

---

# LiveBench: A Challenging, Contamination-Free LLM Benchmark

---

Colin White<sup>\*1</sup>, Samuel Dooley<sup>\*1</sup>, Manley Roberts<sup>\*1</sup>, Arka Pal<sup>\*1</sup>, Benjamin Feuer<sup>2</sup>,  
Siddhartha Jain<sup>3</sup>, Ravid Shwartz-Ziv<sup>2</sup>, Neel Jain<sup>4</sup>, Khalid Saifullah<sup>4</sup>, Siddhartha Venkat Naidu<sup>1</sup>,  
Chinmay Hegde<sup>2</sup>, Yann LeCun<sup>2</sup>, Tom Goldstein<sup>4</sup>, Willie Neiswanger<sup>5</sup>, Micah Goldblum<sup>2</sup>

<sup>1</sup> Abacus.AI, <sup>2</sup> NYU, <sup>3</sup> Nvidia, <sup>4</sup> UMD, <sup>5</sup> USC

## Abstract

Test set contamination can quickly render LLM benchmarks obsolete. Test set contamination occurs when a benchmark ends up in a newer model’s training set, a common occurrence given LLM training from web-scale data. Also, benchmarks that crowdsource new prompts and evaluations from human or LLM judges can introduce significant biases, and they break down when scoring hard questions. In this work, we introduce a new benchmark for LLMs designed to be immune to both test set contamination and the pitfalls of LLM judging and human crowdsourcing. We release LiveBench, the first benchmark that (1) contains frequently-updated questions from recent information sources, in which questions become harder over time, (2) scores answers automatically according to objective ground-truth values, and (3) contains a wide variety of challenging tasks, spanning math, coding, reasoning, language, instruction following, and data analysis. To achieve this, LiveBench contains questions that are based on recently-released math competitions, arXiv papers, and datasets, and it contains harder, contamination-free versions of tasks from previously released benchmarks such as BigBench Hard, AMPS, and IFEval. We evaluate several prominent closed-source models, as well as dozens of open-source models ranging from 0.5B to 8x22B in size, on our benchmark. LiveBench is difficult, with top models achieving below 60% accuracy. We release all questions, code, and model answers. Questions will be added and updated on a monthly basis, and we welcome community engagement and collaboration for expanding the benchmark tasks and models.

## 1 Introduction

In recent years, as large language models (LLMs) have risen in prominence, it has become increasingly clear that traditional machine learning benchmark frameworks are no longer sufficient to evaluate new models. Benchmarks are typically published on the internet, and most modern LLMs include large swaths of the internet in their training data. If the LLM has seen the questions of a benchmark during training, its performance on that benchmark will be artificially inflated [14, 15, 20, 38], hence making LLM evaluations fraught. Recent evidence of test set contamination includes the observation that LLMs’ performance on Codeforces plummet after the training cutoff date of the LLM [23, 38], and before the cutoff date, performance is highly correlated with the number of times the problem appears on GitHub [38]. Similarly, a recent hand-crafted variant of the established math dataset, GSM8K, shows evidence that several models have overfit to this benchmark [12, 51].

To lessen dataset contamination, benchmarks using new “live” questions and LLM or human judging have become increasingly popular [10, 23, 29, 52]. However, using these techniques comes with significant downsides. While LLM judges have multiple advantages, such as their speed and ability

---

<sup>\*</sup>Correspondence to: colin@abacus.ai, samuel@abacus.ai, goldblum@nyu.edu.

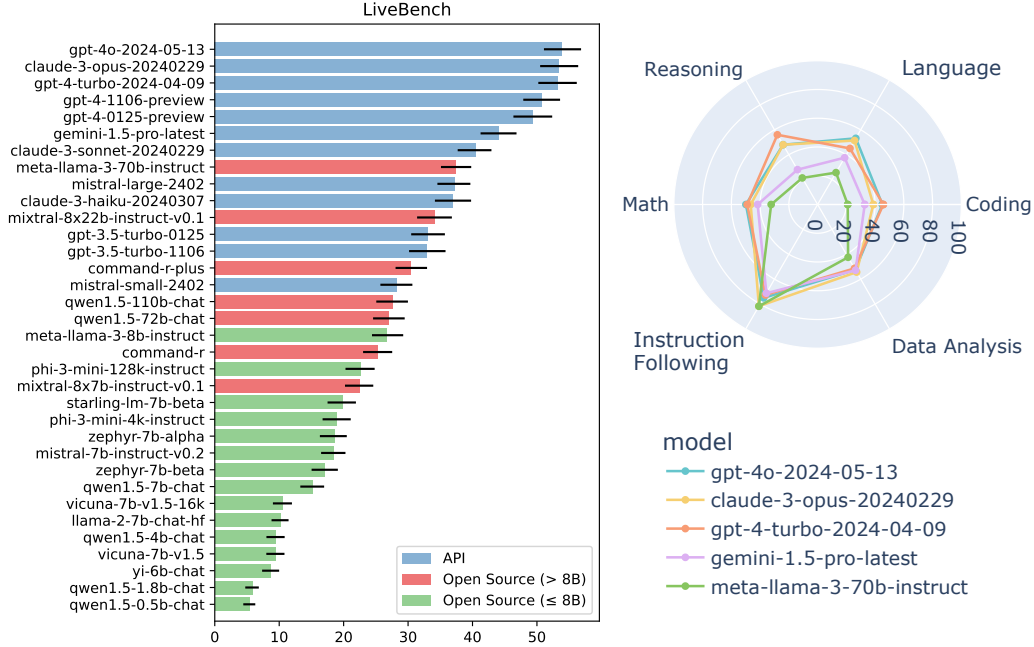


Figure 1: Results on LiveBench for all models, showing 95% bootstrap confidence intervals (left). A radar plot for select models across LiveBench’s six categories demonstrating the that ordering of top models varies between each category (right).

to evaluate open-ended questions, they are prone to make mistakes and can have several biases. For some questions, evaluating a solution is just as hard as producing a solution, resulting in a circular argument. Furthermore, LLMs often favor their own answers over other LLMs [29]. Additionally, using humans to provide evaluations of LLMs can inject biases such as preferences about the length of the output [29], formatting of the output, and the tone and formality of the writing. Using humans to generate questions also presents severe limitations. Human participants might not ask diverse questions, may favor certain topics that do not probe a model’s general capabilities, or may construct their prompts poorly.

