# Position: AI Competitions Provide the Gold Standard for Empirical Rigor in GenAI Evaluation

D. Sculley [1]  Will Cukierski [2]  Phil Culliton [2]  Sohier Dane [2]  Maggie Demkin [2]  Ryan Holbrook [2]
Addison Howard [2]  Paul Mooney [2]  Walter Reade [2]  Megan Risdal [2]  Nate Keating [2]

## Abstract

In this position paper, we observe that empirical evaluation in Generative AI is at a crisis point since traditional ML evaluation and benchmarking strategies are insufficient to meet the needs of evaluating modern GenAI models and systems. There are many reasons for this, including the fact that these models typically have nearly unbounded input and output spaces, typically do not have a well defined ground truth target, and typically exhibit strong feedback loops and prediction dependence based on context of previous model outputs. On top of these critical issues, we argue that the problems of *leakage* and *contamination* are in fact the most important and difficult issues to address for GenAI evaluations. Interestingly, the field of AI Competitions has developed effective measures and practices to combat leakage for the purpose of counteracting cheating by bad actors within a competition setting. This makes AI Competitions an especially valuable (but underutilized) resource. Now is time for the field to view AI Competitions as the gold standard for empirical rigor in GenAI evaluation, and to harness and harvest their results with according value.

## 1. Introduction

As Generative AI (GenAI) models such as Large Language Models (LLMs) become ever more important to the field and to the world at large, it has become clear that performing empirical evaluation of these models and methods is extremely difficult to do in a rigorous and comprehensive way. This difficulty is of course not due to lack of effort or expertise by researchers. Indeed, enormous effort and resources have been poured into creating myriad benchmarks and test cases (Chiang et al., 2024b; Fourrier et al., 2024;

Hendrycks et al., 2021; Cobbe et al., 2021; Zellers et al., 2019; Chen et al., 2021b). However, even accounting for these many important efforts and achievements, our position is that **the current state of evaluation is insufficient to meet the needs of this moment in GenAI for the field and for the world.**

In our view, the root cause of this insufficiency is that the evaluation needs of GenAI models fundamentally break the paradigm of traditional benchmarking that served the field of machine learning (ML) so well during decades of remarkable progress. This breakage goes beyond the familiar difficulty of defining what, exactly, is in the training data for an LLM. **In our view, we need a broader conception of generalization for GenAI** that moves beyond the idea of generalizing to new independently drawn examples from a stationary distribution, and instead refers to performing well on tasks that are entirely novel from a model's perspective. This higher bar is rooted in commonsense standards for human intelligence (Chollet, 2019; Dennett, 1991), but has far reaching consequences, most notably that it implies that the problems of *data leakage* and *contamination* in evaluation are the most pressing concerns.

Together, these factors imply that **rigorous and robust evaluation GenAI models requires a steady source of novel tasks structured to avoids leakage, contamination, and other forms of inadvertent "cheating".** Fortunately, AI Competitions—such as those hosted on Kaggle and similar platforms—offer a ready-made solution, providing a continual source of new tasks for evaluation and significant structures to avoid leakage and related issues.

### 1.1. Summarizing Our Position

Our position can be summarized by the following points:

- Traditional paradigms for ML evaluation are ill-equipped to meet the demands of GenAI Evaluation.

- Leakage should be viewed by the field as the most important pitfall to avoid in evaluations.

- GenAI evaluations should be considered leaked the moment it has been shared online or sent over the wire.

[1]Work done at Kaggle [2]Kaggle, Inc. Correspondence to: D. Sculley <d@sculley.ai>, Nate Keating <natekeating@kaggle.com>, Megan Risdal <meg@kaggle.com>.

- If we have to choose between reproducibility and robustness in GenAI evaluations, we should choose to prioritize robustness.

- We should replace the notion of reproducible static benchmarks with repeatable processes and procedures.

- The field should use established AI Competitions platforms as a renewable stream of novel evaluation tasks.

- The standards and practices developed that help AI Competitions guard against cheating should be viewed by the field as the gold standard for empirical rigor in evaluation.

- Meta-analyses should be valued as highly in the field of AI as they are in fields such as medicine.

### 1.2. Structure of this paper

In the remainder of this paper, we will first review the most typical structure and assumptions in traditional ML evaluation and discuss why it is insufficient for GenAI evaluations. We will examine the nature of generalization for GenAI, how this leads to specific concerns around leakage, and additionally show how goals of reproducibility and robustness in evaluation may be fundamentally at odds. We will then show how difficult the problem of leakage is even for traditional ML evaluations with some brief case studies, and look at current GenAI benchmarks that are aiming to overcome leakage and contamination. We finish with an examination of the ways that AI Competitions address these issues, discuss our recommendations, and examine alternate viewpoints. Our goal is to provide convincing support of the view that AI Competitions do indeed provide a gold standard for empirical rigor in evaluating GenAI models, and that the field should place accordingly high value and attention on their results.

## 2. Background: Revisiting Benchmarking

Traditional ML benchmarking has been founded on the idea of a *test-train split*, in which an evaluation is structured by training a model from scratch on a given portion of training data and then evaluating that trained model on a holdout set of test data (Mitchell, 1997). This conceptual structure is so fundamental to modern ML practice that it may sometimes be taken for granted and not carefully examined. So let us take a moment to reflect this basic structure and its implications.

In classical supervised ML, the most common traditional setup is to evaluate a model $f(\mathbf{x}) \to y$, with $\mathbf{x} \in \Re^d$ as feature vectors in some $d$ dimensional feature space and $y \in Y$ as a space of possible labels, such as $\{0, 1\}$ for binary classification or $(0, 1)$ for regression on probabilities. Labeled examples $(\mathbf{x}, y)$ are assumed to have come

from some distribution $D$. The `training set` $D_{train}$ and `test set` $D_{test}$ are each independently and identically drawn (IID) from $D$, with only the examples in the `training set` used to fit the model $f(\mathbf{x})$ and only the examples in the `test set` used to evaluate the model (Mitchell & Mitchell, 1997).

The IID requirement on test-train splits is often taken as a footnote in practice, but in reality it is a cornerstone of the robustness of this setup. The reason for this is that we fundamentally wish our evaluations to be interpretable as statements on *generalization* ability of our models: we wish to know how the model will perform on future, previously unseen data. But achieving this is harder than it may sound, because the ML models in question are often of extremely high dimensionality and thus may be prone to overfitting.

One approach for assessing generalization ability lies in the classic literature on statistical learning theory, providing generalization bounds for models based on qualities like their VC dimension and observed error during training that do not require the use of an additional holdout set (Vapnik, 1999). However, these theoretical bounds are unfortunately much too loose to be of practical value—all the more so in the age of ever larger models.

A second approach is to use additional data for evaluation. The issue here to be aware of is the classical statistical trap that correlation does not necessarily imply causation, and that trying to assess generalization ability of a model that has no specific mechanism for disambiguating correlation from causal factors may lead to wildly unreliable performance estimates. It is this issue that the IID assumption addresses—when we know that all test data is drawn IID from the same distribution as the training data, we know that all correlations that held at training time will reappear with the same characteristics at test time, and thus we can take performance on holdout test data as a reasonable estimate of generalization ability. The IID assumption in many ways enabled modern machine learning to advance as a research field, because it forms the theoretical underpinning of all evaluations. And indeed, it is a truism that moving ML models from research to deployed production is difficult precisely because of the fact that the IID assumption often does not hold in practice (Chen et al., 2021a).

