
The Benchmark Lottery

Mostafa Dehghani*, Yi Tay*, Alexey A. Gritsenko*, Zhe Zhao, Neil Houlsby,
Fernando Diaz, Donald Metzler†, Oriol Vinyals†
Google Research & DeepMind
{dehghani, yitay, agritsenko}@google.com

Abstract

1 The world of empirical machine learning (ML) strongly relies on benchmarks in
2 order to determine the relative effectiveness of different algorithms and methods.
3 This paper proposes the notion of a *benchmark lottery* that describes the overall
4 fragility of the ML benchmarking process. The benchmark lottery postulates that
5 many factors, other than fundamental algorithmic superiority, may lead to a method
6 being perceived as superior. On multiple benchmark setups that are prevalent in
7 the ML community, we show that the relative performance of algorithms may be
8 altered significantly simply by choosing different benchmark tasks, highlighting the
9 fragility of the current paradigms and potential fallacious interpretation derived from
10 benchmarking ML methods. Given that every benchmark makes a statement about
11 what it perceives to be important, we argue that this might lead to biased progress in
12 the community. We discuss the implications of the observed phenomena and provide
13 recommendations on mitigating them using multiple machine learning domains
14 and communities as use cases, including natural language processing, computer
15 vision, information retrieval, recommender systems, and reinforcement learning.

16 1 Introduction

17 Quantitative evaluation is a cornerstone of machine learning research. As a result, benchmarks,
18 including those based on data sets and simulations, have become fundamental to tracking the progress
19 of machine learning research. Benchmarks have a long history in artificial intelligence research
20 generally. There have been several attempts at designing milestones to capture progress toward
21 artificial intelligence (e.g., human level game performance, the Turing test [Turing, 1950]). Specific
22 system properties are measured through specialized benchmarks (e.g. for vision, natural language
23 processing, robotics). All of these benchmarks, by design, encode values about what is salient
24 and important, both across domains (e.g. natural language processing benchmarks versus robotics
25 benchmarks) and within them (e.g. which languages are considered in an NLP benchmark, which
26 environments are considered in a robotics benchmark).

27 As benchmarks become widely accepted, researchers adopt them, often without questioning their
28 assumptions, and algorithmic development becomes slowly tied to these success metrics. Indeed, over
29 time, the research community makes collective decisions about what shared tasks—and values—are
30 important (through peer review norms and resource investment) and which are not.

31 Because of this, it is important for the research community to understand the individual, community,
32 social, and political pressures that influence why some benchmarks become canonical and others do
33 not. This paper shares some opinions on this topic along with case studies calling for discussion and
34 reconsiderations on several issues with benchmarking in machine learning and argues that a meta-level
35 understanding of benchmarks is a prerequisite for understanding how the progress in machine learning
36 is made. This paper presents analyses on how benchmarks may affect the direction and pace of progress
37 in machine learning and puts forward the notion of a benchmark lottery. We argue that many factors
38 other than the algorithmic superiority of a method may influence the emergence of algorithms that are
39 perceived as better. Moreover, we claim that for a method to emerge successful, it has to first win the

*Equal contribution, †equal advising.

40 *benchmark lottery*. Out of the many potential trials in this lottery, a method has to be first well-aligned
41 with the suite of benchmarks that the community has accepted as canonical. We refer to the alignment
42 between the tasks brought forth by the community and successful algorithms as the *task selection* bias.
43 We empirically show that the task selection process has a great influence over the relative performance
44 of different methods. Moreover, we argue that benchmarks are *stateful*, meaning that the method has to
45 also participate in the lottery at the right moment, and to align well with existing techniques, tricks, and
46 state-of-the-art. Related to this, we also briefly discuss how benchmark reuse may affect the statistical
47 validity of the results of new methods.

48 As a whole, as we researchers continue to participate in the benchmark lottery, there are long-term
49 implications, which we believe are important to be explicitly aware of. As such, the main goals of this
50 paper are to (i) raise awareness of these phenomena and potential issues they create; and to, (ii) provide
51 some recommendations for mitigating these issues. We argue that community forces and task selection
52 biases, if left unchecked, may lead to unwarranted overemphasis of certain types of models and to
53 unfairly hinder the growth of other classes of models - which may be important for making fast and
54 reliable progress in machine learning.

55 The notion of what makes a benchmark canonical, in the sense that is widely accepted by the
56 community, is also diverse depending on the field of study. On one hand, fields like natural language
57 processing (NLP) or computer vision (CV) have well-established benchmarks for certain problems.
58 On the other hand, fields such as recommender systems or reinforcement learning tend to allow
59 researchers more freedom in choosing their own tasks and evaluation criteria for comparing methods.
60 We show how this may act as *rigging the lottery*, where researchers can “make their own luck” by
61 fitting benchmarks and experimental setups to models instead.

62 Overall, this paper explores these aspects of model evaluation in machine learning research. We frame
63 this from a new perspective of the *benchmark lottery*. While there has been recent work that peers
64 deeply into the benchmark tasks themselves [Bowman and Dahl, 2021], this work takes meta- and
65 macro-perspectives to encompass factors that go beyond designing reliable standalone tasks.

66 The remainder of the paper is organized as follow: Section 2 discusses how benchmarks can influence
67 long-term research directions in a given (sub-)field. Section 3 introduces the *task selection bias* and
68 using established benchmarks as examples shows how relative performance of algorithms is affected
69 by the task selection process. Section 4 takes another view of the task selection bias and proposes
70 *community bias* as a higher-level process that influences task selection. We show that forces from
71 the broader research community directly impact the task selection process and as a result, play a
72 substantial role in creating the lottery. Section 5 posits that benchmarks are stateful entities and that
73 participation in a benchmark differs vastly depending upon its state. We also argue continual re-use
74 of the same benchmark may be problematic. Section 6 discusses *rigging the lottery*, the issue that
75 some communities (e.g. recommender systems and reinforcement learning) face, where the lack of
76 well-established community-driven sets of benchmarks or clear guidelines may inadvertently enable
77 researchers to fit benchmarks to model. We highlight the potential drawbacks of such an approach.
78 Finally, in Section 7 we provide recommendations for finding a way out of the lottery by building
79 better benchmarks and rendering more accurate judgments when comparing models.

80 Overall, unified benchmarks have led to incredible progress and breakthroughs in machine learning
81 and artificial intelligence research [Kingma and Welling, 2013, Mikolov et al., 2013, Sutskever et al.,
82 2014, Bahdanau et al., 2014, Goodfellow et al., 2014, Hinton et al., 2015, Silver et al., 2016, He et al.,
83 2016a, Vaswani et al., 2017, Devlin et al., 2018, Brown et al., 2020, Dosovitskiy et al., 2020]. There is
84 certainly a lot of benefits of having the community come together to solve shared tasks and benchmarks.
85 Given that the role of benchmarks is indispensable and highly important for measuring progress, this
86 work seeks to examine, introspect and find ways to improve.

87 2 Background

88 Measuring progress is one of the most difficult aspects of empirical computer science and machine
89 learning. Such questions as “What are the best setup and task to use for evaluation?” [Ponce et al.,
90 2006, Machado et al., 2018, Lin, 2019, Bowman and Dahl, 2021, Recht et al., 2019, Lin et al., 2021,
91 Gulcehre et al., 2020, Perazzi et al., 2016, Vania et al., 2020, Musgrave et al., 2020], “Which data
92 or benchmark are most applicable?” [Metzler and Kurland, 2012, Beyer et al., 2020, Northcutt et al.,
93 2021, Gulcehre et al., 2020, Dacrema et al., 2019], “Which metrics are suitable?” [Machado et al.,
94 2018, Bouthillier et al., 2021, Balduzzi et al., 2018, Bouthillier et al., 2019, Musgrave et al., 2020],
95 or “What are the best practices for fair benchmarking?” [Torralba and Efros, 2011, Armstrong et al.,
96 2009, Machado et al., 2018, Sculley et al., 2018, Lin, 2019, Bowman and Dahl, 2021, Bouthillier et al.,
97 2021, Recht et al., 2019, Lin et al., 2021, Balduzzi et al., 2018, Lipton and Steinhardt, 2018, Bouthillier
98 et al., 2019, Vania et al., 2020, Mishra and Arunkumar, 2021, Marie et al., 2021, Dodge et al., 2019]

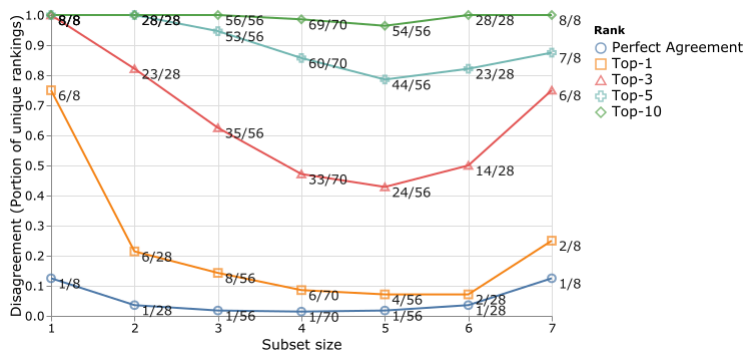


Figure 1: Disagreement of model rankings on the SuperGLUE benchmark as a function of the number of selected benchmark tasks. The x -axis represents the number of tasks in each sub-selection of tasks and each line corresponds to a different value of k for the Top- k in the rankings. Points are labels as A/B , where A is the number of unique model rankings and B is the total number of possible task combinations for this subset size. If $A = 1$, then all rankings are equivalent and consistent across all task selections; higher values of A correspond to higher degrees of disagreement between models rankings.

99 are of utmost importance to correct empirical evaluation of new ideas and algorithms, and have been
 100 extensively studied. Nevertheless, the jury is still out on most of these questions.

