SOTA performance in established benchmark datasets —
277 a valuation incentives individual researchers concentrate efforts on maximizing performance gains on
278 well established benchmarks. However, as discussed in Section 2, over-concentrating research efforts
279 on established benchmark datasets risks distorting measures of progress. Moreover, as the rate of
280 technology transfer has accelerated, benchmarks have been increasingly used by industry practitioners
281 to assess the suitability and robustness of different algorithms for live deployment. This transition
282 has transformed epistemic concerns about overfitting datasets into ethical ones. For example, critical
283 research on face recognition and generation datasets, has repeatedly highlighted the lack of diversity
284 in standard benchmark datasets used to evaluate progress [4], even as the technologies are applied
285 in law enforcement contexts that adversely affect those populations [37]. Figure 4c shows the top

286 datasets in usage within the *Face Recognition* community. Here, we see a significant amount of high
 287 stakes reserch being concentrated on a small number of datasets, many of which contain significant
 288 racial and gender biases [4, 38]. An in depth examination of bias within the top benchmarks datasets
 289 in use within different task communities is outside the scope of this work. However, the systemic
 290 nature of bias concerns in ML datasets compounds the epistemic concerns of highly concentrated
 291 research.

292 Our findings also indicate that datasets regularly transfer between different task communities. On the
 293 most extreme end, the majority of the benchmark datasets in circulation for some task communities
 294 were created for other tasks. For example, Figure 4 plots the dataset usages of *Image Generation*
 295 papers on PWC broken down by dataset name (Figure 4b) and origin task (Figure 4a). We observe
 296 only one of the datasets heavily used in the *Image Generation* community was designed specifically
 297 for this task. The widespread practice of adopting established datasets to train and evaluate models
 298 in new problem domains isn’t inherently a problem. However, this practice does raise potential
 299 concerns regarding the extent to which datasets are appropriately aligned with a given problem space.
 300 Moreover, given the widespread prevalence of systematic biases in the most prominent ML datasets,
 301 adopting existing datasets, rather than investing in careful curation of new datasets, risks further
 302 entrenching existing biases.

303 Our findings relating to creation and adoption rates are quite nuanced, and the extent to which high
 304 adoption rates raise significant concerns to ecological validity are yet to be determined. Furthermore
 305 we believe it is worth distinguishing between at least two forms of dataset adoption that seem to be
 306 conflated in the PWC data. On the one hand, we observe datasets that have been developed for one
 307 task be adopted and *adapted* in some form for a new task through, for example, the addition of new
 308 annotations. On the other hand, we observe some datasets being adopted whole cloth from one task
 309 community to another. Each of these forms of dataset adoption potentially raises unique concerns
 310 regarding the validity of the benchmark in a given context. That said, our results add empirical
 311 support to the growing body of scholarship calling for dataset development and use to be rooted in
 312 context [3, 12], particularly important for application oriented tasks.

313 Our findings also compliment and support the growing calls to include forms of qualitative and
 314 quantitative evaluations beyond top-line benchmark metrics [18, 19, 20, 21]. Given the observed high
 315 concentration of research on a small number of benchmark datasets, we believe diversifying forms of
 316 evaluation is especially important to avoid overfitting to existing datasets and misrepresenting progress
 317 in the field. Reducing the near-singular emphasis on SOTA results on established benchmarks may
 318 also offer more voices the opportunity to shape the culture and trajectory of the field.

319 5.2 Social inequality in MLR

320 The extent of concentration we observe underscores that benchmarking is also a vehicle for inequality
 321 in science. The *prima facie* scientific validity granted by SOTA benchmarking is generically
 322 confounded with the social credibility researchers obtain by showing they can compete on a widely
 323 recognized dataset, even if a more context-specific benchmark might be more technically appropriate.
 324 We posit that this dynamics creates a “Matthew Effect” where successful benchmarks, and the elite