### 2.1. The Rise of Reproducible Benchmarks

One statistical shortcut that immediately became standard practice was instead of drawing a new `training set` and `test set` from $D$ for every new evaluation, researchers would make one draw of each and use them as a canonical train / test set. The primary benefit of this approach (besides convenience) was that these paired test-train splits could now be used as *reproducible benchmarks*. All future researchers could replicate the exact problem
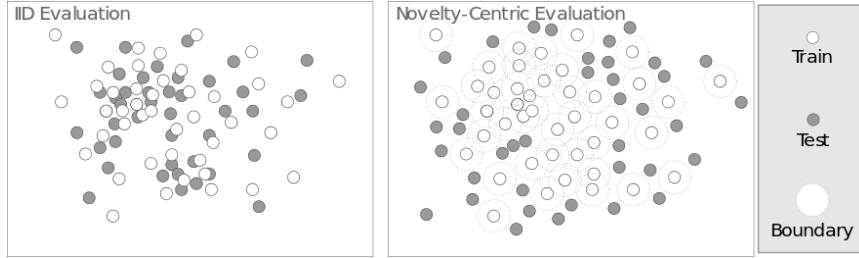
*Figure 1.* **IID Evaluations vs. Novelty-Centric Evaluations.** In the IID evaluation, left, both training and test data are drawn from the same distribution, resulting in significant overlap in examples in each set. In the novelty-centric version, right, no test example is allowed to be too similar to any given training example. We argue that this conceptualization more closely mirrors desired behavior for GenAI evaluations, where generalization is expected to connote the ability to respond well on totally novel inputs.

setup, giving a new untrained model the same training data for training and the same test data for evaluation, allowing for full apples-to-apples comparisons. This approach was wildly successful, with canonical benchmarks such as MNIST (LeCun & Cortes, 2010) and ImageNet (Deng et al., 2009) responsible for driving incredibly rapid progress in computer vision, for instance, and benchmarks such as (Rajpurkar et al., 2016; Marcus et al., 1993; Diemert Eustache, Betlei Artem et al., 2018) moving forward email spam classification, natural language processing, and many others. Websites such as the UCI repository (Kelly et al.) and OpenML (Vanschoren et al., 2014) among many others have been remarkably valuable to the field for this reason.

### 2.2. Surprisingly, Overfitting Was Not the Main Issue

Given that one of the fundamental concerns in evaluating models was ensuring that we could properly address generalization and avoid overfitting, it is reasonable to ask whether having a field re-use the same standard benchmark datasets in thousands of papers might not lead to overfitting. As authors, we were deeply surprised by the work of (Roelofs et al., 2019b) showing that in practice, this appears to actually not have been a problem. In (Recht et al., 2019), these authors shows that the rank ordering of ImageNet models when evaluated on brand new data was remarkably consistent with the rank ordering of those models on the benchmark data, despite its wide reuse. And in follow-up work, it was shown that similarly evaluation on public leaderboard data on Kaggle competitions was a remarkably good indicator of rank ordering on private holdout data, despite the risk of overfitting when many thousands of teams participate in the same challenge (Roelofs et al., 2019a).

### 3. Reconsidering Generalization for GenAI

As we recalled above, the IID assumption in traditional ML evaluations gives a clear conception of the idea of generalization: a model generalizes well if it accurately predicts the true-but-hidden label $y$ for labeled examples $(\mathbf{x}, y)$ drawn IID from the same fixed and stationary distribution $D$ from

which the model's training data was drawn. This was a cornerstone allowing the field of ML to progress effectively by narrowing the problem and enabling tractable statistical theory. But if we reflect on broader notions of intelligence, including those first proposed in the seminal paper by Turing (Turing, 1950), it is clear that this narrow-but-useful notion of generalization does not adequately reflect the deeper goals that GenAI is aiming to deliver on.

Instead, we believe **the form of generalization the field should most care about for GenAI is novelty-based generalization**—that is, generalizing well to problems and tasks that have truly never been seen before by the system in training or development.

More deeply, evaluations of reasoning and understanding that have been in our view easiest to design as tests often have the quality that solving the problems is hard (in a formal sense) or expensive, while answer verification is significantly easier or cheaper. This experience holds true in planning, solving mathematics problems, doing coding problems, solving riddles, and even formulating essays, and mirrors the fundamental quality of NP-hard problems. Once an answer for a problem is known to a subject, the ability to use that problem or very similar problems for that subject again in the future is fundamentally compromised.

In order to assess novelty-centric generalization, we need novelty-based evaluations. Informally, the goal of a novelty-based evaluation is to ensure that no evaluation task or example is too closely similar (for some definition of *similar* and some measure of *too close*) to any instance previously known to the model or system. We illustrate the distinction between IID-based generalization and novelty-based generalization in Figure 1 visually—we can imagine a small conceptual ring around each training example and ensure that no evaluation instance crosses any of those rings.

In our view, this novelty-centric view of generalization has already been implicitly adopted by many in the field as the true aspirational goal, and influences the design of important benchmarks including the LM Arena (Chiang et al., 2024a) among others—we are simply writing this *de facto* standard

down. We will now examine some of the implications.

## 3.1. The IID Assumption is Broken

While the IID assumption has often been broken in practice for traditional ML systems in real-world deployment with only modest harm, we believe that the IID assumption and the overall framework of neatly labeled examples $(\mathbf{x}, y)$ is broken beyond repair for GenAI evaluation. In particular, **the novelty-centric view of generalization strongly implies that evaluation examples should *not* be drawn from some identical distribution used for training**, but should instead be chosen or constructed with the explicit goal of avoiding high similarity with examples or data that the model has previously been exposed to.

We also note that the nature of typical GenAI models themselves leads to other ways that the IID assumption is broken. In particular, GenAI outputs are often far from independent, and instead use context of previous responses (for example, in multi-turn chat-style interfaces) to inform future responses, creating feedback loops that fully break ideas of stationarity. Finally, because the input spaces and output spaces are so vast (such as the space of all possible strings of up to a given size), the very notion of testing distributional equivalence is arguably vacuous.

## 3.2. Leakage and Contamination Are the Biggest Pitfalls

While the potential pitfall of overfitting receives strong attention, practitioners have long understood that *leakage* is an equally important and often more difficult problem in practice (Nisbet et al., 2009; Kaufman et al., 2012) Intuitively, leakage is any issue or structure in the construction of evaluation data that allows a model to "cheat" by using information that it should not have access to. In Section 4 we will look at a number of case studies on leakage and will show how hard it is to prevent leakage and how vigilant we must be even for traditional ML evaluations to avoid this pitfall. Here, we point out that leakage is an especially large problem for novelty-centric GenAI evaluations. This is because **novelty-centric GenAI evaluations have all of the leakage risks that traditional ML evaluations do, but also carry the additional burden of *novelty assurance*.**

A novelty-centric evaluation rests on assurance that a model has never before been exposed to data that are too close to the evaluation problems or tasks. While this may seem obvious, in practice it can be extremely difficult, as GenAI models like LLMs are often trained on enormous amounts of data and it can be extremely difficult to say for sure what similar data may or may not have been included. Indeed, leakage for GenAI is so important that specific forms of it have been given an additional name: *contamination* (Magar & Schwartz, 2022; Oren et al., 2023; Sainz et al., 2023; Balloccu et al., 2024). Contamination is said to occur when

evaluation datasets and benchmarks appear in training data.