101 We argue that some models and algorithms are not inherently superior to their alternatives, but are
 102 instead perceived as such by the research community due to various factors that we discuss in this
 103 paper. One of these factors is the software and hardware support for an idea, as captured in the concept
 104 of hardware lottery by Hooker [2020]. Here however we focus mainly on *benchmarking*-related
 105 factors, and discuss the role they play in the selection of a model as “fashionable” in the research world,
 106 and how this is often conflated with the model being better. When a class of models or algorithms
 107 gets recognition in the community, there will be more follow up research, adaption to more setups,
 108 more tuning and discovery of better configurations, which lead to better results. This is a valid way
 109 of propelling the field further. However, a question that we should also ask is how much progress could
 110 have been made by investing the same amount of time, effort, computational resources and talent in a
 111 different class of models. In other words, assuming model development as a complex high-dimensional
 112 optimization process, in which researchers are exploring a fitness surface, the initial point, as well
 113 as the fitness function, are the key factors for ending up with better optima, and both these factors
 114 are highly affected by the benchmarks used for evaluation.

115 3 Task selection bias

116 As we show in this section, relative model performance is highly sensitive to the choice of tasks and
 117 datasets it is measured on. As a result, the selection of well-established benchmarks plays a more
 118 important role than is perhaps acknowledged, and constitutes a form of partiality and bias - the *task*
 119 *selection bias*.

120 3.1 Case Studies

121 In this section, we study different popular benchmarks and use the data from the leaderboards of these
 122 benchmarks to run analyses that highlight the effect of task selection bias.

123 3.1.1 SuperGLUE

124 In order to study the effect of aggregated scores and how findings change by emphasizing and de-
 125 emphasizing certain tasks, we explore the SuperGLUE dataset [Wang et al., 2019]. To demonstrate the
 126 task selection bias on this benchmark, we re-compute the aggregated scores using different combina-
 127 tions of eight SuperGLUE tasks. We consider over 55 different top performing models that are studied
 128 in [Narang et al., 2021], including transformer-based models with various activation functions, normal-
 129 ization and parameter initialization schemes, and also architectural extensions (e.g., Evolved Trans-
 130 formers [So et al., 2019], Synthesizers [Tay et al., 2020a], Universal Transformer [Dehghani et al., 2019], and
 131 Switch Transformers [Fedus et al., 2021]) as well as convolution-based models (e.g. lightweight and
 132 dynamic convolutions). We consider the fine-grained scores of these models on the 8 individual tasks of
 133 SuperGLUE and their different combinations. For each combination of tasks, we take a mean-aggregate
 134 model performance for all models on the selected tasks and produce a ranking of all 55 models. To
 135 make this ranking more meaningful, we only consider its Top- k entries, where $k \in \{1, 3, 5, 10\}$.

136 **Ranking inconsistency.** Figure 1 gives a concise overview of the number of unique Top- k rankings
 137 produced obtained from fixed-size subsets of tasks. For example among the 70 different possibilities
 138 of selecting 4 out of 8 tasks, there are 6 distinct model ranking orders produced for Top-1 (i.e. there
 139 are 6 different possible top models). Moreover, when considering Top-3 or even Top-5, almost 60
 140 out of 70 rankings do not agree with each other. Overall, the rankings become highly diverse as the
 141 subset of tasks selected from the benchmark is varied. This forms the core of the empirical evidence

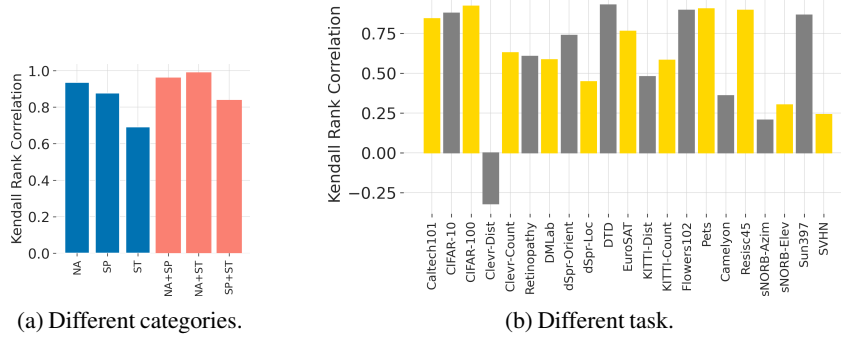


Figure 2: Rank correlation between the full VTAB score and the score for subsets of the benchmark.

142 of the task selection bias. More analyses on ranking of models on all possible combinations of tasks,
 143 rank correlation between SuperGLUE score and individual tasks, effect of relative raking of models
 144 in Appendix A.1,A.2, and A.3.

145 **3.1.2 Visual Task Adaptation Benchmark (VTAB)**

146 A similar situation can be observed for the Visual Task Adaptation Benchmark (VTAB; [Zhai et al.,
 147 2019]) benchmark. VTAB is used for evaluating the quality of representations learned by different
 148 models in terms of their ability to adapt to diverse, unseen tasks with few examples. VTAB defines
 149 a total of 19 tasks, grouped into three categories: *Natural*, *Specialized*, and *Structured*. We have
 150 evaluated 32 different models against all the 19 VTAB tasks. The difference between models is on their
 151 architectures (e.g. WAE-GAN [Tolstikhin et al., 2017] vs. VIVI[Tschannen et al., 2020]), their sizes
 152 (e.g. ResNet-50 vs. ResNet-101 [Kolesnikov et al., 2019]), or the dataset they were pre-trained on (e.g.
 153 ResNet-50 pretrained on ImageNet-21k vs. ResNet-50 pretrained on JFT [Kolesnikov et al., 2019]).
 154 Models we considered in our study are those that are introduced as “representation learning algorithms”
 155 in [Zhai et al., 2019]. More details on the tasks, categories, and models can be found in Appendix A.4.

156 First, we study the agreement of the aggregated score across all 19 tasks with the aggregated scores
 157 obtained from different combinations of the three task categories: natural (NA), specialized (SP), and
 158 structured (ST). Figure 2a shows the Kendall rank correlation, when ranking different models based
 159 on the full VTAB score and based on the category (combination) score. It can be seen that rankings of
 160 models based on different combinations of categories are not always perfectly correlated. For instance,
 161 the structured (ST) subcategory has a correlation of ≈ 0.7 with the full VTAB score, thus highlighting
 162 rather different aspects of the competing models. A more striking point is the full disagreement of
 163 different subcategories on the winning model, i.e. top-1 that is shown in Appendix B, where we future
 164 present the results that show disagreement in the top-1, 2, and 3 rank positions based on different
 165 combinations of sub-categories and tasks. This shows that crowning a model as the winner based on a
 166 single score can be suboptimal, and demonstrates how the random nature of task selection can become
 167 a lottery that algorithms need to win.

168 Figure 2b also presents the correlations between the rankings based on the individual tasks and the
 169 aggregated VTAB score. Unsurprisingly, an even stronger disagreement between rankings is observed
 170 (mean Kendall correlation of ≈ 0.60), including tasks with negative correlation. For more analyses
 171 and additional case studies (Long Range Arena and RL-Unplugged) check Appendix A.

172 **3.2 Score and rank aggregation**

173 So far, we highlighted the issue with reporting a single aggregated score that is supposed to reflect
 174 the performance on multiple tasks, by showcasing the disagreement between different subsets of tasks.
 175 One of the main difficulties for aggregating scores of multiple tasks is the lack of a clear mechanism
 176 for incorporating the difficulty of tasks into account. This is made more complex by the fact that there
 177 are multiple facets to what makes a task difficult. For instance, the size of the training data for different
 178 tasks, the number of prediction classes (and consequently the score for a random baseline for the task),
 179 distribution shift between the pretraining dataset and the downstream tasks, different performance
 180 ranges across tasks, or overrepresenting particular aspects by multiple tasks that introduces biases
 181 into averages [Balduzzi et al., 2018]. As a concrete example, in the case of VTAB some tasks use
 182 the same input data thus upweighting those domains, e.g. CLEVR-Count and CLEVR-Dist use the
 183 same data for different tasks, and for this particular example, given the negative correlation between
 184 CLEVR-Dist and the mean score, this upweighting effect makes the aggregated score even noisier.

185 To address some of these issues, there are alternative ways for ranking models instead of using the mean
 186 score across all tasks as the model performance on the benchmark. For instance, One can grouping
 187 tasks based on their domain) and use macro-averaging to account for the effect upweighting some
 188 domains [Zhai et al., 2019]. Given that using simple averaging for aggregation across multiple tasks,
 189 the maximum score is bounded, this may limit the range of performances, implicitly upweighting tasks
 190 with more headroom. To address this issue, one can use geometric mean instead of arithmetic mean.

191 There are also solutions for rank aggregation that ignore absolute score differences in favor of relative
192 ordering [Dwork et al., 2001, Tabrizi et al., 2015]: For instance, the “average rank” that is obtained by
193 ranking the methods for each task based on their score and then computing the average ranks across tasks.
194 Another alternatives are, for instance, robust average rank, where, before averaging ranks across tasks,
195 the accuracy is binned into buckets of size 1% and all methods in the same bucket get the same rank or
196 elimination ranking (which is equivalent to an exhaustive ballot voting system) [Hao and Ryan, 2016].

197 3.3 Human evaluation bias

198 Related to the task selection bias we discussed in this section, *human evaluation bias* within a task can
199 also play a role in model selection in some tasks like natural language generation. Lack of consistency
200 in how human evaluation, e.g. due to different levels of expertise, cognitive biases, or even inherent
201 ambiguity in the annotation task can introduce a large variability in model comparisons [Schoch et al.,
202 2020]. In the context of measuring the reliability in human annotation, it has been shown that selecting
203 a subset of annotators for evaluation may change the performance of models [Van Der Lee et al., 2019,
204 Amidei et al., 2018, Schoch et al., 2020, Amidei et al., 2020], which can be framed as “annotation
205 bias” that also contribute to the benchmark lottery.

206 4 Community bias

207 Even when viewed as a random process, the task selection bias described in Section 3 alone is sufficient
208 for creating arbitrary selection pressures for machine learning models. We argue however that there is
209 also a higher-level process in which the broader research community influences the task selection, and
210 that counterintuitively leads to the lottery forces not being diminished, but instead more pronounced.
211 This section takes a people perspective of the benchmark lottery and postulates that it is not only the
212 “gamemasters” (benchmark proposers) but also the community that contribute to and reinforce it.