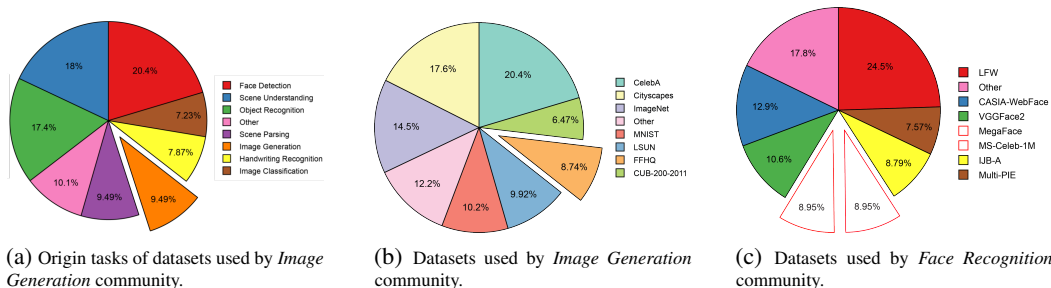


Figure 4: **Top datasets used across *Image Generation* and *Face Recognition* task communities:** (a) Origin task communities of top *Image Generation*. Only 9.49% of *Image Generation* papers in PWC evaluate on datasets developed for *Image Generation*. (b) Names of top *Image Generation*. Only one of the top datasets, FFHQ [33], was developed for the task. (c) The small number of datasets in usage within the high stakes domain of *Face Recognition*. Two of the datasets, MegaFace [34] and MS Celeb-1M [35], have been recently retracted, the latter due to serious ethical violations [36].

325 institutions that introduce them, gain outside stature within the field. To the extent that benchmarks
326 shape the types of questions that get asked and algorithms that get produced, current benchmarking
327 practices thus offer a mechanism through which a small number of elite institutions, both academic
328 and for-profit, can shape the agenda of the field. Moreover, because research trends influence broader
329 public discourse, opinions, and potentially even policy decisions, this influence extends into the
330 broader social world as well.

331 The recently introduced NeurIPS Dataset and Benchmark Track is a clear example of an intervention
332 that shifts incentive structures within the MLR community by rewarding dataset development and
333 other forms of data work. We believe these sorts of interventions can play a critical role in incentiviz-
334 ing careful dataset development that is meaningfully aligned with problem domains. However, our
335 finding that a small number of well-resourced institutions are responsible for most benchmarks in
336 circulation today has implications for data oriented interventions in the field. Our research suggests
337 that simply calling for ML researchers to develop more datasets, and shifting incentive structures
338 so that dataset development is valued and rewarded may not be enough to diversify dataset usage
339 and diversify the perspectives that are ultimately shaping and setting MLR research agendas. In
340 addition to incentivizing dataset development, we advocate for equity oriented policy interventions
341 that prioritize significant funding for people in less resourced institutions to create high-quality
342 datasets. This would diversify — from a social and cultural perspective — the benchmark datasets in
343 rotation.

344 **5.3 Limitations and Future Work**

345 In this paper, we provide the first field-scale analysis on dataset usage in MLR. Because our findings
346 rely on a unique community-curated resource, our findings are contingent on the structure and
347 coverage of PWC. The crowdsourced taxonomy of parent-child task relations in PWC is both noisy
348 and open to interpretation. We have included the full list of parent tasks used in our analysis in the
349 supplementary material, as well as the parent/child relations. We focused our adoption and creation
350 rate analyses on parent-to-parent transfers in an effort to curtail any concerns regarding arbitrariness
351 of task boundaries for fine grained tasks.

352 As with any dataset, PWC also reflects various forms of curatorial bias. To control for spurious labels
353 of dataset usage, we conservatively only considered usages of a dataset valid if they shared a task
354 label with the dataset. Our own curatorial decision influenced the final dataset as well. As noted
355 in Section 3, there were also a large number of datasets that were not assigned an origin task. We
356 manually assigned tasks to the top datasets (assignments and justifications for assignment included in
357 Supplementary Material) from this set and dropped the remaining 14% of uses. Lastly, PWC is likely
358 to reflect recency bias.