To help give intuition for the breadth of this issue, consider that every major LLM we have tested so far (both open and proprietary) shows extensive detailed knowledge of the contents of standard test datasets from Kaggle. Consider the remarkably strong performance on many static benchmarks by LLMs that do not seem to correlate with strong performance on other tasks (Fourrier et al., 2024; Muennighoff et al., 2023; Zheng et al., 2023). Consider the question: if a model does particularly well on qualification exam normally given to humans, is this because the model has gained strong expertise or because example examinations have appeared in its training data, and how would we be able to disambiguate? Consider the difficulty in teasing apart exactly which data sources are or are not part of an openly shared dataset such as the open and widely used Nectar dataset (Zhu et al., 2024), which includes the description:

> Nectar's prompts are an amalgamation of diverse sources, including lmsys-chat-1M, ShareGPT, Antropic/hh-rlhf, UltraFeedback, Evol-Instruct, and Flan. Nectar's 7 responses per prompt are primarily derived from a variety of models, namely GPT-4, GPT-3.5-turbo, GPT-3.5-turbo-instruct, LLama-2-7B-chat, and Mistral-7B-Instruct, alongside other existing datasets and models.

Together, these practical realities and considerations force leakage to the forefront of problems that must be addressed by any serious GenAI evaluation.

## 4. Leakage Case Studies

Because leakage and contamination are the most important hurdles to solve for GenAI evaluations, it is useful to study them in depth, beginning with leakage from traditional ML evaluations. Here, we draw on lessons learned surveying more than a decade of Kaggle competitions, in which a broad range of leakage issues have been identified through intense scrutiny of a large community. Experience has shown that the risk of leakage is compounded in open machine learning challenge benchmarks, where teams will exploit (knowingly or unknowingly) anything that gives an advantage on the leaderboard.

Leakage can occur simply by how observations are ordered. An extreme example occurred during the SETI Breakthrough Listen competition (Siemion et al., 2021), where data was processed in order of its class label. The file timestamps were not reset, and competitors found it trivial to make predictions based on file metadata. A more subtle example occurred during the TalkingData AdTracking Fraud Detection Challenge (Yin et al., 2018), where the data was

mistakenly sorted so that if multiple events were present within the same timestamp, any positive labels occurred after negative labels.

Ironically, randomization can also be a source of leakage. An example occurred during the Predict AI Model Runtime competition (Phothilimthana et al., 2023) where teams had to rank order the runtimes of 5 different subsets of data, each subset requiring a different model. Two of the buckets were randomized using the same seed, and teams discovered that using ordering of one bucket on another improved their scores.

Any data that is synthetically generated is highly prone to having artifacts that leak information. Again, in the SETI Breakthrough Listen competition, synthetically-created "ET" signals were injected into real radio telescope signals. Care was taken with normalization to ensure the averages and standard deviations of the injections matched the background signals. The code that created the injected signals used FP16 while the background signals were FP32. This created a minute difference in the mean and standard deviations between positive and negative samples, but enough to differentiate the classes based on this information alone.

Private evaluation data leaking to the public during an open challenge is a risk that needs to be considered. During the LANL Earthquake Prediction challenge (RL et al., 2019), for example, the dataset was described in a research paper, including some summary statistics and a graph. A few teams discovered this and were able to utilize it to their advantage.

Space precludes a larger set of case studies, but experience from practitioners in preparing competitions shows that each AI Competition has more ways to go wrong than to go right, and that paranoia and vigilance are helpful practices. In addition to the failure modes highlighted above, other broad categories include future data leaks, the many ways metadata can leak information (e.g., the model of a medical machine being correlated to disease incidence, medical images that include a hand-drawn circle around a concerning skin lesion, image aspect ratio, file size on disk, etc.), old versions of the private evaluation dataset that were not kept private, the ability of teams to reverse a synthetic data generation process or to re-assemble data that has been split up, reverse engineering data obfuscation, near duplication between training and test observations, etc. These are not hypothetical; they have all occurred in challenges created by competent, careful teams, and highlight the very real difficulty of creating leak-free competitions and benchmarks.

### 4.1. Reproducibility and Robustness in Conflict

Because of the importance of leakage and the practical difficulty in ensuring that leakage does not impact GenAI evaluations, we argue that it is simplest and safest to adopt a leakage rule of thumb that **an evaluation should be considered** `leaked` **the moment it has been shared online or sent over the wire**. Adopting this rule of thumb significantly improves our ability to trust the results of evaluations and gives them substantially more robustness. However, it also critically weakens the notion of reproducibility. It is the position of this paper that this is a fundamental tension, analogous to the Heisenberg Uncertainty Principle from quantum physics, and that we simply cannot have a published static benchmark that is robust to leakage—no matter the good intentions of the researchers, it is just too hard to avoid contamination and to broadly trust results from such a benchmark.

Instead, we must seek alternative strategies and structures to create leak-proof evaluations.

## 5. Evaluations Aiming to Avoid Leakage

Conscientious researchers have been aware of the issue of leakage in novelty-based evaluations for GenAI and have proposed new benchmarks that attempt to control or mitigate leakage through various design mechanisms. Here we review key examples, the mechanisms used to control for leakage, and briefly discuss their benefits and drawbacks.

### 5.1. Unreleased Holdout Sets

The SEAL Leaderboards (Scale AI, 2024), ARC-AGI (Chollet, 2019), FrontierMath (Glazer et al., 2024), and Humanity's Last Exam (Phan et al., 2025) benchmarks are composed of private test questions manually created by domain experts. The test sets, model responses, and evaluation runs are not published publicly to prevent leakage of test data.

Unreleased holdout sets can be effective at mitigating risks of leakage. However, they may have limitations in evaluating proprietary API-based models where test data must necessarily be sent over the internet to third party servers. While many leading AI model providers grant controls preventing logging or storage of user prompts, this still requires a level of trust. In particular, providers must be trusted not to change their policies, but even more importantly all researchers touching the evaluation data must be trusted to follow these practices without error.

Additionally, as holdout sets and evaluation runs must necessarily be kept private, results are not reproducible by researchers. To mitigate this, some benchmarks take a hybrid approach. For example, the FACTS Grounding Leaderboard (Jacovi et al., 2025) publishes half the test set publicly which enables partial reproducibility and better understanding of the benchmark. A model's performance can be compared on private and public parts of test sets to identify models that may have (intentionally or not) trained on test.
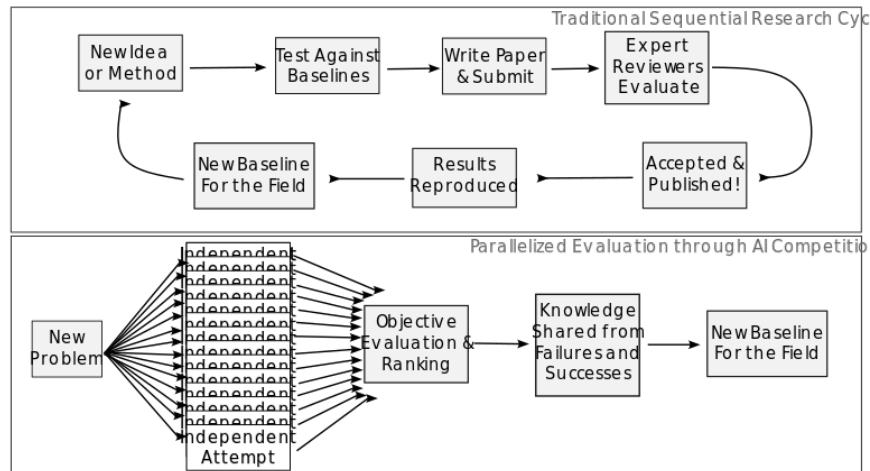
*Figure 2.* **Comparing sequential and parallelized evaluation structures.** In the traditional research structure, top, each new idea is evaluated in a linear sequence that typically requires several months for a single pass. The parallelized structure, bottom, allows hundreds or thousands of approaches to be simultaneously.