We introduce a framework for benchmarking LLMs designed to be immune to both test set contamination and the pitfalls of LLM judging and human crowdsourcing. We use this framework to create LiveBench, the first benchmark with these three desiderata: (1) LiveBench contains frequently-updated questions based on recent information sources, in which questions become harder over time; (2) LiveBench is scored automatically according to the objective ground-truth without the use of an LLM judge; and (3) LiveBench questions are drawn from a diverse set of six categories with regularly released updates of new questions. LiveBench questions are *difficult*; for example, GPT-4-Turbo achieves around 50% accuracy. Questions will be added and updated on a monthly basis, and we will release new tasks and harder versions of tasks over time so that LiveBench can distinguish between the capabilities of LLMs as they improve in the future.

**Categories in LiveBench.** LiveBench currently consists of 18 tasks across 6 categories: reasoning, data analysis, math, coding, language comprehension, and instruction following. Each task falls into one of two types: (1) tasks which use an information source for their questions, e.g., data analysis based on recent Kaggle datasets, or fixing typos in recent arXiv abstracts; and (2) tasks which are more challenging or diverse versions of existing benchmark tasks, e.g., from AMPS [21], Big-Bench Hard [42], IFEval [53], or bAbI [47]. Each category and task included in LiveBench are:

- **Reasoning:** a harder version of Web of Lies from Big-Bench Hard [42], a harder version of positional reasoning from bAbI [47], and Zebra Puzzles (e.g., [24])
- **Data Analysis:** three tasks, all of which use recent datasets from Kaggle and Socrata: table reformatting (among JSON, JSONL, Markdown, CSV, TSV, and HTML), predicting which

columns can be used to join two tables, and predicting the correct type annotation of a data column

- **Math:** questions from high school math competitions from the past 12 months (AMC12, AIME, USAMO, IMO, SMC), as well as harder versions of AMPS [22] questions
- **Coding:** two tasks from Leetcode and AtCoder (via LiveCodeBench [23]): code generation and a novel code completion task
- **Language Comprehension:** three tasks featuring Connections word puzzles, a typo removal task, and a movie synopsis unscrambling task from recent movies on IMDb and Wikipedia
- **Instruction Following:** four tasks to paraphrase, simplify, summarize, or generate stories about recent news articles from The Guardian, subject to one or more instructions such as word limits or incorporating specific elements in the response

We evaluate popular proprietary models as well as dozens of open-source models with sizes ranging from 0.5B to 8x22B. We release all questions, code, and model answers, and we welcome community engagement and collaboration. Our codebase is available at <https://github.com/livebench/livebench>, and our leaderboard is available at <https://livebench.ai>.

## 2 LiveBench Description

In this section, we introduce LiveBench. It currently has six categories: math, coding, reasoning, data analysis, instruction following, and language comprehension. Categories are diverse with two to four tasks per problem. Each task either includes recent information sources (such as very recent news articles, movie synopses, or datasets) or is a more challenging, more diverse version of an existing benchmark task.

Each task is designed to span a range of difficulty, from easy to very challenging, while loosely aiming for 30-70% success on the top models. Prompts are tailored for each category and task but typically include the following: zero-shot chain of thought [26, 46], asking the model to make its best guess if it does not know the answer, and asking the LLM to output its final answer in a way that is easy to parse, such as in `**double asterisks**`. In the following sections, we give an overview description of each task from each category.

### 2.1 Math Category

Evaluating the mathematical abilities of LLMs has been one of the cornerstones of recent research in LLMs, featuring prominently in many releases and reports [6, 7, 34, 37]. Our benchmark includes math questions of three types: questions from recent high school math competitions, fill-in-the-blank questions from recent proof-based USAMO and IMO problems, and questions from our new, harder version of the AMPS dataset [22].

Our first two math tasks, Competitions and Proof Competition, use expert human-designed math problems that offer a wide variety in terms of problem type and solution technique. First, we include questions from AMC12 2023, SMC 2023, AIME 2024 in Competitions and from USAMO 2023 and IMO 2023 in Proof Competitions. These are challenging and prestigious competitions for high school students in the USA (AMC, AIME, USAMO), the UK (SMC), or internationally (IMO). The competitions test mathematical problem solving with arithmetic, algebra, counting, geometry, number theory, probability, and other secondary school math topics [18].

Finally, we release synthetically generated math questions in the AMPS\_Hard task. This task is inspired by the math question generation used to create the MATH and AMPS datasets [22]. We generate harder questions by drawing random primitives, using a larger and more challenging distribution than AMPS across the 10 hardest tasks within AMPS.

### 2.2 Coding Category

The coding ability of LLMs is one of the most widely studied and sought-after skills for LLMs [23, 28, 33]. We include two coding tasks in LiveBench: a modified version of the code generation task from LiveCodeBench (LCB) [23], and a novel code completion task combining LCB problems with partial solutions collected from GitHub sources.

In the LCB Generation task, we assess a model’s ability to parse a competition coding question statement and write a correct answer. We include 50 questions from LiveCodeBench [23] which has several tasks to assess the coding capabilities of large language models.

The Completion task specifically focuses on the ability of models to complete a partially correct solution—assessing whether a model can parse the question, identify the function of the existing code, and determine how to complete it. We use LeetCode medium and hard problems from LiveCodeBench’s [23] April 2024 release, combined with matching solutions from <https://github.com/kamyu104/LeetCode-Solutions>, omitting the last 15% of each solution and asking the LLM to complete the solution.