359 Finally, we emphasize that our findings are highly nuanced. We report trends that our analysis
360 revealed, but refrain from imposing normative judgements on many of these trends. For example, the
361 high rates of adoption raise potential concerns and points to an important future area of examination.
362 The mere fact that datasets travel between task communities is not necessarily problematic, and
363 indeed the widespread sharing of datasets has been central to methodological advancements in the
364 field. We hope this work will offer a foundation for future empirical work examining the details of
365 dataset transfer and the context specific implications of our findings.

366 **6 Conclusion**

367 Benchmark datasets play a powerful role in the social organization of the field of machine learning. In
368 this work, we empirically examine patterns of creation, adoption, and usage within and across MLR
369 task communities. We find that benchmarking practices are heavily concentrated on a small number
370 of datasets for each task community and heavily concentrated on datasets originating from a small
371 number of well resourced institutions across the field as a whole. We also find that many benchmark
372 datasets flow between multiple task communities and are leveraged to evaluate progress on tasks
373 for which the data was not explicitly designed. We hope this analysis will inform community-wide
374 initiatives to shift patterns of dataset development and use so as to enable more rigorous, ethical, and
375 socially informed research.

376 **References**

- 377 [1] David Schlangen. Targeting the benchmark: On methodology in current natural language
378 processing research. *ArXiv*, abs/2007.04792, 2020.
- 379 [2] Ravit Dotan and Smitha Milli. Value-laden disciplinary shifts in machine learning. In *Proceed-*
380 *ings of the 2020 Conference on Fairness, Accountability, and Transparency*, FAT* '20, page
381 294, New York, NY, USA, 2020. Association for Computing Machinery.
- 382 [3] Amandalynne Paullada, Inioluwa Deborah Raji, Emily M. Bender, Emily Denton, and Alex
383 Hanna. Data and its (dis)contents: A survey of dataset development and use in machine learning
384 research. *NeurIPS Workshop on Machine Learning Retrospectives, Surveys, and Meta-analyses*,
385 2020.
- 386 [4] Joy Buolamwini and Timnit Gebru. Gender shades: Intersectional accuracy disparities in
387 commercial gender classification. In Sorelle A. Friedler and Christo Wilson, editors, *Proceedings*
388 *of Machine Learning Research*, volume 81, pages 77–91, New York, NY, USA, 23–24 Feb 2018.
389 PMLR.
- 390 [5] S. Shankar, Yoni Halpern, Eric Breck, J. Atwood, Jimbo Wilson, and D. Sculley. No classifica-
391 tion without representation: Assessing geodiversity issues in open data sets for the developing
392 world. *arXiv: Machine Learning*, 2017.
- 393 [6] Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. Gender bias
394 in coreference resolution: Evaluation and debiasing methods. In *Proceedings of the 2018*
395 *Conference of the North American Chapter of the Association for Computational Linguistics:*
396 *Human Language Technologies, Volume 2 (Short Papers)*, pages 15–20, New Orleans, Louisiana,
397 June 2018. Association for Computational Linguistics.
- 398 [7] Lucas Dixon, John Li, Jeffrey Sorensen, Nithum Thain, and Lucy Vasserman. Measuring
399 and mitigating unintended bias in text classification. In *Proceedings of the 2018 AAAI/ACM*
400 *Conference on AI, Ethics, and Society*, pages 67–73, 2018.
- 401 [8] Kate Crawford and Trevor Paglen. *Excavating AI: The Politics of Images in Machine Learning*
402 *Training Sets*, 2019.
- 403 [9] Abeba Birhane and Vinay Uday Prabhu. Large image datasets: A pyrrhic win for computer
404 vision? In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer*
405 *Vision*, pages 1537–1547, 2021.
- 406 [10] Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna
407 Wallach, Hal Daumé III, and Kate Crawford. Datasheets for datasets. *arXiv preprint*
408 *arXiv:1803.09010*, 2018.
- 409 [11] R. Stuart Geiger, Kevin Yu, Yanlai Yang, Mindy Dai, Jie Qiu, Rebekah Tang, and Jenny Huang.
410 Garbage in, garbage out? do machine learning application papers in social computing report
411 where human-labeled training data comes from? In *Proceedings of the 2020 Conference on*
412 *Fairness, Accountability, and Transparency*, FAT* '20, page 325–336, New York, NY, USA,
413 2020. Association for Computing Machinery.
- 414 [12] Morgan Klaus Scheuerman, Emily Denton, and Alex Hanna. Do datasets have politics?
415 disciplinary values in computer vision dataset development. *Computer Supported Cooperative*
416 *Work (CSCW)*, 2021.
- 417 [13] Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel,
418 Matthias Bethge, and Felix A Wichmann. Shortcut learning in deep neural networks. *arXiv*
419 *preprint arXiv:2004.07780*, 2020.
- 420 [14] Alexander D’Amour, Katherine A. Heller, Dan Moldovan, Ben Adlam, Babak Alipanahi, Alex
421 Beutel, Christina Chen, Jonathan Deaton, Jacob Eisenstein, Matthew D. Hoffman, Farhad
422 Hormozdiari, Neil Houlsby, Shaobo Hou, Ghassen Jerfel, Alan Karthikesalingam, Mario Lucic,
423 Yi-An Ma, Cory Y. McLean, Diana Mincu, Akinori Mitani, Andrea Montanari, Zachary Nado,
424 Vivek Natarajan, Christopher Nielson, Thomas F. Osborne, Rajiv Raman, Kim Ramasamy, Rory