## 5.2. Dynamic Benchmarks

LiveBench (White et al., 2024), LiveCodeBench (Jain et al., 2024), and SWE-Bench (Jimenez et al., 2024) benchmarks frequently update test sets from sources that refresh naturally over time. For example, LiveBench test set questions are updated weekly from sources such as newly published news articles or papers on arXiv. By testing models only with very recent data, benchmarks can mitigate the risk that such data was included in model training. LiveBench also does not release the most recently added test data so that a significant percentage of questions are always heldout.

Dynamic benchmarks have some advantages over unreleased holdout sets. By frequently refreshing data, older test sets can be released publicly to improve reproducibility and trust. However, using data that is publicly available on the internet—even if only very recently—does not pass our rule of thumb in 4.1 and is a potential source of leakage. Moreover, by changing the test set frequently, benchmark creators must be careful to ensure they are not "moving the goal posts." Additionally, dynamic benchmarks come with higher maintenance costs and sources of frequently refreshing data are not available in many domains and may be infeasible to collect.

## 5.3. Community Benchmarks

The LM Arena (Chiang et al., 2024a), formerly known as LMSYS Chatbot Arena, is a collection of benchmarks that draws on community votes of head-to-head match-ups between LLMs on user prompts or tasks. By outsourcing test data collection and evaluation to users at test time, the benchmark has a constant fresh source of novel test questions.

Community benchmarks are difficult to build and maintain. To evaluate many models, the number of votes required can be very large, necessitating a large and constant pool of voters. Community benchmarks also don't work for all tasks, for example, tasks that require specialized knowledge or which might take humans many days to complete will not scale to human rating. Community benchmarks are also necessarily biased by any sampling effects; the diversity and distribution of voters can affect results and great attention and care is required to filter out low quality, duplicate, or contaminated results.

## 6. AI Competitions as Structural Solution

AI Competitions typified by platforms like Kaggle and others offer an "embarrassingly parallel" structure to empirical evaluation shown in Figure 2 that hearkens back to the classic MapReduce structure from parallel computing (Dean & Ghemawat, 2008). In this structure, independent teams of researchers—often numbering in the thousands—each compete to solve a given problem, and in so doing create an evaluation of many different approaches in one massive parallel effort. Here we show ways that this structure offers useful benefits to the the problem of GenAI evaluation writ large.

### 6.1. Parallelization Improves Robustness

The risk of leakage and contamination starts as soon as an evaluation is shared publicly or evaluation data is sent across the wire. This leads to a problem: how can we fairly compare different models and systems in a valid way that ensures robustness and avoids inadvertent invalidation of results from leakage and contamination?

The parallelized structure of AI Competitions provides a useful solution to this issue. **Novelty-centric evaluations can happen simultaneously, in parallel, ensuring that each new task is indeed novel to each of the thousands of models at time of testing.** Because the independent teams each pursue different models, ideas, and approaches, this structure yields direct apples-to-apples benchmarking and a form of real-time reproduction of results.

In addition, competition platforms such as Kaggle can serve as trusted keepers of hidden test data by running isolated code competitions, where competitors submit their models to be run on an isolated, secure backend without access to internet. By evaluating all models securely offline, competitions platforms can guarantee no hidden test data is leaked.

Finally, competitions hosted on large community platforms offer additional non-structural characteristics which represent best practices the industry should adopt to further improve empirical rigor. Competitions encourage or often require open sharing of code, data, and experimental details, including both successes and failures. This transparency facilitates reproduction of results, fosters trust in new baselines, and accelerates the dissemination of knowledge within the research and practitioner communities.

## 6.2. Leak-Proof Competition Structures

While preventing traditional leakage remains a challenge for competition-style evaluation as it does everywhere, competitions can be uniquely structured to mitigate this issue particularly well. Furthermore, the structure of competitions with many thousands of research teams ensures that when issues of leakage do occur, they are rapidly discovered, shared, and addressed simultaneously across all research efforts happening on the task in parallel.

We provide some examples of competitions that demonstrate the feasibility of leak-proof evaluation design. Employing strategies such as prospective ground truth, novel task generation, and post-deadline data collection, generally combined with test data that is directly inaccessible to competitors, competitions can provide a robust and reliable platform for novel evaluation of GenAI models. These best practices should be considered and adapted as blueprints for future competition and benchmark design.

**Prospective Ground Truth**  Prospective ground truth is a strategy for leakage mitigation whereby test set labels are completely unknown to the world during the active training phase of a competition.

The Critical Assessment of protein Function Annotation (CAFA) 5 challenge (Friedberg et al., 2023) is an example of a competition that uses a prospective ground truth to mitigate leakage. The competition took as its test set proteins whose sequences were known, but whose functional annotations had not yet been determined in a wet lab. Nearly two thousand participants across 1,625 independent teams therefore effectively developed models predicting the function of a set of proteins without any ground truth yet available to any human or model during an active training phase. Months later, the final evaluation was determined following a "curation phase" on the basis of newly published protein functions. This novelty makes the competition reasonably leak-proof.

**Novel Task Generation**  Another approach to designing leak-proof competitions is generating novel tasks altogether in which test data doesn't resemble training data and therefore demands meaningful generalization.

The AI Mathematical Olympiad (AIMO) challenges (XTX Investments, 2024; Frieder et al., 2024), designed to motivate open progress on human-level mathematical reasoning capabilities in GenAI systems, used this approach. In these challenges, competitors were tasked with solving national-level math challenges. Because many, if not all, AI models used by competitors were trained on internet-scale data, test-train leakage poses a significant challenge in the evaluation of their mathematical reasoning capabilities. Fresh sets of novel math problems were therefore created specifically for the competition by an international team of mathematicians, making it highly unlikely that the data has been leaked or contaminated.

**Post-Deadline Data Collection**  Post-deadline data collection is a leakage mitigation strategy used in a number of competitions which are similar to prospective ground truth competitions except rather than evaluating on newly available labels, solutions are evaluated on completely newly generated data. There are many examples of this competition design, two of which are described below.

In the WSDM Cup – Multilingual Chatbot Arena competition (Chiang et al., 2024a) hosted by LMSYS.org, participants were tasked with building solutions predicting human preferences between LLMs in head-to-head match-ups based on multilingual conversation and rating data from LM Arena. Similar to CAFA 5, this competition was designed with an active training phase followed by a data collection phase after which final models were evaluated against brand new conversations after the submission deadline in order to prevent leakage.