213 While researchers technically have the freedom to select any dataset to showcase their method, this
214 choice is often moderated by the community. A common feedback in the review process of scientific
215 publications that any ML researcher will face eventually is a criticism of the choice of benchmark. For
216 example “*the method was not evaluated on X or Y dataset*” or “*the method’s performance is not SOTA*
217 *on dataset Z*”. Over time, ML researchers tend to gravitate to safe choices of tasks and benchmarks.
218 For example, most papers proposing new pretrained language models [Lan et al., 2019, Liu et al., 2019,
219 Clark et al., 2020, Yang et al., 2020] evaluate on GLUE even if alternatives exist (see example below for
220 further substantiation). In other words, the selection of tasks commonly used in publication is largely
221 driven by the community. Moreover, whether a benchmark is selected as the canonical testbed or not, is
222 not necessarily governed by the quality of the test examples, metrics, evaluation paradigm, or even what
223 the benchmark truly measures. In fact, an argument that the community is solely responsible for the task
224 selection bias is not without merit, since the community is the final endorser and enforcer of these circum-
225 stances. There can be no task selection bias if there is no one to act upon it. To this end, the community
226 might ‘*double down*’ on a benchmark where it becomes almost an unspoken rule for one to evaluate
227 on a particular benchmark. Once a benchmark builds up a following and becomes well-established, it
228 is not hard to imagine that reviewers would ask for results on these benchmarks, potentially regardless
229 of suitability and/or appropriateness. This makes it difficult to fix potentially broken benchmarks.

230 As foreshadowed, commonly used benchmarks are not immune to containing errors. While these errors
231 are likely to be small (as otherwise they would presumably be noticed early on), they do matter in close
232 calls between competing methods. Northcutt et al. [2021] identified label errors in test sets of 10 of
233 the most commonly-used computer vision, natural language, and audio datasets; for example, there are
234 label errors in 6% of the examples in the ImageNet validation set. They showed that correcting label
235 errors in these benchmarks changes model ranking, especially for models that had similar performance.
236 In the field of NLP, it was later found in SNLI [Bowman et al., 2015], which is a dataset for natural
237 language inference (NLI), a large amount of annotation artifacts exists, and it is possible to simply
238 infer the correct label by only using the premise and not the hypothesis [Gururangan et al., 2018]. It
239 is worth noting that SNLI, being the canonical benchmark for NLI, was easily perceived as mandatory
240 for almost any NLI based research.

241 The possibility of having such an issue is not only restricted to the peer review process, but it may extend
242 to the public perception of papers after they are published regardless of whether they went through the
243 peer review process or not. The community bias problem can be raised as the community collectively
244 assigning a weighted impact score for doing well on arbitrarily selected tasks. Achieving state of the
245 art on task Y is then deemed significantly less meaningful than doing that for task X . Moreover, this
246 is not necessarily done without any explicit reasoning as to why one task is preferred to the other, or
247 even how such a “decision” was made. The main concern with respect to the community bias is that
248 research is becoming too incremental and biased toward the common expectations, since a completely
249 new approach will initially have a hard time competing against established and carefully fine-tuned
250 models. For more discussion and concrete examples on the community bias check Appendix C.

251 5 Benchmarks are stateful

252 With leaderboards and the continuous publication of new methods, it is clear that benchmarks are stateful
253 entities. At any point in time, the attempt of a new idea for beating a particular benchmark depends on
254 the information gathered from previous submissions and publications. This is a natural way of making
255 progress on a given problem. But when viewed from the perspective of the selective pressures it causes,
256 it creates another kind of lottery. For many machine learning benchmarks, researchers have full access
257 to the holdout set. Although not explicitly, this typically leads to the violation of the most basic datum
258 of “one should not train on test/holdout set” by getting inspiration from already published works by
259 others who presumably report only the best of the numerous models they evaluated on the test set.

260 Beyond that, it is common to copy-paste hyper-parameters, use the same code, and more recently to
261 even start from pre-retrained checkpoints of previous successful models². In such setups, where the
262 discovery of new models is built on top of thousands of queries, direct or indirect, to the test set, the error
263 rate on test data does not necessarily reflect the true population error [Arora and Zhang, 2021, Blum and
264 Hardt, 2015, Dwork et al., 2015]. The adaptive data analysis framework [Dwork et al., 2015] provides
265 evaluation mechanisms with guaranteed upper bounds on the difference between average error on the
266 test examples and the expected error on the full distribution (population error rates). Based on this
267 framework, if the test set has size N , and the designer of a new model can see the error of the first $i - 1$
268 models on the test set before designing the i -th model, one can ensure the accuracy of the i -th model on
269 the test set is as high as $\Omega(\sqrt{i/N})$ by using the boosting attack [Blum and Hardt, 2015]. In other words,
270 Dwork et al. [2015] state that once we have $i \gg N$ the results on the test set are no longer an indication
271 of model quality. It has been argued that what matters is not only the number of times that a test set
272 has been accessed as stated by adaptive data analysis, but also how it is accessed. Some empirical
273 studies on some popular datasets [Recht et al., 2018, Yadav and Bottou, 2019, Recht et al., 2019]
274 demonstrated that overfitting to holdout data is less of a concern than reasoning from what has been
275 suggested in [Blum and Hardt, 2015]. Roelofs et al. [2019] also studied the holdout reuse by analyzing
276 data from machine learning competitions on the Kaggle and show no significant adaptive overfitting
277 on the classification competitions. Other studies showed that additional factors may prevent adaptive
278 overfitting to happen in practice. For instance, [Feldman et al., 2019b,a] show that in multi-class
279 classification, the large number of classes makes it substantially harder to overfit due to test set reuse. In
280 a recent study, Arora and Zhang [2021] argue that empirical studies that are based on creating or using
281 new test sets (e.g. [Recht et al., 2018, Yadav and Bottou, 2019, Recht et al., 2019]), although reassuring
282 in some level, are not always possible especially in datasets concerning rare or one-time phenomena.
283 They emphasize the need for computing an effective upper bound for the difference between the test
284 and population errors. They propose an upper bound using the description length of models that is
285 based on the knowledge available to model designers before and after the creation of a test set.

286 From the benchmark lottery point of view, the most important aspect of the above phenomena is that the
287 development of new models is shaped by the knowledge of the test errors of all models before it. First of
288 all, there had been events in the past where accessing the test set more than others, intentionally, secured
289 a margin for victory in the race³. In other words, having the ability to access the test set more than others
290 can be interpreted as buying more lottery tickets. Besides, even when there is no explicit intention, the
291 tempting short-term rewards of incremental research polarize people and reinforce the echo chamber
292 effect - leading models are quickly adapted by re-using their code, pre-trained weights, and hyper-
293 parameters are re-used to build something on top of them even faster. Unfortunately, this process makes
294 no time for considering how it affects the statistical validity of results reported on the benchmark.

295 Another aspect of benchmarks being stateful is that participating in shared tasks at a later stage is vastly
296 different from the time of its inception. By then the landscape of research with respect to the specific
297 benchmark is filled with tricks, complicated and specialized strategies, and know-how for obtaining top
298 performance on the task. The adapted recipes for scoring high are not necessarily universal and may be
299 applicable only to a single narrow task or setup. For example, a publication might discover that a niche
300 twist to the loss function produces substantially better results on the task. It is common for all papers
301 subsequently to follow suit. As an example, the community realized that pre-training on MNLI is
302 necessary for obtaining strong performance on RTE and STS datasets [Liu et al., 2019, Clark et al., 2020],
303 and this became common practice later on. Experience shows that it is not uncommon for benchmark

²This is in particular common when a paper provides results based on large scale experiments that are not necessarily feasible to redo for many researchers. For instance, the majority of the papers that propose follow up ideas to Vision Transformer [Dosovitskiy et al., 2020] start by initializing weights from the released pretrained models and follow the setups of the original paper. Similarly, several NLP papers use BERT pretrained models and the same hyper-parameters as BERT in their experimental setup.

³<https://image-net.org/challenges/LSVRC/announcement-June-2-2015>

304 tasks to accumulate lists of best practices and tricks that are dataset- and task-specific⁴. Whether a
305 novel algorithm is able to make use of these tricks (or whether they are available at all) is again a form of
306 lottery, in which models that cannot incorporate *any* of the earlier tricks are significantly disadvantaged.

307 **6 Rigging the lottery: making your own luck**

308 For some tasks and problems, there are already standard benchmarks and established setups that
309 are followed by most of the community. However, for some others, inconsistencies in the employed
310 benchmarks or reported metrics can be observed. This diversity of evaluation paradigms makes
311 comparisons between publications extremely difficult. Alternatively, in some cases, there is simply
312 no standard benchmark or setup, either because the problem is still young, or because there has never
313 been an effort to unify the evaluation. Sometimes this is due to the high computational cost of proper
314 evaluation, like when reporting variance over multiple random seeds is important [Bouthillier et al.,
315 2019]. While in other instances, the root cause is of behavioral nature, where researchers prefer
316 to showcase only what their method shines at - oftentimes to avoid negative reviews, unsuccessful
317 experiments, although performed, are simply not reported. Here, we study two known examples of
318 this issue, which we refer to as *rigging the lottery*.

319 **6.1 Recommender systems and benchmark inconsistencies**

320 Unlike the fields of NLP or CV, there are no well-established evaluation setups for recommender
321 systems [Zhang et al., 2019] that provide canonical ranked lists of model performance. While there
322 has been a famous Netflix prize challenge⁵, this dataset has not been extensively used in academic
323 research or for benchmarking new models. Moreover, even popular datasets like MovieLens [Harper
324 and Konstan, 2015] or Amazon Reviews [He and McAuley, 2016] generally do not have a canonical
325 test split, metric or evaluation method. Therefore, it is still quite unclear about which modern RecSys
326 method one should adopt, as model comparisons are difficult to interpret [Dacrema et al., 2019].