- 425 Sayres, Jessica Schrouff, Martin Seneviratne, Shannon Sequeira, Harini Suresh, Victor Veitch,
426 Max Vladymyrov, Xuezhi Wang, Kellie Webster, Steve Yadlowsky, Taedong Yun, Xiaohua
427 Zhai, and D. Sculley. Underspecification presents challenges for credibility in modern machine
428 learning. *CoRR*, abs/2011.03395, 2020.
- 429 [15] Amandalynne Paullada Emily Denton Alex Hanna Deborah I Raji, Emily M. Bender. Ai and
430 the everything in the whole wide world benchmark. *NeurIPS Workshop on Machine Learning*
431 *Retrospectives, Surveys, and Meta-analyses*, 2020.
- 432 [16] Randall Collins. Why the social sciences won't become high-consensus, rapid-discovery science.
433 In *Sociological forum*, volume 9, pages 155–177. Springer, 1994.
- 434 [17] Tom Simonite. *Google's AI Guru Wants Computers to Think More Like Brains*, 2018.
- 435 [18] D. Sculley, Jasper Snoek, Alexander B. Wiltschko, and A. Rahimi. Winner's curse? on pace,
436 progress, and empirical rigor. In *ICLR*, 2018.
- 437 [19] Roy Schwartz, Jesse Dodge, Noah A. Smith, and Oren Etzioni. Green AI. *CoRR*,
438 abs/1907.10597, 2019.
- 439 [20] Jesse Dodge, Suchin Gururangan, Dallas Card, Roy Schwartz, and Noah A. Smith. Show your
440 work: Improved reporting of experimental results. In *Proceedings of the 2019 Conference on*
441 *Empirical Methods in Natural Language Processing and the 9th International Joint Conference*
442 *on Natural Language Processing (EMNLP-IJCNLP)*, pages 2185–2194, Hong Kong, China,
443 November 2019. Association for Computational Linguistics.
- 444 [21] Kawin Ethayarajh and Dan Jurafsky. Utility is in the eye of the user: A critique of nlp
445 leaderboards. In *arXiv:2009.13888*, 2020.
- 446 [22] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for
447 benchmarking machine learning algorithms. *CoRR*, abs/1708.07747, 2017.
- 448 [23] L. Beyer, Olivier J. H'enauff, Alexander Kolesnikov, Xiaohua Zhai, and Aäron van den Oord.
449 Are we done with imagenet? *ArXiv*, abs/2006.07159, 2020.
- 450 [24] Dimitris Tsipras, Shibani Santurkar, Logan Engstrom, Andrew Ilyas, and Aleksander Madry.
451 From imagenet to image classification: Contextualizing progress on benchmarks. *International*
452 *Conference on Machine Learning (ICML)*, 2020.
- 453 [25] Benjamin Heinzerling. NLP's Clever Hans moment has arrived. *The Gradient*, 2019.
- 454 [26] Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. On
455 the dangers of stochastic parrots: Can language models be too big? . In *Proceedings of FAccT*
456 *2021*, 2021.
- 457 [27] Nithya Sambasivan, Shivani Kapania, Hannah Highfill, Diana Akrong, Praveen Kumar Paritosh,
458 and Lora Mois Aroyo. "everyone wants to do the model work, not the data work": Data cascades
459 in high-stakes ai. 2021.
- 460 [28] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale
461 Hierarchical Image Database. In *CVPR*, 2009.
- 462 [29] Kuansan Wang, Zhihong Shen, Chiyuan Huang, Chieh-Han Wu, Yuxiao Dong, and Anshul
463 Kanakia. Microsoft academic graph: When experts are not enough. *Quantitative Science*
464 *Studies*, 1(1):396–413, 2020.
- 465 [30] James B. McDonald and Michael Ransom. *The Generalized Beta Distribution as a Model*
466 *for the Distribution of Income: Estimation of Related Measures of Inequality*, pages 147–166.
467 Springer New York, New York, NY, 2008.
- 468 [31] George Deltas. The small-sample bias of the gini coefficient: results and implications for
469 empirical research. *Review of economics and statistics*, 85(1):226–234, 2003.
- 470 [32] Michael Smithson and Jay Verkuilen. A better lemon squeezer? maximum-likelihood regression
471 with beta-distributed dependent variables. *Psychological methods*, 11(1):54, 2006.