### 2.3 Reasoning Category

The reasoning ability of large language models is another highly-benchmarked and analyzed skill of LLMs [42, 46, 48]. In LiveBench, we include three reasoning tasks: our harder versions of tasks from Big-Bench Hard [42] and bAbI [47], and Zebra puzzles.

The Web of Lies v2 task is an advancement of the similarly named task included in Big-Bench [5] and Big-Bench Hard [42]. The task is to evaluate the truth value of a random Boolean function expressed as a natural-language word problem. Already by October 2022, LLMs achieved near 100% on this task, and furthermore, there are concerns that Big-Bench tasks leaked into the training data of LLMs such as GPT-4, despite using canary strings [34]. For LiveBench, we create a new, significantly harder version by including additional deductive components and red herrings.

Next, we include a harder version of the PathFinding task from bAbI [47] that we call ‘House Traversal’. The original task consists of sentences of the form, ‘The bedroom is West of the kitchen. The kitchen is South of the garden’, asking the model where a room is in relation to another room. We make this task significantly harder by adding a distinct person in each room and asking the LLM who a person would see if they traverse in a few directions.

The final reasoning task we include is Zebra Puzzles, a well-known [24] reasoning task that tests the ability of the model to follow a set of statements that set up constraints, and then logically deduce the requested information. We build on an existing repository for procedural generation of Zebra puzzles [35]; the repository allows for randomizing the number of people, the number of attributes, and the set of constraint statements provided. Below, we provide an example question from the Zebra Puzzles task.

#### An example question from the Zebra Puzzle task.

There are 3 people standing in a line numbered 1 through 3 in a left to right order.  
Each person has a set of attributes: Food, Nationality, Hobby.  
The attributes have the following possible values:  
- Food: nectarine, garlic, cucumber  
- Nationality: chinese, japanese, thai  
- Hobby: magic-tricks, filmmaking, puzzles  
and exactly one person in the line has a given value for an attribute.  
Given the following premises about the line of people:  
- the person that likes garlic is on the far left  
- the person who is thai is somewhere to the right of the person who likes magic-tricks  
- the person who is chinese is somewhere between the person that likes cucumber and the person who likes puzzles  
Answer the following question: What is the hobby of the person who is thai? Return your answer as a single word, in the following format: **\*\*X\*\***, where X is the answer.

### 2.4 Data Analysis Category

While LLMs’ abilities in the previous three categories have been widely studied, we also include a category that is significantly less studied but is a practical application of LLMs that requires more attention: data analysis tasks. We include three tasks in which the LLM assists in data analysis or data science: column type annotation, table join prediction, and table reformatting.

The first common task a data scientist uses an LLM for is to predict the type of a column of a data table. To create a question for the column table annotation task (CTA), we randomly sample a table and randomly sample a column from that table. We use the actual column name of that column as our ground-truth, and then retrieve some column samples from that column. We provide the name of all the columns from that table and ask the LLM to select the column name from those options.

Data science users often also require a table to be reformatted from one type to another, e.g., json to CSV, XML to TSV. We emulate that task in TableReformat by providing a table in one format and asking the LLM to reformat it into the target format.

Another common application of LLMs in data science is to perform table joins. So in the Join task, we ask an LLM to do just that. Each question prompts an LLM to join two tables by a given column name, where each table is given as a CSV in text form.

## 2.5 Instruction Following Category

An important desideratum of an LLM is its capability to following instructions. To this end, we include instruction-following questions in our benchmark, inspired by IFEval [53] which is an instruction-following evaluation for LLMs containing verifiable instructions such as “write more than 300 words” or “mention the keyword <X> at least 5 times.”. In contrast to IFEval, which presents only the task and instructions with a simple prompt like “write a travel blog about Japan”, in LiveBench the models are provided an article from the Guardian newspaper and must complete the task related to the document in line with the constraints provided by the instructions. We find that by including a document that is referenced, the difficulty of the task increases substantially. We develop four tasks Paraphrase, Simplify, Story Generation, and Summarize. Each of these tasks requires performing what their title suggests e.g. ‘Paraphrase’ requires paraphrasing the document, ‘Summarize’ requires summarizing it, etc. These tasks must be completed while adhering to a set of instructions sampled from Table 4. Care is taken to ensure the instructions are not self-contradictory.

## 2.6 Language Comprehension Category

Finally, we include multiple language comprehension tasks. These tasks assess the language model’s ability to reason about language itself by, (1) completing word puzzles, (2) fixing misspellings but leaving other stylistic changes in place, and (3) reordering scrambled plots of unknown movies.

First, we include the Connections category. Connections is a word puzzle popularized by the New York Times (although similar ideas have existed previously). In this task, we present questions of varying levels of difficulty with 8, 12, and 16-word varieties. The objective of the game is to sort the words into sets of four words, such that each set has a ‘connection’ between them, e.g., types of fruits, homophones, or words that come after the word ‘fire’. Due to the variety of possible connection types, this task is challenging for LLMs, as shown by prior work [43] which tested the task on the GPT family of models.

Next, we include details about the Typos task. The idea behind this task is inspired by the common use case for LLMs where a user will ask the system to identify typos and misspellings in some written text but to leave other aspects of the text unchanged. It is common for the LLM to impose its own writing style onto that of the input text, which is not desirable. Examples of these undesirable changes include switching from US to British spellings or toggling between using or not using the serial comma. We create the questions for this task from recent ArXiv abstracts, which we ensure have no typos, by programmatically injecting common human typos into the text. Below is an example question from the Typos task.

### An example question from the Typos task.

Please output this exact text, with no changes at all except for fixing the misspellings. Please leave all other stylistic decisions like commas and US vs British spellings as in the original text.