327 Furthermore, RecSys evaluation is also very challenging for a number of reasons. (i) Different
328 recommendation platforms tackle slightly different problems (e.g retrieval [Yi et al., 2019], ranking
329 ([Pei et al., 2019]), or multitask learning ([Zhao et al., 2019])), and each requires their own evaluation
330 setup. (ii) As is common for user interacting systems, user’s reaction towards different algorithms
331 can be different. Constructing offline datasets of user behaviors from an existing system creates an
332 off-policy evaluating challenge [Swaminathan et al., 2016]. (iii) A real-world recommendation system
333 trains on billions of users and items, the scale of user-item interactions makes it extremely difficult
334 to create a complete dataset containing all possible user-item interactions [He et al., 2016b]. As a
335 result, evaluation setups in many recommender system papers tend to be arbitrary.

336 There exists a small number of public datasets (see Appendix E), such as MovieLens [Harper and
337 Konstan, 2015] or Amazon Product Review [He and McAuley, 2016] that are commonly used for
338 evaluating recommender systems. However, even these datasets are tweaked differently in various
339 publications, leading sometimes to contradictory results [Rendle et al., 2020, Zheng et al., 2019].
340 For example, some papers use Hit Ratio and NDCG as evaluation metrics [He et al., 2017], while
341 others resort to using Recall@K [Zheng et al., 2019]. Interestingly, in this particular example, the
342 same methods reverse their performance when a different metric is used. Holdout test sets can also
343 be created differently, with some papers for example using random split [Beutel et al., 2017] and others
344 using an out-of-time split [Zhang et al., 2020].

345 While the majority of this paper discusses cases where a standardized benchmark may lead to biased
346 progress in the ML community, here we instead discuss the *exact opposite* - implications of having
347 no consensus datasets or evaluation setups. Having no unified benchmark for the community to make
348 progress on has numerous flaws. To name a few, (i) this hinders progress in the field, while possibly
349 (ii) creating an illusion of progress. It is not surprising that under these circumstances researchers
350 (potentially unknowingly) tend to find good experimental setups that fit their models. For a case study
351 on inconsistencies of the evaluation setup in ALE benchmark check Appendix D.

352 **7 What can we do?**

353 While the previous sections of the paper focused on the challenges that arise from the lottery-like inter-
354 action between ML benchmarks and the research community, here we would like to show that there are
355 reasons to be optimistic about future developments in this regard. We present suggestions for improving
356 the idea benchmarking process in ways that make it less of a lottery. These recommendations can be also

⁴As an example, for achieving scores that are comparable to top-ranked models on the GLUE benchmark, there are a series of extremely specific actions and setups used in pretraining/finetuning that are known as “standard GLUE tricks” introduced/used by submissions to the leaderboard [Liu et al., 2019, Yang et al., 2019, Lan et al., 2019]. Check the Pre-training and fine-tuning details in the appendix of [Clark et al., 2020].

⁵<https://netflixprize.com/index.html>

357 framed as checklists⁶ for different parts of the process, like making benchmarks, using benchmarks, eval-
358 uation of a new ideas. Appendix F presents a proposed benchmarking checklist for the review process.

359 **7.1 Investing in making guidelines**

360 We believe that the risks of “rigging the lottery” that is described in Section 6 can be minimized by
361 standardizing the recipe for creating and using benchmarks.

362 **Guidelines for creating benchmarks.** Investing into shared guidelines for creating new benchmarks
363 can be extremely beneficial to the long-term health of the research community. In our view, such
364 guidelines should include the current best practices and aspects that require special attention; and should
365 highlight potential concerns for issues that may emerge in the future when different models and algo-
366 rithms are applied to the benchmarks. Fortunately, there have been some efforts in providing guidelines
367 and best practices for making new benchmarks. For example, Zhang [2021] discusses the need for how
368 robotic warehouse picking benchmarks should be designed. Kiela et al. [2021] proposed a framework
369 for benchmarking in NLP that sets clear standards for making new tasks and benchmarks. Denton et al.
370 [2020] look at the dataset construction process with respect to the concerns along the ethical and politi-
371 cal dimensions of what has been taken for granted, and discuss how thinking about data within a dataset
372 must be holistic, future-looking, and aligned with ethical principles and values. Bender and Friedman
373 [2018] also proposed using data statements for NLP datasets in order to provide context that allows
374 users to better understand how experimental results on that dataset might generalize, how software
375 might be appropriately deployed, and what biases might be reflected in systems built on the software.

376 In Section 6 we pointed out that sometimes the blocking factor or conduction rigorous evaluation is
377 the high computational costs, in particular for the academic environment. For instance, analyzing
378 all possible sources of variance in the performance is prohibitively expensive. As a potential solution to
379 this problem, the community can invest more in setting up initiatives like reproducibility challenges⁷
380 or specialized tracks at conferences that offer help in terms of expertise, infrastructure, and computational
381 resources for extensive evaluation to the papers submitted to that conference.

382 **Guidelines for benchmark usage.** Besides the necessity of making guidelines for “how to make new
383 benchmark”, it is important to have clear guidelines for “how to use a benchmark”, which for instance
384 includes the exact setup that the benchmark should be used for evaluation or how the results should
385 be reported. This would be a great help with reducing the instances of rigging the lottery prevalent
386 in some domains (Section 6). There are also several efforts targeting this goal. For instance, Albrecht
387 et al. [2015], Machado et al. [2018] propose specific standards for the ALE benchmark (discussed
388 in Section D.1). Ethayarajh and Jurafsky [2020] argue against ranking models merely based on their
389 performance and propose to always report *model size*, *energy efficiency*, *inference latency*, and metrics
390 indicating model *robustness* and *generalization to the out-of-distribution data*. Gebru et al. [2018]
391 proposed that every dataset be accompanied by a datasheet that documents its motivation, composition,
392 collection process, recommended uses, etc with the goal of increasing transparency and accountability,
393 mitigating unwanted biases in ML systems, facilitating greater reproducibility, and helping researchers
394 and practitioners select more appropriate datasets for their chosen tasks.

395 Another important problem that can benefit from established regulation is the hyper-parameter tuning
396 budget used by researchers to improve their model performance. Spending enough time and compute
397 to precisely tune hyper-parameters of the model or the training process can improve the results a great
398 deal [Li et al., 2018, Bello et al., 2021, Steiner et al., 2021]. Given that, a guideline on limiting the
399 budget for the hyper-parameter tuning can curb the improvements that are solely based on exhausting
400 hyper-parameter search and gives a chance to have comparisons that are tied less to the computational
401 budget of the proposing entity, but more to the merits of the methods themselves.

402 **Guidelines for conferences and reviewers.** There have been attempts to ameliorate the problems
403 related to the benchmark lottery, especially its community biases and the statefulness aspects
404 (Sections 4 and 5). For example, NLP conferences have specially called out “*not being SOTA*” as
405 an invalid basis for paper rejection⁸ We believe it is possible to leverage education through the review
406 process in order to alleviate many negative aspects of benchmark lottery.

407 As an example, we can make sure that in the review process, scores on a particular benchmark are
408 not used for immediate comparison with the top-ranking method on that benchmark, but rather as
409 a sanity check for new models and simply an efficient way of comparing against multiple baselines.
410 This way, fundamentally new approaches will have a chance to develop and mature instead of being
411 forced to compete for top performance right away or get rejected if not succeeded in the early attempts.

⁶Similar to the reproducibility checklist <https://www.cs.mcgill.ca/~jpineau/ReproducibilityChecklist.pdf> [Dodge et al., 2019]

⁷For instance <https://paperswithcode.com/rc2020>, https://reproducibility-challenge.github.io/iclr_2019/, or <https://reproducibility-challenge.github.io/neurips2019/>

⁸<https://2020.emnlp.org/blog/2020-05-17-write-good-reviews>.

412 7.2 Statistical significance testing

413 The presence of established benchmarks and metrics alone does not necessarily lead to a steady
414 improvement of research ideas; it should also be accompanied by rigorous procedures for comparing
415 these ideas on the said benchmarks. For example, Armstrong et al. [2009] discuss the importance
416 of comparing improvements to the strongest available baselines, however, the question of how do
417 we know that if a new model B is *significantly* better than its predecessor model A remains anything
418 but solved 10 years later [Lin et al., 2021].

419 **Benchmark results as random samples.** Machine learning models are usually trained on a training
420 set and evaluated on the corresponding held out test set, where some performance metric m is computed.
421 Because model training is subject to sources of uncontrolled variance, the resulting metric m should
422 be viewed as a single sample from the distribution describing the model’s performance. Because of
423 that, deciding which of the two models is better based on point estimates of their performances m_A and
424 m_B may be unreliable due to chance alone. Instead distributions of these metrics $p(m_A)$ and $p(m_B)$
425 can be compared using statistical significance testing to determine whether the chance that model
426 A is at least as good as model B is low, i.e. $p(A \leq B) < \alpha$ for some *a priori* chosen significance level
427 α . Estimation of $p(A \leq B)$ forms for the crux of statistical significance testing. It can be done either
428 by using parametric tests that make assumptions on distributions $p(m_A)$ and $p(m_B)$ and thus often
429 need fewer samples from these distributions, or by using non-parametric tests that rely on directly
430 estimating the metric distributions and require more samples.

431 The popularity of standardized benchmarks and exponential growth in the amount of research that
432 the ML community has experienced in recent years⁹ exacerbate the risk of inadvertently misguiding
433 research through lax standards on declaring a model as an improvement on the SOTA. Indeed, if point
434 estimates are used in place of statistical significance testing procedures, sampling $m'_A \sim p(m_A)$ and
435 $m'_B \sim p(m_B)$ such that $m'_B > m'_A$ is only a matter of time, even if performance of the two models
436 is not actually different. Note that this is *not* the same as the issue described in Section 5, but could
437 instead be thought of as winning a lottery if you purchase enough lottery tickets.

438 **Beyond a single train-test split.** Unfortunately, researchers rarely go through the process of collecting
439 strong empirical evidence that model B significantly outperforms model A . This is not surprising.
440 As discussed in Bouthillier et al. [2021], obtaining such evidence amounts to running multiple trials
441 of hyper-parameter optimization over sources of variation such as dataset splits, data ordering, data
442 augmentation, stochastic regularisation (e.g. dropout), and random initialization to understand the
443 models’ variance, and is prohibitively expensive¹⁰. If studied at all, mean model performance across
444 several random parameter initializations is used for declaring that the proposed model is a significant
445 improvement. This is vastly sub-optimal because dataset split contributes the most to model variance
446 compared to other sources of variation [Bouthillier et al., 2021]. However, providing multiple dataset
447 splits to estimate this variance is not standard practice in benchmark design.