The Konwinski Prize (Konwinski et al., 2024) is another form of post-deadline data collection. This competition, hosted by Andy Konwinski, is a contamination-free version of SWE-Bench which evaluates LLMs on their ability to resolve real-world GitHub issues. It uses a time-based hold-out strategy in which submitted models are frozen for three

months and then evaluated on fresh GitHub issues that have been collected during the intervening time.

# 7. Recommendations for the Field

As a field, we need to overhaul our standard practices to ensure that GenAI evaluations are rigorous and reliable— and that they continue to be viewed as such by the field and the broader world.

**Move away from static benchmarks and towards evergreen repeatable processes.** Due to the risks of leakage and contamination, we believe that static benchmarks should be de-emphasized in importance for GenAI evaluations. (Indeed, anecdotally we see that both researchers and practitioners are taking results from such benchmarks with ever larger grains of salt.) Instead, we need a steady renewable pipeline of novel tasks and problems, and we need to evaluate hundreds or thousands of models in parallel on each of them so that the results are directly comparable and avoid the risks of later contamination and leakage. In this way, evaluations are best viewed as results from a point in time rather than an an immutable final conclusion.

**View the steady stream of AI Competitions as a resource for the field.** Using the pipeline of high quality AI Competitions hosted on platforms like Kaggle is one way to create a renewable pipeline. These structures already exist and are already being used to some degree in this way. However, as a field, we can do more to distill, analyze, and share findings from these competitions through meta-analyses. Indeed, while meta-analysis is a common and highly valued form of academic contribution in fields such as medicine, such papers are extremely rare in our field. We can and should change this through mechanisms that include specialized workshops, conference tracks, journal special topics, and though updated language in calls for papers emphasizing the value of meta-analyses.

**Adopt and improve on the anti-cheating structures from AI Competitions to improve standard practice for GenAI evaluations.** Furthermore, as a field, we can learn from the best practices that have been developed by AI Competitions – the techniques and practices that have been created to combat intentional cheating by bad actors are equally valuable in creating evaluation structures that combat unintentional issues such as leakage and contamination that may invalidate empirical results. A cheat-proof structure is one that provides assurance to researchers that they will not accidentally cheat themselves. We also need to augment and further improve these structures, for example by creating a field-wide standard that major API-based models agree to follow to explicitly avoid collecting or training on data that may appear in evaluations.

# 8. Alternative Views

All position papers should consider opposing views, and ours is no exception. One reasonable alternative view is that the current state of benchmarking is proceeding well without the need for additional intervention. The many new static benchmarks are appearing on platforms like Hugging Face, OpenML, and Kaggle on a near-daily basis and may serve as the steady stream of novel tasks that we described as necessary for the field. While we applaud all efforts to create new benchmarks, we do fundamentally believe that static benchmarks should be considered to have been effectively invalidated once they have been published, and thus it is the *time component* of AI Competitions that provides unique additional value.

Another reasonable viewpoint is that current existing benchmarks that attempt to be leak-proof are sufficient. The most notable one to consider for this viewpoint are the Elo-based side-by-side rankings produced by human raters via LMSYS.org's LMArena. Having an open-loop for the community to provide an unbounded stream of new inputs and judgments is indeed appealing and is a strong step towards solving many of these issues. However, we believe there are limits to what can be achieved in terms of novelty and rigor with an anonymous crowd-based source of tasks and problems, and that AI Competitions allow for the injection of specific domain expertise and carefully crafted test cases that will fully stress test the next generation of GenAI models.

A third reasonable viewpoint is that the metaphorical ship has sailed on the value of academic evaluations for GenAI models. In this paradigm, performance on literal real-world tasks in production deployments may offer the most valid test of GenAI capabilities. In this alternative viewpoint, independent evaluations have little value and each practitioner or group should evaluate fully on their own terms. While this approach is unavoidable for highly specialized domains and applications, we do believe that there is compelling reason to continue independent evaluations of models in general, as the history of the field has shown that these forms of evaluation drive progress in the broadest and most rapid ways. Without controlled, empirical study we as a field risk losing broadly shared knowledge into *why* models perform well or poorly on certain tasks. Openly sharing this understanding is critical for unlocking paths to further progress in this rapidly advancing field.

# References

Balloccu, S., Schmidtová, P., Lango, M., and Dušek, O. Leak, cheat, repeat: Data contamination and evaluation malpractices in closed-source llms, 2024. URL https://arxiv.org/abs/2402.03927.

Chen, C., Murphy, N. R., Parisa, K., Sculley, D., and Underwood, T. *Reliable Machine Learning.* " O'Reilly Media, Inc.", 2021a.

Chen, M., Tworek, J., Jun, H., Yuan, Q., de Oliveira Pinto, H. P., Kaplan, J., Edwards, H., Burda, Y., Joseph, N., Brockman, G., Ray, A., Puri, R., Krueger, G., Petrov, M., Khlaaf, H., Sastry, G., Mishkin, P., Chan, B., Gray, S., Ryder, N., Pavlov, M., Power, A., Kaiser, L., Bavarian, M., Winter, C., Tillet, P., Such, F. P., Cummings, D., Plappert, M., Chantzis, F., Barnes, E., Herbert-Voss, A., Guss, W. H., Nichol, A., Paino, A., Tezak, N., Tang, J., Babuschkin, I., Balaji, S., Jain, S., Saunders, W., Hesse, C., Carr, A. N., Leike, J., Achiam, J., Misra, V., Morikawa, E., Radford, A., Knight, M., Brundage, M., Murati, M., Mayer, K., Welinder, P., McGrew, B., Amodei, D., McCandlish, S., Sutskever, I., and Zaremba, W. Evaluating large language models trained on code, 2021b. URL https://arxiv.org/abs/2107.03374.

Chiang, W.-L., Zheng, L., Dunlap, L., Gonzalez, J. E., Stoica, I., Mooney, P., Dane, S., Howard, A., and Keating, N. LMSYS - Chatbot Arena Human Preference Predictions. https://kaggle.com/competitions/lmsys-chatbot-arena, 2024a. Kaggle.

Chiang, W.-L., Zheng, L., Sheng, Y., Angelopoulos, A. N., Li, T., Li, D., Zhang, H., Zhu, B., Jordan, M., Gonzalez, J. E., and Stoica, I. Chatbot arena: An open platform for evaluating llms by human preference, 2024b. URL https://arxiv.org/abs/2403.04132.

Chollet, F. On the measure of intelligence, 2019. URL https://arxiv.org/abs/1911.01547.

Cobbe, K., Kosaraju, V., Bavarian, M., Chen, M., Jun, H., Kaiser, L., Plappert, M., Tworek, J., Hilton, J., Nakano, R., Hesse, C., and Schulman, J. Training verifiers to solve math word problems, 2021. URL https://arxiv.org/abs/2110.14168.

Dean, J. and Ghemawat, S. Mapreduce: simplified data processing on large clusters. *Commun. ACM*, 51(1): 107–113, January 2008. ISSN 0001-0782. doi: 10. 1145/1327452.1327492. URL https://doi.org/10.1145/1327452.1327492.

Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. Imagenet: A large-scale hierarchical image database. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, pp. 248–255. IEEE, 2009. URL https://ieeexplore.ieee.org/abstract/document/5206848/.

Dennett, D. C. *Consciousness Explained.* Penguin Books, 1991.