We introduce a Bayesian estimation approach for the passive localization of an acoustic source in shallow water using a single mobile receiver. The proposed probabilistic focalization method estimates the time-varying source location in the presence of measurement-origin

Table 1: **LiveBench results across the 15 top-performing models.** We display in this table the highest-performing models on LiveBench, outputting the results on each main category, as well as each model’s overall performance. See Table 3 for the results on all 34 models.

Model	LiveBench Score	Coding	Data Analysis	Instruction Following	Math	Reasoning	Language
gpt-4o-2024-05-13	<b>53.9</b>	45.1	52.4	75.0	<b>49.9</b>	48.0	<b>53.2</b>
claude-3-opus-20240229	53.4	38.7	<b>54.3</b>	81.7	46.5	48.0	51.4
gpt-4-turbo-2024-04-09	53.3	<b>45.7</b>	51.3	72.7	49.0	<b>56.0</b>	45.0
gpt-4-1106-preview	50.8	43.1	51.3	72.9	47.6	42.7	47.4
gpt-4-0125-preview	49.3	42.7	54.1	74.2	42.7	40.7	41.5
gemini-1.5-pro-latest	44.0	32.8	52.8	71.3	41.7	28.0	37.6
claude-3-sonnet-20240229	40.6	23.9	44.6	73.9	29.6	33.3	38.1
meta-llama-3-70b-instruct	37.4	20.9	42.4	<b>81.8</b>	32.3	21.3	25.6
mistral-large-2402	37.2	16.3	44.6	63.7	31.9	38.0	28.7
claude-3-haiku-20240307	37.0	24.5	41.5	70.6	25.7	30.7	28.9
mixtral-8x22b-instruct-v0.1	34.1	31.8	31.7	61.3	26.9	26.7	26.3
gpt-3.5-turbo-0125	33.0	29.2	41.2	54.4	25.5	24.0	23.9
gpt-3.5-turbo-1106	32.9	26.8	41.7	49.3	27.8	24.0	27.8
command-r-plus	30.4	15.0	24.6	69.2	24.9	25.3	23.8
mistral-small-2402	28.3	16.3	31.9	53.9	26.8	20.0	20.9

uncertainty. In particular, probabilistic data association is performed to match time-differences-of-arrival (TDOA) observations extracted from the acoustic signal to TDOA predictions provided by the statistical model. The performance of our approach is evaluated using real acoustic data recorded by a single mobile receiver.

Finally, we include the Plot Unscrambling task, which takes the plot synopses of current movies from IMDb or Wikipedia. We randomly shuffle the synopsis sentences and then ask the LLM to simply reorder the sentences into the original plot. We find that this task is very challenging for LLMs, as it measures their abilities to reason through plausible sequences of events.

### 3 Experiments

In this section, first we describe our experimental setup and present full results for 34 LLMs on all 18 tasks of LiveBench. Next, we give an empirical comparison of LiveBench to existing prominent LLM benchmarks, and finally, we present ablation studies.

**Experimental setup.** Our experiments include 34 LLMs total, with a mix of top proprietary models, large open-source models, and small open-source models. In particular, for proprietary models, we include six GPT models: gpt-4o-2024-05-13, gpt-4-turbo-2024-04-09, gpt-4-1106-preview, gpt-4-0125-preview, gpt-3.5-turbo-1106, gpt-3.5-turbo-0125 [6, 34]; three Anthropic models: claude-3-opus-20240229, claude-3-sonnet-20240229, claude-3-haiku-20240307 [2]; two Mistral models: mistral-large-2402, mistral-small-2402 [25]; and gemini-1.5-pro-latest (the API version) [37].

For large open-source models, we include command-r, command-r-plus [13], meta-llama-3-70b-instruct [32], mixtral-8x22b-instruct-v0.1, mixtral-8x7b-instruct-v0.1 [25], qwen1.5-110b-chat, and qwen1.5-72b-chat [3].

For small open-source models, we include llama-2-7b-chat-hf [44], llama-3-8b-instruct [32], mistral-7b-instruct-v0.2 [25], phi-3-mini-128k-instruct, phi-3-mini-4k-instruct [1], qwen1.5-0.5b-chat, qwen1.5-1.8b-chat, qwen1.5-4b-chat, qwen1.5-7b-chat [3], starling-lm-7b-beta [54], vicuna-7b-v1.5, vicuna-7b-v1.5-16k [9], yi-6b-chat [49], zephyr-7b-alpha, and zephyr-7b-beta [45].

For all models and tasks, we perform single-turn evaluation with temperature 0. All models run with their respective templates from FastChat [52]. We run all open-source models with bfloat16. For each question, a model receives a score from 0 to 1. For each model, we compute the score on each



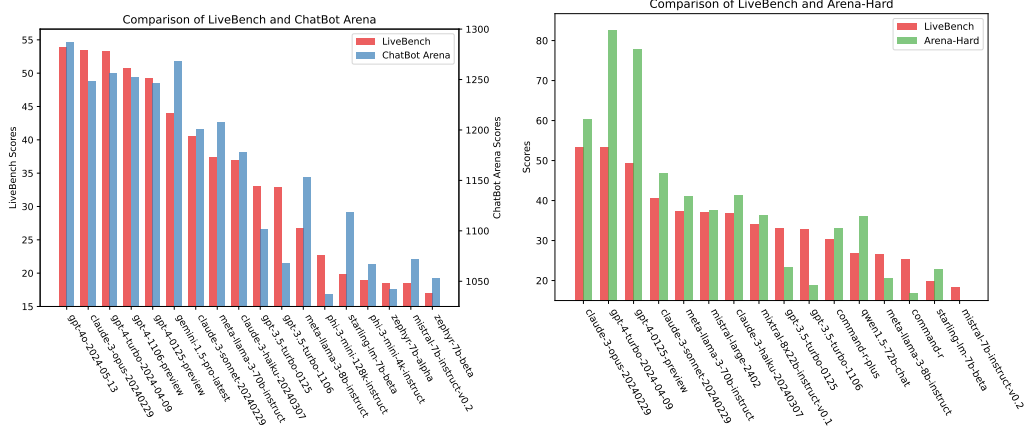


Figure 2: **Comparison of LiveBench to other LLM benchmarks.** We compare LiveBench to ChatBot Arena (left) and Arena-Hard (right). We see that while there are generally similar trends, some models are noticeably stronger on one benchmark vs. the other indicating some of the downfalls of LLM judging.

task as the average of all questions, we compute the score on each of the six categories as the average of all their tasks, and we compute the final LiveBench score as the average of all six categories.