448 Benchmarks typically come with a single fixed test set, and thus could even be said to unintentionally
449 discourage the use of accurate statistical testing procedures. This is particularly problematic for mature
450 benchmarks, where the magnitude of model improvements may become comparable to the model vari-
451 ance. Systematic variance underestimation may lead to a series of false positives (i.e. incorrectly declar-
452 ing a model to be a significant improvement) that stall research progress, or worse - lead the research com-
453 munity astray by innovating on “improved overfitting” in place of algorithmic improvements. Going
454 forward, one way of addressing this limitation is to design benchmarks with *multiple* fixed dataset splits.
455 As an added benefit, model performance reported across such standardized splits would also enable the
456 application of a variety of statistical tests not only within the same study, but also across publications.

457 **Benchmark design with statistical testing in mind.** The choice of a suitable statistical testing
458 procedure is non-trivial. It must consider the distribution of the metric m that is being compared, the
459 assumption that can be safely made about the distribution (i.e whether a parametric test is applicable or
460 a non-parametric test should be used), the number of statistical tests performed (i.e. whether multiple
461 testing correction is employed) and can also change as the understanding of the metric evolves [Demšar,
462 2006, Bouthillier et al., 2021, Lin et al., 2021]. We, therefore, recommend that benchmark design
463 is accompanied by the recommendation of the suitable statistical testing procedures, including the
464 number dataset splits discussed above, number of replicates experiments, known sources of variance
465 that should be randomized, the statistic to be computed across these experiments and the significance
466 level that should be used for determining statistically significant results. This would not only help
467 the adoption of statistical testing for ML benchmarks, but also serve as a centralized source for best

⁹<https://neuripsconf.medium.com/what-we-learned-from-neurips-2020-reviewing-process-e24549eea38f>

¹⁰Although Bouthillier et al. [2021] also propose a pragmatic alternative to the exhaustive study of all source of variation.

468 practices that are allowed to evolve. A detailed discussion of statistical testing is outside of the scope of
469 this paper, and we refer interested readers to [Bouthillier et al., 2021, Dror et al., 2017] for an overview
470 of statistical testing procedures for ML.

471 **Beyond a single dataset.** Often we are interested in understanding whether model B is significantly
472 better than model A *across a range of tasks*. These kinds of comparisons are facilitated by benchmarks
473 that span multiple datasets (e.g. VTAB or GLUE). Already the question of what it means to do better on a
474 multi-task benchmark is non-trivial due to the task selection bias (see Section 3) - is it sufficient for model
475 B to do better on average; or should it outperform model A on all tasks? It is not surprising that the statisti-
476 cal testing procedures for such benchmarks are also more nuanced - the answer to this question leads to
477 different procedures. It is unclear whether the average metric across datasets, a popular choice for report-
478 ing model performance, is meaningful¹¹ because the errors on different datasets may not be commensu-
479 rable, and because models can have vastly different performance and variances across these datasets. For
480 this reason, more elaborate procedures are required. For example, for the case when we are interested in
481 seeing whether B outperforms A on average Demšar [2006] propose to ignore the variance on individual
482 datasets and treat the model A and B 's performance across datasets as samples from two distributions
483 that should be compared. They recommend that the Wilcoxon signed-rank should be used in such a
484 setup; but the recommended can have limited statistical power when the number of datasets in the bench-
485 mark is small. Alternatively, for cases when we are interested in seeing whether B is better than A on all
486 datasets Dror et al. [2017] propose to perform statistical testing on each of the datasets separately while
487 performing multiple testing corrections. Here again the "right" statistical testing procedure depends on
488 the benchmark, its composition, and the criteria for preferring one model over another; and we believe
489 that the community would benefit if these questions were explicitly answered during benchmark design.

490 7.3 Rise of living benchmarks

491 Another major issue for many popular benchmarks is "creeping overfitting", as algorithms over time
492 become too adapted to the dataset, essentially memorizing all its idiosyncrasies, and losing the ability to
493 generalize. This is essentially related to the statefulness of benchmarks discussed in Section 5. Besides
494 that, measuring progress can be sometimes chasing a moving target since the meaning of progress might
495 change as the research landscape evolves. This problem can be greatly alleviated by for instance chang-
496 ing the dataset that is used for evaluation regularly, as it is done by many annual competitions or reoccur-
497 ing evaluation venues, like WMT¹² or TREC¹³. Besides that, withholding the test set and limiting the
498 number of times a method can query the test set for evaluation on it can also potentially reduce the effect
499 of adaptive overfitting and benchmark reuse. In a more general term, an effective approach is to turn our
500 benchmarks into "living entities". If a benchmark constantly evolves, for instance, adds new examples,
501 adds new tasks, deprecates older data, and fixes labeling mistakes, it is less prone to "tricks" and highly
502 robust models would find themselves consistently doing well across versions of the benchmark. As
503 examples of a benchmark with such a dynamic nature, GEM is a living benchmark for natural language
504 generation [Gehrmann et al., 2021] or Dynabench [Kiela et al., 2021] proposes putting humans and mod-
505 els in the data collection loop where we continuously reevaluate the problem that we really care about.

506 8 Epilogue

507 Ubiquitous access to benchmarks and datasets has been responsible for much of the recent progress in
508 machine learning. We are observing the constant emergence of new benchmarks. And on the one hand,
509 the development of benchmarks is perhaps a sign of continued progress, but on the other hand, there
510 is a danger of getting stuck in a vicious cycle of investing in making static benchmarks that soon will be
511 rejected due to the inflexible flaws in their setup, or lack of generality and possibility for expansion and
512 improvements. We are in the midst of a data revolution and have an opportunity to make faster progress
513 towards the grand goals of artificial intelligence if we understand the pitfalls of the current state of
514 benchmarking in machine learning. The "benchmark lottery" provides just one of the narratives of
515 struggling against benchmark-induced model selection bias. Several topics we touched upon in this
516 paper are discussed in the form of opinions or with a minimum depth as a call for further discussion.
517 We believe each subtopic deserves a dedicated study, like how to better integrate checks for ethical
518 concerns in the mainstream evaluation of every existing benchmark, how to develop tools and libraries
519 that facilitate the rigorous testing of the claimed improvements, or a deep investigation of the social
520 dynamics of the review process and how to improve it. In the end, there are many reasons to be excited
521 about the future - the community is continuously taking positive delta changes that contribute to fixing
522 issues with measuring progress in the empirical machine learning.

¹¹In fact for that reason it was not a popular choice until recently [Demšar, 2006].

¹²<http://statmt.org/>

¹³<https://trec.nist.gov/>

523 **References**

- 524 Stefano V Albrecht, J Christopher Beck, David L Buckeridge, Adi Botea, Cornelia Caragea, Chi-hung
525 Chi, Theodoros Damoulas, Bistra Dilkina, Eric Eaton, Pooyan Fazli, et al. Reports on the 2015
526 aaai workshop program. *Ai Magazine*, 36(2):90–101, 2015.
- 527 Jacopo Amidei, Paul Piwek, and Alistair Willis. Evaluation methodologies in automatic question
528 generation 2013-2018. 2018.
- 529 Jacopo Amidei, Paul Piwek, and Alistair Willis. Identifying annotator bias: A new irt-based method
530 for bias identification. 2020.
- 531 Timothy G Armstrong, Alistair Moffat, William Webber, and Justin Zobel. Improvements that
532 don’t add up: ad-hoc retrieval results since 1998. In *Proceedings of the 18th ACM conference on*
533 *Information and knowledge management*, pages 601–610, 2009.
- 534 Sanjeev Arora and Yi Zhang. Rip van winkle’s razor: A simple estimate of overfit to test data. *arXiv*
535 *preprint arXiv:2102.13189*, 2021.
- 536 Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly
537 learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.
- 538 David Balduzzi, Karl Tuyls, Julien Perolat, and Thore Graepel. Re-evaluating evaluation. *arXiv*
539 *preprint arXiv:1806.02643*, 2018.
- 540 Charles Beattie, Joel Z Leibo, Denis Teplyashin, Tom Ward, Marcus Wainwright, Heinrich Küttler,
541 Andrew Lefrancq, Simon Green, Víctor Valdés, Amir Sadik, et al. Deepmind lab. *arXiv preprint*
542 *arXiv:1612.03801*, 2016.
- 543 Marc G Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The arcade learning environment:
544 An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47:253–279,
545 2013.
- 546 Irwan Bello, William Fedus, Xianzhi Du, Ekin D Cubuk, Aravind Srinivas, Tsung-Yi Lin, Jonathon
547 Shlens, and Barret Zoph. Revisiting resnets: Improved training and scaling strategies. *arXiv*
548 *preprint arXiv:2103.07579*, 2021.
- 549 Emily M Bender and Batya Friedman. Data statements for natural language processing: Toward mit-
550 igating system bias and enabling better science. *Transactions of the Association for Computational*
551 *Linguistics*, 6:587–604, 2018.
- 552 Alex Beutel, Ed H Chi, Zhiyuan Cheng, Hubert Pham, and John Anderson. Beyond globally optimal:
553 Focused learning for improved recommendations. In *Proceedings of the 26th International*
554 *Conference on World Wide Web*, pages 203–212, 2017.
- 555 Lucas Beyer, Olivier J Hénaff, Alexander Kolesnikov, Xiaohua Zhai, and Aäron van den Oord. Are
556 we done with ImageNet? *arXiv preprint arXiv:2006.07159*, 2020.
- 557 Avrim Blum and Moritz Hardt. The ladder: A reliable leaderboard for machine learning competitions.
558 In *International Conference on Machine Learning*, pages 1006–1014. PMLR, 2015.
- 559 Xavier Bouthillier, César Laurent, and Pascal Vincent. Unreproducible research is reproducible. In
560 *International Conference on Machine Learning*, pages 725–734. PMLR, 2019.
- 561 Xavier Bouthillier, Pierre Delaunay, Mirko Bronzi, Assya Trofimov, Brennan Nichyporuk, Justin Szeto,
562 Nazanin Mohammadi Sepahvand, Edward Raff, Kanika Madan, Vikram Voleti, et al. Accounting for
563 variance in machine learning benchmarks. *Proceedings of Machine Learning and Systems*, 3, 2021.
- 564 Samuel R Bowman and George E Dahl. What will it take to fix benchmarking in natural language
565 understanding? *arXiv preprint arXiv:2104.02145*, 2021.
- 566 Samuel R Bowman, Gabor Angeli, Christopher Potts, and Christopher D Manning. A large annotated
567 corpus for learning natural language inference. *arXiv preprint arXiv:1508.05326*, 2015.
- 568 Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal,
569 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
570 few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.