Diemert Eustache, Betlei Artem, Renaudin, C., and Massih-Reza, A. A large scale benchmark for uplift modeling. In *Proceedings of the AdKDD and TargetAd Workshop, KDD, London,United Kingdom, August, 20, 2018*. ACM, 2018.

Fourrier, C., Habib, N., Lozovskaya, A., Szafer, K., and Wolf, T. Open llm leaderboard v2. https://huggingface.co/spaces/open-llm-leaderboard/open_llm_leaderboard, 2024.

Friedberg, I., Radivojac, P., Paolis, C. D., Piovesan, D., Joshi, P., Reade, W., and Howard, A. CAFA 5 Protein Function Prediction. https://kaggle.com/competitions/cafa-5-protein-function-prediction, 2023. Kaggle.

Frieder, S., Bealing, S., Nikolaiev, A., Smith, G. C., Buzzard, K., Gowers, T., Liu, P. J., Loh, P.-S., Mackey, L., de Moura, L., Roberts, D., Sculley, D., Tao, T., Balduzzi, D., Coyle, S., Gerko, A., Holbrook, R., Howard, A., and Markets, X. Ai mathematical olympiad - progress prize 2. https://kaggle.com/competitions/ai-mathematical-olympiad-progress-prize-2, 2024. Kaggle.

Glazer, E., Erdil, E., Besiroglu, T., Chicharro, D., Chen, E., Gunning, A., Olsson, C. F., Denain, J.-S., Ho, A., de Oliveira Santos, E., Järviniemi, O., Barnett, M., Sandler, R., Vrzala, M., Sevilla, J., Ren, Q., Pratt, E., Levine, L., Barkley, G., Stewart, N., Grechuk, B., Grechuk, T., Enugandla, S. V., and Wildon, M. Frontiermath: A benchmark for evaluating advanced mathematical reasoning in ai, 2024. URL https://arxiv.org/abs/2411.04872.

Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D., and Steinhardt, J. Measuring massive multitask language understanding, 2021. URL https://arxiv.org/abs/2009.03300.

Jacovi, A., Wang, A., Alberti, C., Tao, C., Lipovetz, J., Olszewska, K., Haas, L., Liu, M., Keating, N., Bloniarz, A., Saroufim, C., Fry, C., Marcus, D., Kukliansky, D., Tomar, G. S., Swirhun, J., Xing, J., Wang, L., Gurumurthy, M., Aaron, M., Ambar, M., Fellinger, R., Wang, R., Zhang, Z., Goldshtein, S., and Das, D. The facts grounding leaderboard: Benchmarking llms' ability to ground responses to long-form input, 2025. URL https://arxiv.org/abs/2501.03200.

Jain, N., Han, K., Gu, A., Li, W.-D., Yan, F., Zhang, T., Wang, S., Solar-Lezama, A., Sen, K., and Stoica, I.

Livecodebench: Holistic and contamination free evaluation of large language models for code. *arXiv preprint arXiv:2403.07974*, 2024.

Jimenez, C. E., Yang, J., Wettig, A., Yao, S., Pei, K., Press, O., and Narasimhan, K. Swe-bench: Can language models resolve real-world github issues?, 2024. URL https://arxiv.org/abs/2310.06770.

Kaufman, S., Rosset, S., Perlich, C., and Stitelman, O. Leakage in data mining: Formulation, detection, and avoidance. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 6(4), December 2012. ISSN 1556-4681. doi: 10.1145/2382577.2382579. URL https://doi.org/10.1145/2382577.2382579.

Kelly, M., Longjohn, R., and Nottingham, K. UCI machine learning repository. URL https://archive.ics.uci.edu.

Konwinski, A., Rytting, C., Shaw, J. F. A., Dane, S., Reade, W., and Demkin, M. Konwinski Prize. https://kaggle.com/competitions/konwinski-prize, 2024. Kaggle.

LeCun, Y. and Cortes, C. MNIST handwritten digit database. 2010. URL http://yann.lecun.com/exdb/mnist/.

Magar, I. and Schwartz, R. Data contamination: From memorization to exploitation. In Muresan, S., Nakov, P., and Villavicencio, A. (eds.), *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pp. 157–165, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-short.18. URL https://aclanthology.org/2022.acl-short.18/.

Marcus, M. P., Santorini, B., and Marcinkiewicz, M. A. Building a large annotated corpus of English: The Penn Treebank. *Computational Linguistics*, 19(2):313–330, 1993. URL https://aclanthology.org/J93-2004/.

Mitchell, T. M. *Machine learning*, volume 1. McGraw-hill New York, 1997.

Mitchell, T. M. and Mitchell, T. M. *Machine learning*, volume 1. McGraw-hill New York, 1997.

Muennighoff, N., Tazi, N., Magne, L., and Reimers, N. MTEB: Massive Text Embedding Benchmark, 2023. URL https://arxiv.org/abs/2210.07316.

Nisbet, R., Elder, J., and Miner, G. *Handbook of Statistical Analysis & Data Mining Applications*. Elsevier, Inc, 2009.

Oren, Y., Meister, N., Chatterji, N., Ladhak, F., and Hashimoto, T. B. Proving test set contamination in black box language models, 2023. URL https://arxiv.org/abs/2310.17623.

Phan, L., Gatti, A., Han, Z., Li, N., Hu, J., Zhang, H., Shi, S., Choi, M., Agrawal, A., Chopra, A., Khoja, A., Kim, R., Hausenloy, J., Zhang, O., Mazeika, M., Anderson, D., Nguyen, T., Mahmood, M., Feng, F., Feng, S. Y., Zhao, H., Yu, M., Gangal, V., Zou, C., Wang, Z., Wang, J. P., Kumar, P., Pokutnyi, O., Gerbicz, R., Popov, S., Levin, J.-C., Kazakov, M., Schmitt, J., Galgon, G., Sanchez, A., Lee, Y., Yeadon, W., Sauers, S., Roth, M., Agu, C., Riis, S., Giska, F., Utpala, S., Giboney, Z., Goshu, G. M., of Arc Xavier, J., Crowson, S.-J., Naiya, M. M., Burns, N., Finke, L., Cheng, Z., Park, H., Fournier-Facio, F., Wydallis, J., Nandor, M., Singh, A., Gehrunger, T., Cai, J., McCarty, B., Duclosel, D., Nam, J., Zampese, J., Hoerr, R. G., Bacho, A., Loume, G. A., Galal, A., Cao, H., Garretson, A. C., Sileo, D., Ren, Q., Cojoc, D., Arkhipov, P., Qazi, U., Li, L., Motwani, S., de Witt, C. S., Taylor, E., Veith, J., Singer, E., Hartman, T. D., Rissone, P., Jin, J., Shi, J. W. L., Willcocks, C. G., Robinson, J., Mikov, A., Prabhu, A., Tang, L., Alapont, X., Uro, J. L., Zhou, K., de Oliveira Santos, E., Maksimov, A. P., Vendrow, E., Zenitani, K., Guillod, J., Li, Y., Vendrow, J., Kuchkin, V., Ze-An, N., Marion, P., Efremov, D., Lynch, J., Liang, K., Gritsevskiy, A., Martinez, D., Pageler, B., Crispino, N., Zvonkine, D., Fraga, N. W., Soori, S., Press, O., Tang, H., Salazar, J., Green, S. R., Brüssel, L., Twayana, M., Dieuleveut, A., Rogers, T. R., Zhang, W., Li, B., Yang, J., Rao, A., Loiseau, G., Kalinin, M., Lukas, M., Manolescu, C., Mishra, S., Kamdoum, A. G. K., Kreiman, T., Hogg, T., Jin, A., Bosio, C., Sun, G., Coppola, B. P., Tarver, T., Heidinger, H., Sayous, R., Ivanov, S., Cavanagh, J. M., Shen, J., Imperial, J. M., Schwaller, P., Senthilkuma, S., Bran, A. M., Dehghan, A., Algaba, A., Verbeken, B., Noever, D., V, R. P., Schut, L., Sucholutsky, I., Zheltonozhskii, E., Lim, D., Stanley, R., Sivarajan, S., Yang, T., Maar, J., Wykowski, J., Oller, M., Sandlin, J., Sahu, A., Hu, Y., Fish, S., Heydari, N., Apronti, A., Rawal, K., Vilchis, T. G., Zu, Y., Lackner, M., Koppel, J., Nguyen, J., Antonenko, D. S., Chern, S., Zhao, B., Arsene, P., Goldfarb, A., Ivanov, S., Poświata, R., Wang, C., Li, D., Crisostomi, D., Achilleos, A., Myklebust, B., Sen, A., Perrella, D., Kaparov, N., Inlow, M. H., Zang, A., Thornley, E., Orel, D., Poritski, V., Ben-David, S., Berger, Z., Whitfill, P., Foster, M., Munro, D., Ho, L., Hava, D. B., Kuchkin, A., Lauff, R., Holmes, D., Sommerhage, F., Schneider, K., Kazibwe, Z., Stambaugh, N., Singh, M., Magoulas, I., Clarke, D., Kim, D. H., Dias, F. M., Elser, V., Agarwal, K. P., Vilchis, V. E. G., Klose, I., Demian, C., Anantheswaran, U., Zweiger, A., Albani, G., Li, J., Daans, N., Radionov, M., Rozhoň, V., Ma, Z., Stump,