### 3.1 Discussion of Results

We compare all 34 models on LiveBench according to the experimental setup described above; see Table 3. We find that gpt-4o-2024-05-13 performs the best overall, in a virtual tie with claude-3-opus-20240229 and gpt-4-turbo-2024-04-09. The best-performing open-source model is meta-llama-3-70b-instruct, and meta-llama-3-8b-instruct is the best-performing open-source model that is 8B or smaller. Interestingly, although the top three models nearly tie, each one performs significantly better or worse than the others in certain categories.

We find that the GPT-4 models, particularly gpt-4-turbo-2024-04-09, are the highest-performing models in coding, which is in line with recent existing results [23, 40]. gpt-4-turbo-2024-04-09 is similarly highest-performing in reasoning.

Surprisingly, in the instruction-following category, meta-llama-3-70b-instruct achieves the top performance, closely followed by claude-3-opus-20240229, with GPT models significantly behind. To obtain additional validation for this result, we ran two additional tests (see Appendix B.1 for the full details): we ran claude-3-opus-20240229 and gpt-4-turbo-2024-04-09 on IFEval [53], finding that the former outperformed the latter across all categories. We also looked at examples in which the former outperformed GPT models. We find that GPT models have trouble with challenging instructions such as, ‘In your response, the letter i should appear less than n times.’; see Table 5. We also note that since LiveBench does not use few-shot prompting – only giving specific instructions – as with other zero-shot benchmarks, at least a partial component of scoring well on all the tasks requires good instruction-following.

### 3.2 Comparison to Other LLM Benchmarks

Next, we compare LiveBench to two prominent benchmarks, ChatBot Arena [10] and Arena-Hard [29]. In Figure 2, we show a bar plot comparison among models that are common to both benchmarks, and in Figure 3, we compare the performance of these models to a best-fit line. We also compute the correlation coefficient of model scores among the benchmarks: LiveBench has a 0.92 and 0.90 correlation with ChatBot Arena and Arena-Hard, respectively.

Based on the plots and the correlation coefficients, we see that there are generally similar trends to LiveBench, yet some models are noticeably stronger on one benchmark vs. the other. For example, gpt-4-0125-preview and gpt-4-turbo-2024-04-09 perform substantially better on Arena-Hard compared to LiveBench – likely due to the known bias from using gpt-4 itself as the LLM judge [29].

Table 2: **LLM judges cannot accurately evaluate challenging math and reasoning questions.** Pearson’s correlation coefficient of objective ground truth scoring vs. LLM-as-a-judge scoring, on challenging math (AMC, AIME, SMC) and reasoning (Zebra puzzles) tasks. The judge is gpt-4-turbo-2024-04-09. If the LLM judge was highly accurate, we would expect a correlation close to 1. However, on all tasks, the correlation is surprisingly low, and sometimes close to 0, showing that LLMs are not reliable judges for these tasks.

Model	AMC12 2024	AIME 2024	SMC 2023	Zebra Puzzles
GPT-4-Turbo	0.227	0.548	0.247	0.272
Claude-3-Opus	0.25	0.596	0.408	0.098

We hypothesize that the strong performance of some models such as gemini-1.5-pro-latest and starling-lm-7b-beta on ChatBot Arena compared to LiveBench may be due to having an output style that is preferred by humans. These observations emphasize the benefit of using ground-truth judging, which is immune to biases based on the style of the output.

### 3.3 Comparison between Ground-Truth and LLM-Judging

In this section, we run an ablation study to compare the result of ground-truth judging with LLM judging, by taking three math sub-tasks and one reasoning task and scoring them by either matching with the ground-truth answer or by asking an LLM judge to score the answer as either correct or incorrect. We use a judge prompt based on the MT-Bench judge prompt (see Appendix B.2 for details), and we use gpt-4-turbo-2024-04-09 as the judge. We judge the model outputs of both gpt-4-turbo-2024-04-09 and claude-3-opus-20240229 in Table 2. If the LLM judge was highly accurate, we would expect a Pearson’s correlation coefficient  $\rho$  to be close to 1. However, we find that  $\rho$  for all tasks is far below a reasonable value, indicating that LLM judges are not appropriate for challenging math and logic tasks. Interestingly, the highest correlations are on AIME 2024, which is also the task with the lowest overall success rate according to ground-truth judgment.

## 4 Related Work

We describe the most prominent LLM benchmarks and the ones that are most related to our work. For a comprehensive survey, see [8]. The Huggingface Open LLM Leaderboard [4, 19] is a widely-used benchmark suite that consists of six static datasets: ARC [11], GSM8K [12], HellaSwag [50], MMLU [21], TruthfulQA [31], and Winogrande [39]. While this has been incredibly useful in tracking the performance of LLMs, its static nature has left it prone to test set contamination by models.

**LLMs-as-a-judge.** AlpacaEval [16, 17, 30], MT-Bench [10], and Arena-Hard [29] are benchmarks that employ LLM judges on a fixed set of questions. Using an LLM-as-a-judge is fast and relatively cheap. Furthermore, this strategy has the flexibility of being able to evaluate open-ended questions, instruction-following questions, and chatbots. However, LLM judging also has downsides. First, LLMs have biases towards their own answers [29]. In addition to favoring their own answers [29], GPT-4 judges have a noticeable difference in terms of variance and favorability of other models compared to Claude judges. Additionally, LLMs make errors. As one concrete example, question 2 in Arena-Hard asks a model to write a C++ program to compute whether a given string can be converted to ‘abc’ by swapping two letters. GPT-4 incorrectly judges gpt-4-0314’s solution as incorrect [29].