- 472 [33] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for genera-
 473 tive adversarial networks. In *2019 IEEE/CVF Conference on Computer Vision and Pattern
 474 Recognition (CVPR)*, pages 4396–4405, 2019.
- 475 [34] Ira Kemelmacher-Shlizerman, Steven M Seitz, Daniel Miller, and Evan Brossard. The megaface
 476 benchmark: 1 million faces for recognition at scale. In *Proceedings of the IEEE Conference on
 477 Computer Vision and Pattern Recognition*, pages 4873–4882, 2016.
- 478 [35] Yandong Guo, Lei Zhang, Yuxiao Hu, Xiaodong He, and Jianfeng Gao. Ms-celeb-1m: A dataset
 479 and benchmark for large-scale face recognition. volume 9907, pages 87–102, 10 2016.
- 480 [36] Jules. Harvey, Adam. LaPlace. Exposing.ai, 2021.
- 481 [37] Clare Garvie, Alvaro Bedoya, and Jonathan Frankle. The perpetual line-up: Unregulated police
 482 face recognition in america, 2016.
- 483 [38] Mei Wang, Weihong Deng, Jiani Hu, Xunqiang Tao, and Yaohai Huang. Racial faces in the
 484 wild: Reducing racial bias by information maximization adaptation network. *2019 IEEE/CVF
 485 International Conference on Computer Vision (ICCV)*, pages 692–702, 2019.