C., Berkani, M., Platnick, J., Nevirkovets, V., Basler, L., Piccardo, M., Jeanplong, F., Cohen, N., Tkadlec, J., Rosu, P., Padlewski, P., Barzowski, S., Montgomery, K., Menezes, A., Patel, A., Wang, Z., Tucker-Foltz, J., Stade, J., Goertzen, T., Kazemi, F., Milbauer, J., Ambay, J. A., Shukla, A., Labrador, Y. C. L., Givré, A., Wolff, H., Rossbach, V., Aziz, M. F., Kaddar, Y., Chen, Y., Zhang, R., Pan, J., Terpin, A., Muennighoff, N., Schoelkopf, H., Zheng, E., Carmi, A., Jones, A., Shah, J., Brown, E. D. L., Zhu, K., Bartolo, M., Wheeler, R., Ho, A., Barkan, S., Wang, J., Stehberger, M., Kretov, E., Sridhar, K., EL-Wasif, Z., Zhang, A., Pyda, D., Tam, J., Cunningham, D. M., Goryachev, V., Patramanis, D., Krause, M., Redenti, A., Bugas, D., Aldous, D., Lai, J., Coleman, S., Bahaloo, M., Xu, J., Lee, S., Zhao, S., Tang, N., Cohen, M. K., Carroll, M., Paradise, O., Kirchner, J. H., Steinerberger, S., Ovchynnikov, M., Matos, J. O., Shenoy, A., de Oliveira Junior, B. A., Wang, M., Nie, Y., Giordano, P., Petersen, P., Sztyber-Betley, A., Shukla, P., Crozier, J., Pinto, A., Verma, S., Joshi, P., Yong, Z.-X., Tee, A., Andréoletti, J., Weller, O., Singhal, R., Zhang, G., Ivanov, A., Khoury, S., Mostaghimi, H., Thaman, K., Chen, Q., Khánh, T. Q., Loader, J., Cavalleri, S., Szlyk, H., Brown, Z., Roberts, J., Alley, W., Sun, K., Stendall, R., Lamparth, M., Reuel, A., Wang, T., Xu, H., Raparthi, S. G., Hernández-Cámara, P., Martin, F., Malishev, D., Preu, T., Korbak, T., Abramovitch, M., Williamson, D., Chen, Z., Bálint, B., Bari, M. S., Kassani, P., Wang, Z., Ansarinejad, B., Goswami, L. P., Sun, Y., Elgnainy, H., Tordera, D., Balabanian, G., Anderson, E., Kvistad, L., Moyano, A. J., Maheshwari, R., Sakor, A., Eron, M., McAlister, I. C., Gimenez, J., Enyekwe, I., O., A. F. D., Shah, S., Zhou, X., Kamalov, F., Clark, R., Abdoli, S., Santens, T., Meer, K., Wang, H. K., Ramakrishnan, K., Chen, E., Tomasiello, A., Luca, G. B. D., Looi, S.-Z., Le, V.-K., Kolt, N., Mündler, N., Semler, A., Rodman, E., Drori, J., Fossum, C. J., Jagota, M., Pradeep, R., Fan, H., Shah, T., Eicher, J., Chen, M., Thaman, K., Merrill, W., Harris, C., Gross, J., Gusev, I., Sharma, A., Agnihotri, S., Zhelnov, P., Usawasutsakorn, S., Mofayezi, M., Bogdanov, S., Piperski, A., Carauleanu, M., Zhang, D. K., Ler, D., Leventov, R., Soroko, I., Jansen, T., Lauer, P., Duersch, J., Taamazyan, V., Morak, W., Ma, W., Held, W., Huy, T. D., Xian, R., Zebaze, A. R., Mohamed, M., Leser, J. N., Yuan, M. X., Yacar, L., Lengler, J., Shahrtash, H., Oliveira, E., Jackson, J. W., Gonzalez, D. E., Zou, A., Chidambaram, M., Manik, T., Haffenden, H., Stander, D., Dasouqi, A., Shen, A., Duc, E., Golshani, B., Stap, D., Uzhou, M., Zhidkovskaya, A. B., Lewark, L., Vincze, M., Wehr, D., Tang, C., Hossain, Z., Phillips, S., Muzhen, J., Ekström, F., Hammon, A., Patel, O., Remy, N., Farhidi, F., Medley, G., Mohammadzadeh, F., Peñaflor, M., Kassahun, H., Friedrich, A., Sparrow, C., Sakal, T., Dhamane, O., Mirabadi, A. K., Hallman, E., Battaglia, M., Maghsoudimehrabani, M., Hoang, H., Amit, A., Hulbert, D., Pereira, R., Weber, S., Mensah, S., Andre, N., Peristyy, A., Harjadi, C., Gupta, H., Malina, S., Albanie, S., Cai, W., Mehkary, M., Reidegeld, F., Dick, A.-K., Friday, C., Sidhu, J., Kim, W., Costa, M., Gurdogan, H., Weber, B., Kumar, H., Jiang, T., Agarwal, A., Ceconello, C., Vaz, W. S., Zhuang, C., Park, H., Tawfeek, A. R., Aggarwal, D., Kirchhof, M., Dai, L., Kim, E., Ferret, J., Wang, Y., Yan, M., Burdzy, K., Zhang, L., Franca, A., Pham, D. T., Loh, K. Y., Robinson, J., Gul, S., Chhablani, G., Du, Z., Cosma, A., White, C., Riblet, R., Saxena, P., Votava, J., Vinnikov, V., Delaney, E., Halasyamani, S., Shahid, S. M., Mourrat, J.-C., Vetoshkin, L., Bacho, R., Ginis, V., Maksapetyan, A., de la Rosa, F., Li, X., Malod, G., Lang, L., Laurendeau, J., Adesanya, F., Portier, J., Hollom, L., Souza, V., Zhou, Y. A., Yalın, Y., Obikoya, G. D., Arnaboldi, L., Rai, Bigi, F., Bacho, K., Clavier, P., Recchia, G., Popescu, M., Shulga, N., Tanwie, N. M., Lux, T. C. H., Rank, B., Ni, C., Yakimchyk, A., Huanxu, Liu, Häggström, O., Verkama, E., Narayan, H., Gundlach, H., Brito-Santana, L., Amaro, B., Vajipey, V., Grover, R., Fan, Y., e Silva, G. P. R., Xin, L., Kratish, Y., Łucki, J., Li, W.-D., Xu, J., Scaria, K. J., Vargus, F., Habibi, F., Long, Lian, Rodolà, E., Robins, J., Cheng, V., Grabb, D., Bosio, I., Fruhauff, T., Akov, I., Lo, E. J. Y., Qi, H., Jiang, X., Segev, B., Fan, J., Martinson, S., Wang, E. Y., Hausknecht, K., Brenner, M. P., Mao, M., Jiang, Y., Zhang, X., Avagian, D., Scipio, E. J., Siddiqi, M. R., Ragoler, A., Tan, J., Patil, D., Plecnik, R., Kirtland, A., Montecillo, R. G., Durand, S., Bodur, O. F., Adoul, Z., Zekry, M., Douville, G., Karakoc, A., Santos, T. C. B., Shamseldeen, S., Karim, L., Liakhovitskaia, A., Resman, N., Farina, N., Gonzalez, J. C., Maayan, G., Hoback, S., Pena, R. D. O., Sherman, G., Mariji, H., Pouriamanesh, R., Wu, W., Demir, G., Mendoza, S., Alarab, I., Cole, J., Ferreira, D., Johnson, B., Milliron, H., Safdari, M., Dai, L., Arthornthurasuk, S., Pronin, A., Fan, J., Ramirez-Trinidad, A., Cartwright, A., Pottmaier, D., Taheri, O., Outevsky, D., Stepanic, S., Perry, S., Askew, L., Rodríguez, R. A. H., Dendane, A., Ali, S., Lorena, R., Iyer, K., Salauddin, S. M., Islam, M., Gonzalez, J., Ducey, J., Campbell, R., Somrak, M., Mavroudis, V., Vergo, E., Qin, J., Borbás, B., Chu, E., Lindsey, J., Radhakrishnan, A., Jallon, A., McInnis, I. M. J., Hoover, A., Möller, S., Bian, S., Lai, J., Patwardhan, T., Yue, S., Wang, A., and Hendrycks, D. Humanity's last exam, 2025. URL https://arxiv.org/abs/2501.14249.