571 Gong Cheng, Junwei Han, and Xiaoqiang Lu. Remote sensing image scene classification: Benchmark
572 and state of the art. *Proceedings of the IEEE*, 105(10):1865–1883, 2017.

573 Mircea Cimpoi, Subhransu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi.
574 Describing textures in the wild. In *Proceedings of the IEEE Conference on Computer Vision and*
575 *Pattern Recognition*, pages 3606–3613, 2014.

576 Kevin Clark, Minh-Thang Luong, Quoc V Le, and Christopher D Manning. Electra: Pre-training
577 text encoders as discriminators rather than generators. *arXiv preprint arXiv:2003.10555*, 2020.

578 Maurizio Ferrari Dacrema, Paolo Cremonesi, and Dietmar Jannach. Are we really making much
579 progress? a worrying analysis of recent neural recommendation approaches. In *Proceedings of*
580 *the 13th ACM Conference on Recommender Systems*, pages 101–109, 2019.

581 Mostafa Dehghani, Stephan Gouws, Oriol Vinyals, Jakob Uszkoreit, and Lukasz Kaiser. Universal
582 transformers. In *Proceedings of the 7th International Conference on Learning Representations,*
583 *ICLR’19*, 2019. URL <https://arxiv.org/abs/1807.03819>.

584 Janez Demšar. Statistical comparisons of classifiers over multiple data sets. *The Journal of Machine*
585 *Learning Research*, 7:1–30, 2006.

586 Emily Denton, Alex Hanna, Razvan Amironesei, Andrew Smart, Hilary Nicole, and Morgan Klaus
587 Scheuerman. Bringing the people back in: Contesting benchmark machine learning datasets. *arXiv*
588 *preprint arXiv:2007.07399*, 2020.

589 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep
590 bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.

591 Jesse Dodge, Suchin Gururangan, Dallas Card, Roy Schwartz, and Noah A Smith. Show your work:
592 Improved reporting of experimental results. *arXiv preprint arXiv:1909.03004*, 2019.

593 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
594 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is
595 worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*,
596 2020.

597 Rotem Dror, Gili Baumer, Marina Bogomolov, and Roi Reichart. Replicability analysis for natural
598 language processing: Testing significance with multiple datasets. *Transactions of the Association*
599 *for Computational Linguistics*, 5:471–486, 2017.

600 Gabriel Dulac-Arnold, Daniel Mankowitz, and Todd Hester. Challenges of real-world reinforcement
601 learning. *arXiv preprint arXiv:1904.12901*, 2019.

602 Cynthia Dwork, Ravi Kumar, Moni Naor, and Dandapani Sivakumar. Rank aggregation methods for the
603 web. In *Proceedings of the 10th international conference on World Wide Web*, pages 613–622, 2001.

604 Cynthia Dwork, Vitaly Feldman, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Aaron Roth. The
605 reusable holdout: Preserving validity in adaptive data analysis. *Science*, 349(6248):636–638, 2015.

606 Kawin Ethayarajh and Dan Jurafsky. Utility is in the eye of the user: A critique of nlp leaderboards.
607 *arXiv preprint arXiv:2009.13888*, 2020.

608 William Fedus, Barret Zoph, and Noam Shazeer. Switch transformers: Scaling to trillion parameter
609 models with simple and efficient sparsity. *arXiv preprint arXiv:2101.03961*, 2021.

610 Li Fei-Fei, Rob Fergus, and Pietro Perona. One-shot learning of object categories. *IEEE transactions*
611 *on pattern analysis and machine intelligence*, 28(4):594–611, 2006.

612 Vitaly Feldman, Roy Frostig, and Moritz Hardt. The advantages of multiple classes for reducing
613 overfitting from test set reuse. In *International Conference on Machine Learning*, pages 1892–1900.
614 PMLR, 2019a.

615 Vitaly Feldman, Roy Frostig, and Moritz Hardt. Open problem: How fast can a multiclass test set
616 be overfit? In *Conference on Learning Theory*, pages 3185–3189. PMLR, 2019b.

617 Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal
618 Daumé III, and Kate Crawford. Datasheets for datasets. *arXiv preprint arXiv:1803.09010*, 2018.

- 619 Sebastian Gehrmann, Tosin Adewumi, Karmanya Aggarwal, Pawan Sasanka Ammanamanchi,
620 Aremu Anuoluwapo, Antoine Bosselut, Khyathi Raghavi Chandu, Miruna Clinciu, Dipanjan Das,
621 Kaustubh D Dhole, et al. The GEM benchmark: Natural language generation, its evaluation and
622 metrics. *arXiv preprint arXiv:2102.01672*, 2021.
- 623 Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The kitti
624 dataset. *The International Journal of Robotics Research*, 32(11):1231–1237, 2013.
- 625 Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,
626 Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information
627 processing systems*, 27, 2014.
- 628 Caglar Gulcehre, Ziyu Wang, Alexander Novikov, Thomas Paine, Sergio Gómez, Konrad Zolna,
629 Rishabh Agarwal, Josh S Merel, Daniel J Mankowitz, Cosmin Paduraru, et al. Rl unplugged: A
630 collection of benchmarks for offline reinforcement learning. *Advances in Neural Information
631 Processing Systems*, 33, 2020.
- 632 Ruiqi Guo, Quan Geng, David Simcha, Felix Chern, Sanjiv Kumar, and Xiang Wu. New
633 loss functions for fast maximum inner product search. *CoRR*, abs/1908.10396, 2019. URL
634 <http://arxiv.org/abs/1908.10396>.
- 635 Suchin Gururangan, Swabha Swayamdipta, Omer Levy, Roy Schwartz, Samuel Bowman, and
636 Noah A. Smith. Annotation artifacts in natural language inference data. In *Proceedings of the
637 2018 Conference of the North American Chapter of the Association for Computational Linguistics:
638 Human Language Technologies, Volume 2 (Short Papers)*, pages 107–112, New Orleans, Louisiana,
639 June 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-2017. URL
640 <https://www.aclweb.org/anthology/N18-2017>.
- 641 Feng Hao and Peter YA Ryan. *Real-world electronic voting: Design, analysis and deployment*. CRC
642 Press, 2016.
- 643 F. Maxwell Harper and Joseph A. Konstan. The movielens datasets: History and context. *ACM
644 Trans. Interact. Intell. Syst.*, 5(4), December 2015. ISSN 2160-6455. doi: 10.1145/2827872. URL
645 <https://doi.org/10.1145/2827872>.
- 646 Matthew Hausknecht, Joel Lehman, Risto Miikkulainen, and Peter Stone. A neuroevolution approach
647 to general atari game playing. *IEEE Transactions on Computational Intelligence and AI in Games*,
648 6(4):355–366, 2014.
- 649 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
650 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
651 pages 770–778, 2016a.
- 652 Ruining He and Julian McAuley. Ups and downs: Modeling the visual evolution of fashion trends
653 with one-class collaborative filtering. In *proceedings of the 25th international conference on world
654 wide web*, pages 507–517, 2016.
- 655 Xiangnan He, Hanwang Zhang, Min-Yen Kan, and Tat-Seng Chua. Fast matrix factorization for online
656 recommendation with implicit feedback. In *Proceedings of the 39th International ACM SIGIR
657 conference on Research and Development in Information Retrieval*, pages 549–558, 2016b.
- 658 Xiangnan He, Lizi Liao, Hanwang Zhang, Liqiang Nie, Xia Hu, and Tat-Seng Chua. Neural
659 collaborative filtering. In *Proceedings of the 26th international conference on world wide web*,
660 pages 173–182, 2017.
- 661 Nicolas Heess, Dhruva TB, Srinivasan Sriram, Jay Lemmon, Josh Merel, Greg Wayne, Yuval Tassa,
662 Tom Erez, Ziyu Wang, SM Eslami, et al. Emergence of locomotion behaviours in rich environments.
663 *arXiv preprint arXiv:1707.02286*, 2017.
- 664 Patrick Helber, Benjamin Bischke, Andreas Dengel, and Damian Borth. Eurosat: A novel dataset
665 and deep learning benchmark for land use and land cover classification. *IEEE Journal of Selected
666 Topics in Applied Earth Observations and Remote Sensing*, 12(7):2217–2226, 2019.
- 667 Peter Henderson, Riashat Islam, Philip Bachman, Joelle Pineau, Doina Precup, and David Meger.
668 Deep reinforcement learning that matters. In *Proceedings of the AAAI Conference on Artificial
669 Intelligence*, volume 32, 2018.

670 Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir
671 Mohamed, and Alexander Lerchner. beta-vae: Learning basic visual concepts with a constrained
672 variational framework. 2016.

673 Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv*
674 *preprint arXiv:1503.02531*, 2015.

675 Sara Hooker. The hardware lottery. *arXiv preprint arXiv:2009.06489*, 2020.

676 Max Jaderberg, Volodymyr Mnih, Wojciech Marian Czarnecki, Tom Schaul, Joel Z Leibo, David
677 Silver, and Koray Kavukcuoglu. Reinforcement learning with unsupervised auxiliary tasks. *arXiv*
678 *preprint arXiv:1611.05397*, 2016.

679 Justin Johnson, Bharath Hariharan, Laurens Van Der Maaten, Li Fei-Fei, C Lawrence Zitnick, and
680 Ross Girshick. Clevr: A diagnostic dataset for compositional language and elementary visual
681 reasoning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*,
682 pages 2901–2910, 2017.

683 Kaggle and EyePacs. Kaggle diabetic retinopathy detection., 2015. URL <https://www.kaggle.com/c/diabetic-retinopathy-detection/data>.