486 Checklist

- 487 1. For all authors...
- 488 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
 489 contributions and scope? [Yes]
- 490 (b) Did you describe the limitations of your work? [Yes] We describe limitations in
 491 Section 5.3.
- 492 (c) Did you discuss any potential negative societal impacts of your work? [Yes] As
 493 discussed in Section 5.3, our findings are highly nuanced. One potential negative
 494 societal impact of this work would be if the nuances of our analysis are lost. We have
 495 taken great care to articulate the extent to which our analysis is contingent on the
 496 particularities of the PWC repository and importance of future work examining in more
 497 depth the implications of this work for research validity in different task communities.
- 498 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
 499 them? [Yes]
- 500 2. If you are including theoretical results...
- 501 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 502 (b) Did you include complete proofs of all theoretical results? [N/A]
- 503 3. If you ran experiments...
- 504 (a) Did you include the code, data, and instructions needed to reproduce the main ex-
 505 perimental results (either in the supplemental material or as a URL)? [No] We have
 506 included our data in the supplementary material. Code will be released upon paper
 507 acceptance.
- 508 (b) Did you specify all the training details [Yes] Parameters of regression model described
 509 in Section 4.1.1.
- 510 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
 511 ments multiple times)? [Yes]
- 512 (d) Did you include the total amount of compute and the type of resources used (e.g., type
 513 of GPUs, internal cluster, or cloud provider)? [N/A]
- 514 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 515 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 516 (b) Did you mention the license of the assets? [Yes]
- 517 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
 518 The PWC data is open source. We have included additional coding we added in the
 519 supplement.
- 520 (d) Did you discuss whether and how consent was obtained from people whose data you’re
 521 using/curating? [Yes]
- 522 (e) Did you discuss whether the data you are using/curating contains personally identifiable
 523 information or offensive content? [N/A]
- 524 5. If you used crowdsourcing or conducted research with human subjects...

- 525 (a) Did you include the full text of instructions given to participants and screenshots, if
 526 applicable? [N/A]
 527 (b) Did you describe any potential participant risks, with links to Institutional Review
 528 Board (IRB) approvals, if applicable? [N/A]
 529 (c) Did you include the estimated hourly wage paid to participants and the total amount
 530 spent on participant compensation? [N/A]

531 **A Appendix**

Table 1: "Exponentiated coefficients for fixed effects in Regression Model 1"

Term	Estimate	Std Error	Statistic	P-value
1 (Intercept)	1.21	1.22	0.96	0.34
2 Year	1.10	1.04	2.37	0.02
3 Task Size	2.73	1.16	6.63	0.00
4 Task Age	1.03	1.11	0.26	0.79
5 CV	0.96	1.21	-0.19	0.85
6 NLP	1.03	1.22	0.16	0.87
7 Methodology	0.68	1.21	-2.04	0.04
8 Year:Task Size	0.84	1.03	-5.95	0.00
9 Year:Task Age	0.98	1.02	-1.06	0.29
10 Task Size:Task Age	1.36	1.17	1.96	0.05
11 Year:CV	0.95	1.04	-1.33	0.18
12 Year:NLP	0.90	1.04	-2.57	0.01
13 Year:Methodology	1.06	1.04	1.48	0.14
14 Year:Task Size:Task Age	0.95	1.03	-1.80	0.07
15 SD(Task Random Intercepts)	1.66			

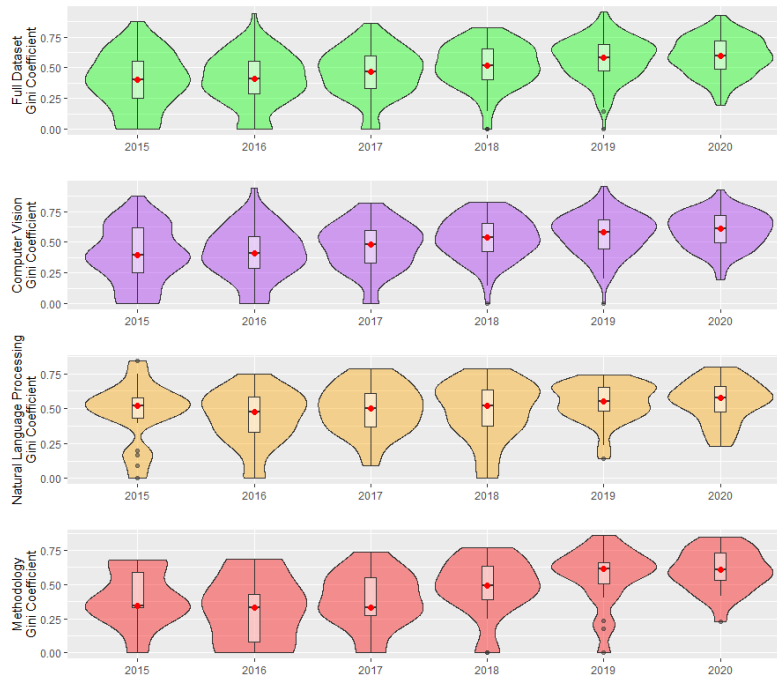


Figure 5: **Increases in concentration on datasets within task communities over time.** Higher Gini coefficient indicates greater concentration on fewer datasets. We observe significant spread of Gini across different task communities, with the median trending upwards over time for all modalities. Green is the full dataset, other colors indicate subsets of the data by PWC task category.