Phothilimthana, M., Abu-El-Haija, S., Perozzi, B., Reade, W., and Chow, A. Google - Fast or Slow? Predict AI Model Runtime. https://kaggle.com/competitions/predict-ai-model-runtime, 2023. Kaggle.

Rajpurkar, P., Zhang, J., Lopyrev, K., and Liang, P.

SQuAD: 100,000+ questions for machine comprehension of text, 2016. URL https://arxiv.org/abs/1606.05250.

Recht, B., Roelofs, R., Schmidt, L., and Shankar, V. Do ImageNet classifiers generalize to ImageNet? In Chaudhuri, K. and Salakhutdinov, R. (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 5389–5400. PMLR, 09–15 Jun 2019. URL https://proceedings.mlr.press/v97/recht19a.html.

RL, B., Pyrak-Nolte, L., Reade, W., and Howard, A. LANL Earthquake Prediction. https://kaggle.com/competitions/LANL-Earthquake-Prediction, 2019. Kaggle.

Roelofs, R., Fridovich-Keil, S., Miller, J., Shankar, V., Hardt, M., Recht, B., and Schmidt, L. *A meta-analysis of overfitting in machine learning*. Curran Associates Inc., Red Hook, NY, USA, 2019a.

Roelofs, R., Shankar, V., Recht, B., Fridovich-Keil, S., Hardt, M., Miller, J., and Schmidt, L. A meta-analysis of overfitting in machine learning. In Wallach, H., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019b. URL https://proceedings.neurips.cc/paper_files/paper/2019/file/ee39e503b6bedf0c98c388b7e8589aca-Paper.pdf.

Sainz, O., Campos, J., García-Ferrero, I., Etxaniz, J., de Lacalle, O. L., and Agirre, E. NLP evaluation in trouble: On the need to measure LLM data contamination for each benchmark. In Bouamor, H., Pino, J., and Bali, K. (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2023*, pp. 10776–10787, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-emnlp.722. URL https://aclanthology.org/2023.findings-emnlp.722/.

Scale AI. SEAL leaderboards. https://https://scale.com/leaderboard, 2024.

Siemion, A., Alonso, D. D., Reade, W., Wang, S., Croft, S., and Chen, Y. SETI Breakthrough Listen - E.T. Signal Search. https://kaggle.com/competitions/seti-breakthrough-listen, 2021. Kaggle.

Turing, A. M. Computing machinery and intelligence. *Mind*, 59(236):433–460, 1950. ISSN 00264423. URL http://www.jstor.org/stable/2251299.

Vanschoren, J., van Rijn, J. N., Bischl, B., and Torgo, L. Openml: networked science in machine learning. *CoRR*, abs/1407.7722, 2014. URL http://arxiv.org/abs/1407.7722.

Vapnik, V. *The Nature of Statistical Learning Theory*. Springer: New York, 1999.

White, C., Dooley, S., Roberts, M., Pal, A., Feuer, B., Jain, S., Shwartz-Ziv, R., Jain, N., Saifullah, K., Naidu, S., Hegde, C., LeCun, Y., Goldstein, T., Neiswanger, W., and Goldblum, M. Livebench: A challenging, contamination-free llm benchmark. 2024. URL arXivpreprintarXiv:2406.19314.

XTX Investments. AI mathematical olympiad - progress prize 1. https://kaggle.com/competitions/ai-mathematical-olympiad-prize, 2024. Kaggle.

Yin, A., Kleinman, J., Yana, T., Reade, W., and Elliott, J. TalkingData AdTracking Fraud Detection Challenge. https://kaggle.com/competitions/talkingdata-adtracking-fraud-detection, 2018. Kaggle.

Zellers, R., Holtzman, A., Bisk, Y., Farhadi, A., and Choi, Y. Hellaswag: Can a machine really finish your sentence?, 2019. URL https://arxiv.org/abs/1905.07830.

Zheng, L., Chiang, W.-L., Sheng, Y., Zhuang, S., Wu, Z., Zhuang, Y., Lin, Z., Li, Z., Li, D., Xing, E. P., Zhang, H., Gonzalez, J. E., and Stoica, I. Judging LLM-as-a-Judge with MT-bench and Chatbot Arena, 2023. URL https://arxiv.org/abs/2306.05685.

Zhu, B., Frick, E., Wu, T., Zhu, H., and Jiao, J. Nectar. https://huggingface.co/datasets/berkeley-nest/Nectar, 2024.