685 Douwe Kiela, Max Bartolo, Yixin Nie, Divyansh Kaushik, Atticus Geiger, Zhengxuan Wu, Bertie
686 Vidgen, Grusha Prasad, Amanpreet Singh, Pratik Ringshia, et al. Dynabench: Rethinking
687 benchmarking in nlp. *arXiv preprint arXiv:2104.14337*, 2021.

688 Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint*
689 *arXiv:1312.6114*, 2013.

690 Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Joan Puigcerver, Jessica Yung, Sylvain Gelly,
691 and Neil Houlsby. Big transfer (bit): General visual representation learning. *arXiv preprint*
692 *arXiv:1912.11370*, 2019.

693 Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.

694 Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu
695 Soricut. Albert: A lite bert for self-supervised learning of language representations. *arXiv preprint*
696 *arXiv:1909.11942*, 2019.

697 Yann LeCun, Fu Jie Huang, and Leon Bottou. Learning methods for generic object recognition with
698 invariance to pose and lighting. In *Proceedings of the 2004 IEEE Computer Society Conference on*
699 *Computer Vision and Pattern Recognition, 2004. CVPR 2004.*, volume 2, pages II–104. IEEE, 2004.

700 Liam Li, Kevin Jamieson, Afshin Rostamizadeh, Ekaterina Gonina, Moritz Hardt, Benjamin Recht,
701 and Ameet Talwalkar. A system for massively parallel hyperparameter tuning. *arXiv preprint*
702 *arXiv:1810.05934*, 2018.

703 Yitao Liang, Marlos C Machado, Erik Talvitie, and Michael Bowling. State of the art control of atari
704 games using shallow reinforcement learning. *arXiv preprint arXiv:1512.01563*, 2015.

705 Jimmy Lin. The neural hype and comparisons against weak baselines. *ACM SIGIR Forum*, 52(2):
706 40–51, 2019.

707 Jimmy Lin, Daniel Campos, Nick Craswell, Bhaskar Mitra, and Emine Yilmaz. Significant
708 improvements over the state of the art? a case study of the ms marco document ranking leaderboard.
709 *arXiv preprint arXiv:2102.12887*, 2021.

710 Nir Lipovetzky, Miquel Ramirez, and Hector Geffner. Classical planning with simulators: Results
711 on the atari video games. In *Proc. IJCAI*, 2015.

712 Zachary C Lipton and Jacob Steinhardt. Troubling trends in machine learning scholarship. *arXiv*
713 *preprint arXiv:1807.03341*, 2018.

714 Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis,
715 Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach.
716 *arXiv preprint arXiv:1907.11692*, 2019.

717 Marlos C Machado, Marc G Bellemare, Erik Talvitie, Joel Veness, Matthew Hausknecht, and Michael
718 Bowling. Revisiting the arcade learning environment: Evaluation protocols and open problems
719 for general agents. *Journal of Artificial Intelligence Research*, 61:523–562, 2018.

- 720 Benjamin Marie, Atsushi Fujita, and Raphael Rubino. Scientific credibility of machine translation
721 research: A meta-evaluation of 769 papers. *arXiv preprint arXiv:2106.15195*, 2021.
- 722 Jarryd Martin, Suraj Narayanan Sasikumar, Tom Everitt, and Marcus Hutter. Count-based exploration
723 in feature space for reinforcement learning. *arXiv preprint arXiv:1706.08090*, 2017.
- 724 Donald Metzler and Oren Kurland. Experimental methods for information retrieval. In *Proceedings*
725 *of the 35th international ACM SIGIR conference on Research and development in information*
726 *retrieval*, pages 1185–1186, 2012.
- 727 Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations
728 of words and phrases and their compositionality. In *Advances in neural information processing*
729 *systems*, pages 3111–3119, 2013.
- 730 Swaroop Mishra and Anjana Arunkumar. How robust are model rankings: A leaderboard customization
731 approach for equitable evaluation. *arXiv preprint arXiv:2106.05532*, 2021.
- 732 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan
733 Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint*
734 *arXiv:1312.5602*, 2013.
- 735 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare,
736 Alex Graves, Martin Riedmiller, Andreas K Fiedel, Georg Ostrovski, et al. Human-level control
737 through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- 738 Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim
739 Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement
740 learning. In *International conference on machine learning*, pages 1928–1937. PMLR, 2016.
- 741 Kevin Musgrave, Serge Belongie, and Ser-Nam Lim. A metric learning reality check. In *European*
742 *Conference on Computer Vision*, pages 681–699. Springer, 2020.
- 743 Arun Nair, Praveen Srinivasan, Sam Blackwell, Cagdas Alcicek, Rory Fearon, Alessandro De Maria,
744 Vedavyas Panneershelvam, Mustafa Suleyman, Charles Beattie, Stig Petersen, et al. Massively
745 parallel methods for deep reinforcement learning. *arXiv preprint arXiv:1507.04296*, 2015.
- 746 Sharan Narang, Hyung Won Chung, Yi Tay, William Fedus, Thibault Fevry, Michael Matena, Karishma
747 Malkan, Noah Fiedel, Noam Shazeer, Zhenzhong Lan, Yanqi Zhou, Wei Li, Nan Ding, Jake Marcus,
748 Adam Roberts, and Colin Raffel. Do transformer modifications transfer across implementations
749 and applications?, 2021.
- 750 Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. Reading
751 digits in natural images with unsupervised feature learning. 2011.
- 752 Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number
753 of classes. In *2008 Sixth Indian Conference on Computer Vision, Graphics & Image Processing*,
754 pages 722–729. IEEE, 2008.
- 755 Curtis G Northcutt, Anish Athalye, and Jonas Mueller. Pervasive label errors in test sets destabilize
756 machine learning benchmarks. *arXiv preprint arXiv:2103.14749*, 2021.
- 757 Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, and CV Jawahar. Cats and dogs. In *2012 IEEE*
758 *conference on computer vision and pattern recognition*, pages 3498–3505. IEEE, 2012.
- 759 Changhua Pei, Yi Zhang, Yongfeng Zhang, Fei Sun, Xiao Lin, Hanxiao Sun, Jian Wu, Peng Jiang,
760 Junfeng Ge, Wenwu Ou, et al. Personalized re-ranking for recommendation. In *Proceedings of*
761 *the 13th ACM Conference on Recommender Systems*, pages 3–11, 2019.
- 762 Federico Perazzi, Jordi Pont-Tuset, Brian McWilliams, Luc Van Gool, Markus Gross, and Alexander
763 Sorkine-Hornung. A benchmark dataset and evaluation methodology for video object segmentation.
764 In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 724–732,
765 2016.
- 766 Jean Ponce, Tamara L Berg, Mark Everingham, David A Forsyth, Martial Hebert, Svetlana Lazebnik,
767 Marcin Marszalek, Cordelia Schmid, Bryan C Russell, Antonio Torralba, et al. Dataset issues in
768 object recognition. In *Toward category-level object recognition*, pages 29–48. Springer, 2006.

- 769 Alexander Pritzel, Benigno Uria, Sriram Srinivasan, Adria Puigdomenech Badia, Oriol Vinyals,
770 Demis Hassabis, Daan Wierstra, and Charles Blundell. Neural episodic control. In *International*
771 *Conference on Machine Learning*, pages 2827–2836. PMLR, 2017.
- 772 Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do cifar-10 classifiers
773 generalize to cifar-10? *arXiv preprint arXiv:1806.00451*, 2018.
- 774 Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do ImageNet classifiers
775 generalize to ImageNet? In *International Conference on Machine Learning*, pages 5389–5400.
776 PMLR, 2019.
- 777 Steffen Rendle, Walid Krichene, Li Zhang, and John Anderson. Neural collaborative filtering vs.
778 matrix factorization revisited. In *Fourteenth ACM Conference on Recommender Systems*, pages
779 240–248, 2020.
- 780 Rebecca Roelofs, Vaishaal Shankar, Benjamin Recht, Sara Fridovich-Keil, Moritz Hardt, John Miller,
781 and Ludwig Schmidt. A meta-analysis of overfitting in machine learning. *Advances in Neural*
782 *Information Processing Systems*, 32:9179–9189, 2019.
- 783 Stephanie Schoch, Diyi Yang, and Yangfeng Ji. “this is a problem, don’t you agree?” framing and
784 bias in human evaluation for natural language generation. In *Proceedings of the 1st Workshop on*
785 *Evaluating NLG Evaluation*, pages 10–16, 2020.
- 786 Julian Schrittwieser, Thomas Hubert, Amol Mandhane, Mohammadamin Barekatain, Ioannis
787 Antonoglou, and David Silver. Online and offline reinforcement learning by planning with a learned
788 model. *arXiv preprint arXiv:2104.06294*, 2021.
- 789 David Sculley, Jasper Snoek, Alex Wiltschko, and Ali Rahimi. Winner’s curse? on pace, progress,
790 and empirical rigor. 2018.
- 791 Minjoon Seo, Tom Kwiatkowski, Ankur Parikh, Ali Farhadi, and Hannaneh Hajishirzi. Phrase-indexed
792 question answering: A new challenge for scalable document comprehension. In *Proceedings*
793 *of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 559–564,
794 Brussels, Belgium, October–November 2018. Association for Computational Linguistics. doi:
795 10.18653/v1/D18-1052. URL <https://www.aclweb.org/anthology/D18-1052>.
- 796 David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche,
797 Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering
798 the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.
- 799 David So, Quoc Le, and Chen Liang. The evolved transformer. In *International Conference on*
800 *Machine Learning*, pages 5877–5886. PMLR, 2019.
- 801 Andreas Steiner, Alexander Kolesnikov, Xiaohua Zhai, Ross Wightman, Jakob Uszkoreit, and Lucas
802 Beyer. How to train your vit? data, augmentation, and regularization in vision transformers. *arXiv*
803 *preprint arXiv:2106.10270*, 2021.
- 804 Ilya Sutskever, Oriol Vinyals, and Quoc V Le. Sequence to sequence learning with neural networks.
805 *arXiv preprint arXiv:1409.3215*, 2014.
- 806 Adith Swaminathan, Akshay Krishnamurthy, Alekh Agarwal, Miroslav Dudík, John Langford,
807 Damien Jose, and Imed Zitouni. Off-policy evaluation for slate recommendation. *arXiv preprint*
808 *arXiv:1605.04812*, 2016.
- 809 Shayan A Tabrizi, Javid Dadashkarimi, Mostafa Dehghani, Hassan Nasr Esfahani, and Azadeh
810 Shakery. Revisiting optimal rank aggregation: A dynamic programming approach. In *Proceedings*
811 *of the 2015 International Conference on The Theory of Information Retrieval*, pages 353–356, 2015.
- 812 Yuval Tassa, Yotam Doron, Alistair Muldal, Tom Erez, Yazhe Li, Diego de Las Casas, David Budden,
813 Abbas Abdolmaleki, Josh Merel, Andrew Lefrancq, et al. Deepmind control suite. *arXiv preprint*
814 *arXiv:1801.00690*, 2018.
- 815 Yi Tay, Dara Bahri, Donald Metzler, Da-Cheng Juan, Zhe Zhao, and Che Zheng. Synthesizer:
816 Rethinking self-attention in transformer models. *arXiv preprint arXiv:2005.00743*, 2020a.
- 817 Yi Tay, Mostafa Dehghani, Samira Abnar, Yikang Shen, Dara Bahri, Philip Pham, Jinfeng Rao,
818 Liu Yang, Sebastian Ruder, and Donald Metzler. Long range arena: A benchmark for efficient
819 transformers. *arXiv preprint arXiv:2011.04006*, 2020b.