**Humans-as-a-judge.** ChatBot Arena [10, 52] leverages human prompting and feedback on a large scale. Users ask questions and receive outputs of two randomly selected models and have to pick which output they prefer. This preference feedback is aggregated into Elo scores for the different models. While human evaluation is great for capturing the preferences of a crowd, using a human-as-a-judge has many disadvantages. First, human-judging can be quite labor-intensive, especially for certain tasks included in LiveBench such as complex math, coding, or long-context reasoning problems. Whenever humans are involved in annotation (of which judging is a sub-case), design choices or factors can cause high error rates [27], and even in well-designed human-annotation setups, high variability from human to human leads to unpredictable outcomes [36].

**Other benchmarks** Perhaps the most-related benchmark to ours is LiveCodeBench [23], which also regularly releases new questions and makes use of ground-truth judging. However, it is limited to only coding tasks. Concurrent work, the SEAL Benchmark [40], uses private questions with expert human



scorers, however, the benchmark currently only contains the following categories: Math, Coding, Instruction Following, and Spanish. In [41], the authors modify the original MATH dataset [22] by changing numbers in the problem setup. They find drastic declines in model performance for all LLMs including the frontier ones. However, while such work can evaluate LLMs on data that is not in the pretraining set, the data still ends up being highly similar to the kind of data likely seen in the pretraining set. In addition, the hardness of the benchmark remains the same over time.

## 5 Conclusions, Limitations, and Future Work

In this work, we introduced LiveBench, an LLM benchmark designed to mitigate both test set contamination and the pitfalls of LLM judging and human crowdsourcing. LiveBench is the first benchmark that (1) contains frequently updated questions from new information sources, in which questions become harder over time, (2) scores answers automatically according to objective ground-truth values, without the use of LLM judges, and (3) contains a wide variety of challenging tasks, spanning math, coding, reasoning, writing, instruction following, and data analysis. LiveBench contains questions that are based on recently released math competitions, arXiv papers, and datasets, and it contains harder, ‘contamination-proof’ versions of previously released benchmarks. We released all questions, code, and model answers, and questions will be added and updated on a monthly basis. We welcome community engagement and collaboration for expanding the benchmark tasks and models.

**Limitations and Future Work.** While we attempted to make LiveBench as diverse as possible, there are still additions from which it would benefit. For example, we hope to add non-English language tasks in the future. Furthermore, while ground truth scoring is beneficial in many ways, it still cannot be used for certain use cases, such as ‘write an email to my boss’, or ‘write a travel guide to Hawaii’ in which it is hard to define a ground truth. Finally, while we attempted to make all tasks and categories fair for all models, there are still biases due to certain LLM families favoring certain prompt types. We plan to update the prompts (at the start and end of each question) in the future, as new prompt strategies are developed. Similarly, plan to update the LiveBench leaderboard as new LLMs are released.

## References

- [1] Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat Behl, et al. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*, 2024.
- [2] AI Anthropic. The claude 3 model family: Opus, sonnet, haiku. *Claude-3 Model Card*, 2024.
- [3] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
- [4] Edward Beeching, Cl  mentine Fourrier, Nathan Habib, Sheon Han, Nathan Lambert, Nazneen Rajani, Omar Sanseviero, Lewis Tunstall, and Thomas Wolf. Open llm leaderboard. [https://huggingface.co/spaces/HuggingFaceH4/open\\_llm\\_leaderboard](https://huggingface.co/spaces/HuggingFaceH4/open_llm_leaderboard), 2023.
- [5] BIG bench authors. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *Transactions on Machine Learning Research*, 2023.
- [6] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [7] S  bastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*, 2023.
- [8] Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, et al. A survey on evaluation of large language models. *ACM Transactions on Intelligent Systems and Technology*, 15(3):1–45, 2024.
- [9] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. Vicuna: An open-source chatbot impressing gpt-4 with 90%\* chatgpt quality. See <https://vicuna.lmsys.org> (accessed 14 April 2023), 2(3):6, 2023.
- [10] Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li, Dacheng Li, Hao Zhang, Banghua Zhu, Michael Jordan, Joseph E Gonzalez, et al. Chatbot arena: An open platform for evaluating llms by human preference. *arXiv preprint arXiv:2403.04132*, 2024.
- [11] Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge, 2018.
- [12] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- [13] Cohere. Command r: Retrieval-augmented generation at production scale. <https://txt.cohere.com/command-r>, March 2024.
- [14] Chunyuan Deng, Yilun Zhao, Xiangru Tang, Mark Gerstein, and Arman Cohan. Investigating data contamination in modern benchmarks for large language models. *arXiv preprint arXiv:2311.09783*, 2023.
- [15] Yihong Dong, Xue Jiang, Huanyu Liu, Zhi Jin, and Ge Li. Generalization or memorization: Data contamination and trustworthy evaluation for large language models. *arXiv preprint arXiv:2402.15938*, 2024.
- [16] Yann Dubois, Bal  zs Galambosi, Percy Liang, and Tatsunori B Hashimoto. Length-controlled alpacaeval: A simple way to debias automatic evaluators. *arXiv preprint arXiv:2404.04475*, 2024.