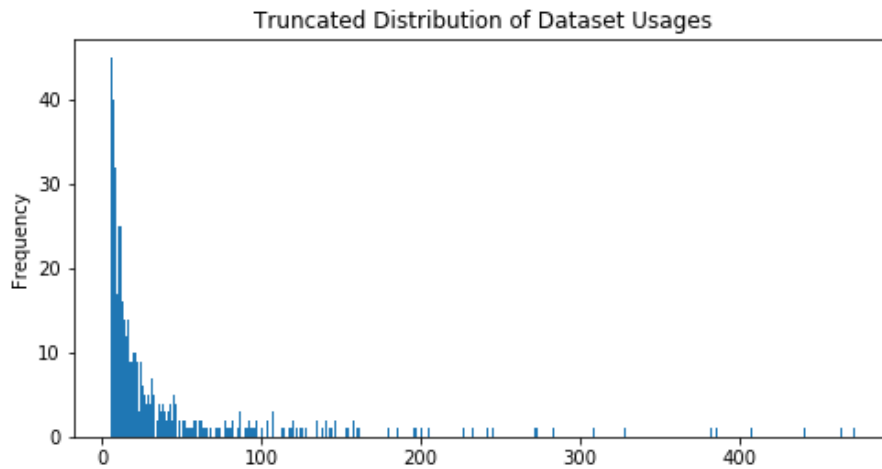


Figure 6: **Truncated distribution of usages per dataset in PWC.** Usages measured conservatively by only allowing usages from tasks the dataset was labeled for. 3760 datasets with less than 5 papers and 8 dataset with over 500 uses dropped for clarity. 8 datasets are Penn Treebank, CelebA, SQuAD, KITTI, MNIST, Cityscapes, ImageNet, COCO.

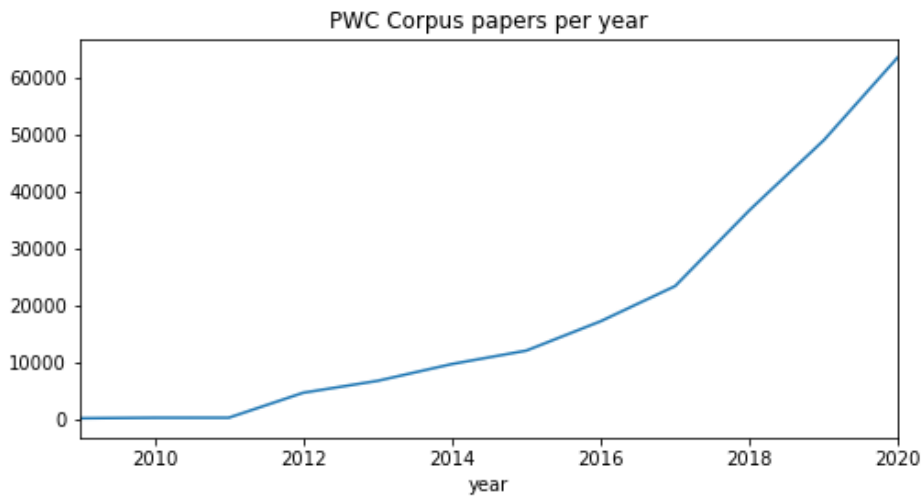


Figure 7: **Number of Papers in the Papers with Code Corpus.** Full set of "Papers with Abstracts" on Papers with Code as of June 2021. Total dataset size is 137,510 papers. Daily snapshots of this dataset are available at github.com/paperswithcode.