- 820 Yi Tay, Mostafa Dehghani, Dara Bahri, and Donald Metzler. Efficient transformers: A survey. *arXiv*
821 *preprint arXiv:2009.06732*, 2020c.
- 822 Yi Tay, Mostafa Dehghani, Jai Gupta, Dara Bahri, Vamsi Aribandi, Zhen Qin, and Donald Metzler. Are
823 pre-trained convolutions better than pre-trained transformers? *arXiv preprint arXiv:2105.03322*,
824 2021.
- 825 Ilya Tolstikhin, Olivier Bousquet, Sylvain Gelly, and Bernhard Schoelkopf. Wasserstein auto-encoders.
826 *arXiv preprint arXiv:1711.01558*, 2017.
- 827 Antonio Torralba and Alexei A Efros. Unbiased look at dataset bias. In *CVPR 2011*, pages 1521–1528.
828 IEEE, 2011.
- 829 Michael Tschannen, Josip Djolonga, Marvin Ritter, Aravindh Mahendran, Neil Houlsby, Sylvain Gelly,
830 and Mario Lucic. Self-supervised learning of video-induced visual invariances. In *Proceedings of*
831 *the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13806–13815, 2020.
- 832 Alan M. Turing. Computing machinery and intelligence. *Mind*, LIX(236):433–460, 1950.
- 833 Chris Van Der Lee, Albert Gatt, Emiel Van Miltenburg, Sander Wubben, and Emiel Kraahmer. Best
834 practices for the human evaluation of automatically generated text. In *Proceedings of the 12th*
835 *International Conference on Natural Language Generation*, pages 355–368, 2019.
- 836 Hado Van Hasselt, Arthur Guez, and David Silver. Deep reinforcement learning with double q-learning.
837 In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 30, 2016.
- 838 Clara Vania, Phu Mon Htut, William Huang, Dhara Mungra, Richard Yuanzhe Pang, Jason Phang,
839 Haokun Liu, Kyunghyun Cho, and Samuel R. Bowman. Comparing test sets with item response
840 theory. *arXiv preprint arXiv:2106.00840*, 2020.
- 841 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
842 Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing*
843 *systems*, pages 5998–6008, 2017.
- 844 Bastiaan S Veeling, Jasper Linmans, Jim Winkens, Taco Cohen, and Max Welling. Rotation
845 equivariant cnns for digital pathology. In *International Conference on Medical image computing*
846 *and computer-assisted intervention*, pages 210–218. Springer, 2018.
- 847 Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer
848 Levy, and Samuel R Bowman. Superglue: A stickier benchmark for general-purpose language
849 understanding systems. *arXiv preprint arXiv:1905.00537*, 2019.
- 850 Ziyu Wang, Tom Schaul, Matteo Hessel, Hado Hasselt, Marc Lanctot, and Nando Freitas. Dueling
851 network architectures for deep reinforcement learning. In *International conference on machine*
852 *learning*, pages 1995–2003. PMLR, 2016.
- 853 Adina Williams, Nikita Nangia, and Samuel R Bowman. A broad-coverage challenge corpus for
854 sentence understanding through inference. *arXiv preprint arXiv:1704.05426*, 2017.
- 855 Jianxiong Xiao, James Hays, Krista A Ehinger, Aude Oliva, and Antonio Torralba. Sun database:
856 Large-scale scene recognition from abbey to zoo. In *2010 IEEE computer society conference on*
857 *computer vision and pattern recognition*, pages 3485–3492. IEEE, 2010.
- 858 Chhavi Yadav and Léon Bottou. Cold case: The lost mnist digits. *arXiv preprint arXiv:1905.10498*,
859 2019.
- 860 Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le.
861 Xlnet: Generalized autoregressive pretraining for language understanding. *Advances in neural*
862 *information processing systems*, 32, 2019.
- 863 Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Ruslan Salakhutdinov, and Quoc V. Le.
864 Xlnet: Generalized autoregressive pretraining for language understanding, 2020.
- 865 Xinyang Yi, Ji Yang, Lichan Hong, Derek Zhiyuan Cheng, Lukasz Heldt, Aditee Kumthekar, Zhe Zhao,
866 Li Wei, and Ed Chi. Sampling-bias-corrected neural modeling for large corpus item recommenda-
867 tions. In *Proceedings of the 13th ACM Conference on Recommender Systems*, pages 269–277, 2019.

- 868 Xiaohua Zhai, Joan Puigcerver, Alexander Kolesnikov, Pierre Ruysen, Carlos Riquelme, Mario
869 Lucic, Josip Djolonga, Andre Susano Pinto, Maxim Neumann, Alexey Dosovitskiy, et al. A
870 large-scale study of representation learning with the visual task adaptation benchmark. *arXiv*
871 *preprint arXiv:1910.04867*, 2019.
- 872 Shuai Zhang, Lina Yao, Aixin Sun, and Yi Tay. Deep learning based recommender system: A survey
873 and new perspectives. *ACM Computing Surveys (CSUR)*, 52(1):1–38, 2019.
- 874 Tianhao Zhang. The need for performance based assessments, 2021. URL <https://covariant.ai/news/performance-based-assessments>.
875
- 876 Yin Zhang, Derek Zhiyuan Cheng, Tiansheng Yao, Xinyang Yi, Lichan Hong, and Ed H Chi. A model
877 of two tales: Dual transfer learning framework for improved long-tail item recommendation. *arXiv*
878 *preprint arXiv:2010.15982*, 2020.
- 879 Zhe Zhao, Lichan Hong, Li Wei, Jilin Chen, Aniruddh Nath, Shawn Andrews, Aditee Kumthekar,
880 Maheswaran Sathiamoorthy, Xinyang Yi, and Ed Chi. Recommending what video to watch next:
881 a multitask ranking system. In *Proceedings of the 13th ACM Conference on Recommender Systems*,
882 pages 43–51, 2019.
- 883 Lei Zheng, Chun-Ta Lu, Lifang He, Sihong Xie, Huang He, Chaozhuo Li, Vahid Noroozi, Bowen Dong,
884 and S Yu Philip. Mars: Memory attention-aware recommender system. In *2019 IEEE International*
885 *Conference on Data Science and Advanced Analytics (DSAA)*, pages 11–20. IEEE, 2019.

886 **Checklist**

- 887 1. For all authors...
- 888 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
889 contributions and scope? [Yes]
- 890 (b) Did you describe the limitations of your work? [Yes] This paper, touches upon several
891 sub-topics that are connected to the benchmark lottery phenomena. Each of these
892 subtopics deserves a dedicated study, with for instance more empirical investigations
893 and deeper analysis on the social aspects of the problem. However, for some of these
894 topics, this paper solely shares some opinions based on limited observation. While we
895 believe all these connected subtopics are extremely important, we found them outside
896 the scope and focus of this paper. We briefly discuss this point in Section 8.
- 897 (c) Did you discuss any potential negative societal impacts of your work? [N/A] This paper
898 targets calling for more discussion on a topic that has societal impact on the academic
899 community, and how the shape and pace of progress can be affected by the benchmarking
900 process. We highlight some of the existing problems and share our opinion on potential
901 ways that can address the issue. While several parts that we discussed relate to how
902 progress in ML may impact the broader society, we believe the content of our paper
903 itself has no specific point with potentially negative impacts.
- 904 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
905 them? [Yes]
- 906 2. If you are including theoretical results...
- 907 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 908 (b) Did you include complete proofs of all theoretical results? [N/A]
- 909 3. If you ran experiments (e.g. for benchmarks)...
- 910 (a) Did you include the code, data, and instructions needed to reproduce the main
911 experimental results (either in the supplemental material or as a URL)? [N/A]
- 912 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were
913 chosen)? [N/A]
- 914 (c) Did you report error bars (e.g., with respect to the random seed after running experiments
915 multiple times)? [N/A]
- 916 (d) Did you include the total amount of compute and the type of resources used (e.g., type
917 of GPUs, internal cluster, or cloud provider)? [N/A]
- 918 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 919 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 920 (b) Did you mention the license of the assets? [Yes]
- 921 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 922 (d) Did you discuss whether and how consent was obtained from people whose data you’re
923 using/curating? [N/A]
- 924 (e) Did you discuss whether the data you are using/curating contains personally identifiable
925 information or offensive content? [N/A]
- 926 5. If you used crowdsourcing or conducted research with human subjects...
- 927 (a) Did you include the full text of instructions given to participants and screenshots, if
928 applicable? [N/A]
- 929 (b) Did you describe any potential participant risks, with links to Institutional Review Board
930 (IRB) approvals, if applicable? [N/A]
- 931 (c) Did you include the estimated hourly wage paid to participants and the total amount
932 spent on participant compensation? [N/A]