- [17] Yann Dubois, Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. AlpacaFarm: A simulation framework for methods that learn from human feedback, 2023.
- [18] J Douglas Faires and David Wells. *The Contest Problem Book VIII: American Mathematics Competitions (AMC 10) 2000–2007*, volume 19. American Mathematical Society, 2022.
- [19] Leo Gao, Jonathan Tow, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Kyle McDonell, Niklas Muennighoff, Jason Phang, Laria Reynolds, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. A framework for few-shot language model evaluation, September 2021.
- [20] Shahriar Golchin and Mihai Surdeanu. Time travel in llms: Tracing data contamination in large language models. *arXiv preprint arXiv:2308.08493*, 2023.
- [21] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*, 2020.
- [22] Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*, 2021.
- [23] Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free evaluation of large language models for code. *arXiv preprint arXiv:2403.07974*, 2024.
- [24] S Jeremy. Einstein’s riddle: Riddles, paradoxes, and conundrums to stretch your mind. *Bloomsbury USA*, pages 10–11, 2009.
- [25] Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- [26] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. *Proceedings of the Annual Conference on Neural Information Processing Systems (NeurIPS)*, 2022.
- [27] Matthew Lease. On quality control and machine learning in crowdsourcing. In *Workshops at the twenty-fifth AAAI conference on artificial intelligence*, 2011.
- [28] Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou, Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, et al. Starcoder: may the source be with you! *arXiv preprint arXiv:2305.06161*, 2023.
- [29] Tianle Li, Wei-Lin Chiang, Evan Frick, Lisa Dunlap, Banghua Zhu, Joseph E. Gonzalez, and Ion Stoica. From live data to high-quality benchmarks: The arena-hard pipeline, April 2024.
- [30] Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. AlpacaEval: An automatic evaluator of instruction-following models. [https://github.com/tatsu-lab/alpaca\\_eval](https://github.com/tatsu-lab/alpaca_eval), 2023.
- [31] Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic human falsehoods. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2022.
- [32] Meta. Introducing meta llama 3: The most capable openly available llm to date. <https://ai.meta.com/blog/meta-llama-3/>, April 2024. Accessed: June 4, 2024.
- [33] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Bellemare, Alex Graves, Martin Riedmiller, Andreas K. Fidjeland, Georg Ostrovski, Stig Petersen, Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharshan Kumaran, Daan Wierstra, Shane Legg, and Demis Hassabis. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, February 2015.

- [34] OpenAI. Gpt-4 technical report. *Technical Report*, 2023.
- [35] quint t. Puzzle generator and puzzle solver. <https://github.com/quint-t/Puzzle-Generator-and-Solver>, 2023.
- [36] Hannah Rashkin, Vitaly Nikolaev, Matthew Lamm, Lora Aroyo, Michael Collins, Dipanjan Das, Slav Petrov, Gaurav Singh Tomar, Iulia Turc, and David Reitter. Measuring attribution in natural language generation models. *Computational Linguistics*, 49(4):777–840, 2023.
- [37] Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024.
- [38] Manley Roberts, Himanshu Thakur, Christine Herlihy, Colin White, and Samuel Dooley. To the cutoff... and beyond? a longitudinal perspective on llm data contamination. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2024.
- [39] Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial winograd schema challenge at scale. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(05):8732–8740, Apr. 2020.
- [40] Scale AI. Seal leaderboards. <https://scale.com/leaderboard>, May 2024.
- [41] Saurabh Srivastava, Anto PV, Shashank Menon, Ajay Sukumar, Alan Philipose, Stevin Prince, Sooraj Thomas, et al. Functional benchmarks for robust evaluation of reasoning performance, and the reasoning gap. *arXiv preprint arXiv:2402.19450*, 2024.
- [42] Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc Le, Ed Chi, Denny Zhou, et al. Challenging big-bench tasks and whether chain-of-thought can solve them. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 13003–13051, 2023.
- [43] Graham Todd, Tim Merino, Sam Earle, and Julian Togelius. Missed connections: Lateral thinking puzzles for large language models, 2024.
- [44] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- [45] Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada, Shengyi Huang, Leandro von Werra, Clémentine Fourier, Nathan Habib, Nathan Sarrazin, Omar Sanseviero, Alexander M. Rush, and Thomas Wolf. Zephyr: Direct distillation of lm alignment, 2023.
- [46] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Proceedings of the Annual Conference on Neural Information Processing Systems (NeurIPS)*, 2022.
- [47] Jason Weston, Antoine Bordes, Sumit Chopra, Alexander M Rush, Bart Van Merriënboer, Armand Joulin, and Tomas Mikolov. Towards ai-complete question answering: A set of prerequisite toy tasks. In *4th International Conference on Learning Representations, ICLR 2016*, 2016.
- [48] Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *Proceedings of the Annual Conference on Neural Information Processing Systems (NeurIPS)*, 36, 2024.
- [49] Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng Zhu, Jianqun Chen, Jing Chang, et al. Yi: Open foundation models by 01. ai. *arXiv preprint arXiv:2403.04652*, 2024.

- [50] Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. HellaSwag: Can a machine really finish your sentence? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4791–4800, Florence, Italy, July 2019. Association for Computational Linguistics.
- [51] Hugh Zhang, Jeff Da, Dean Lee, Vaughn Robinson, Catherine Wu, Will Song, Tiffany Zhao, Pranav Raja, Dylan Slack, Qin Lyu, et al. A careful examination of large language model performance on grade school arithmetic. *arXiv preprint arXiv:2405.00332*, 2024.
- [52] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36, 2024.
- [53] Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and Le Hou. Instruction-following evaluation for large language models. *arXiv preprint arXiv:2311.07911*, 2023.
- [54] Banghua Zhu, Evan Frick, Tianhao Wu, Hanlin Zhu, and Jiantao Jiao. Starling-7b: Improving llm helpfulness & harmlessness with rlaiif, November 2023.

## Checklist

1. For all authors...
  - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? [Yes] Our abstract and introduction both state the main claims in the paper.
  - (b) Did you describe the limitations of your work? [Yes] We describe the limitations of our work in [Section 5](#).
  - (c) Did you discuss any potential negative societal impacts of your work? [Yes] We describe the broader societal impacts in [Appendix A](#).
  - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes] We have reviewed the code of ethics and made sure to adhere to it.
2. If you are including theoretical results...
  - (a) Did you state the full set of assumptions of all theoretical results? [N/A] Our work does not include theoretical results.
  - (b) Did you include complete proofs of all theoretical results? [N/A] Our work does not include theoretical results.
3. If you ran experiments (e.g. for benchmarks)...
  - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] We include all code, data, and instructions needed to run and score our benchmark at <https://livebench.ai>. The only part that is not public is the code used to generate questions, since this would raise the issue of test data contamination for future LLM.
  - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [N/A] Our work does not involve training models.
  - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes] We include confidence intervals in [Figure 1](#).
  - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes]
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
  - (a) If your work uses existing assets, did you cite the creators? [Yes] We cite the creators of all assets that we use, in [Section 2](#) and [Section 3](#).
  - (b) Did you mention the license of the assets? [Yes]
  - (c) Did you include any new assets either in the supplemental material or as a URL? [Yes] We include our full benchmark at <https://livebench.ai>
  - (d) Did you discuss whether and how consent was obtained from people whose data you’re using/curating? [N/A] We do not use/curate people’s data.
  - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A] We do not use/curate any data that contains personally identifiable information or offensive content.
5. If you used crowdsourcing or conducted research with human subjects...
  - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A] We do not crowdsource or conduct research with human subjects.
  - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A] We do not crowdsource or conduct research with human subjects.
  - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A] We do not crowdsource or conduct research with human subjects.

## A Broader Societal Impact Statement

Our paper introduces a new benchmark for LLMs, which contains frequently-updated questions from new information sources, scores answers according to objective ground-truth values, and contains a wide variety of tasks. We do not see any inherently negative broader societal impacts of our work.



Table 3: LiveBench **Results across 34 models**. We run 14 proprietary and 20 open-source models on LiveBench, outputting the results on each main category, as well as each model’s overall performance.

Model	LiveBench Score	Coding	Data Analysis	Instruction Following	Math	Reasoning	Language
gpt-4o-2024-05-13	<b>53.9</b>	45.1	52.4	75.0	<b>49.9</b>	48.0	<b>53.2</b>
claude-3-opus-20240229	53.4	38.7	<b>54.3</b>	81.7	46.5	48.0	51.4
gpt-4-turbo-2024-04-09	53.3	<b>45.7</b>	51.3	72.7	49.0	<b>56.0</b>	45.0
gpt-4-1106-preview	50.8	43.1	51.3	72.9	47.6	42.7	47.4
gpt-4-0125-preview	49.3	42.7	54.1	74.2	42.7	40.7	41.5
gemini-1.5-pro-latest	44.0	32.8	52.8	71.3	41.7	28.0	37.6
claude-3-sonnet-20240229	40.6	23.9	44.6	73.9	29.6	33.3	38.1
meta-llama-3-70b-instruct	37.4	20.9	42.4	<b>81.8</b>	32.3	21.3	25.6
mistral-large-2402	37.2	16.3	44.6	63.7	31.9	38.0	28.7
claude-3-haiku-20240307	37.0	24.5	41.5	70.6	25.7	30.7	28.9
mistral-8x22b-instruct-v0.1	34.1	31.8	31.7	61.3	26.9	26.7	26.3
gpt-3.5-turbo-0125	33.0	29.2	41.2	54.4	25.5	24.0	23.9
gpt-3.5-turbo-1106	32.9	26.8	41.7	49.3	27.8	24.0	27.8
command-r-plus	30.4	15.0	24.6	69.2	24.9	25.3	23.8
mistral-small-2402	28.3	16.3	31.9	53.9	26.8	20.0	20.9
qwen1.5-110b-chat	27.6	22.2	31.5	53.9	25.2	22.7	10.2
qwen1.5-72b-chat	26.9	22.9	33.0	51.2	26.8	16.7	11.0
meta-llama-3-8b-instruct	26.7	18.3	23.3	64.5	17.6	18.0	18.7
command-r	25.3	12.3	31.7	56.5	16.9	20.7	13.6
phi-3-mini-128k-instruct	22.7	11.6	26.7	47.0	21.5	24.0	5.5
mistral-8x7b-instruct-v0.1	22.5	10.0	26.8	47.3	19.0	18.7	13.0
starling-lm-7b-beta	19.8	18.3	22.0	34.4	13.8	23.3	7.0
phi-3-mini-4k-instruct	19.0	14.9	14.7	37.6	19.9	22.0	4.8
zephyr-7b-alpha	18.5	11.3	22.1	47.0	9.6	14.7	6.5
mistral-7b-instruct-v0.2	18.5	11.6	20.6	42.3	16.0	11.3	9.1
zephyr-7b-beta	17.1	8.3	19.7	42.1	11.2	16.7	4.3
qwen1.5-7b-chat	15.1	6.6	18.2	33.0	12.9	14.7	5.3
vicuna-7b-v1.5-16k	10.6	1.3	7.9	35.2	6.6	4.7	7.6
llama-2-7b-chat-hf	10.2	0.0	12.7	35.6	4.8	1.3	6.9
qwen1.5-4b-chat	9.4	4.0	9.1	26.9	6.7	4.0	5.8
vicuna-7b-v1.5	9.4	1.0	7.3	31.3	4.3	4.0	8.7
yi-6b-chat	8.7	1.3	18.4	18.6	7.1	2.0	4.7
qwen1.5-1.8b-chat	5.8	0.0	9.3	20.4	2.1	0.0	3.2
qwen1.5-0.5b-chat	5.4	0.0	4.7	20.1	3.4	1.3	2.9

Our hope is that our work will have a positive impact for both practitioners and researchers: by providing a new benchmark with frequently-updated questions, our work has the potential to both accelerate future research and enable more comprehensive and rigorous evaluations of existing and future models. Furthermore, we hope that the general framework of our benchmark – frequently-updated questions with new information sources – will catch on, mitigating the negative effects of contamination in future LLM evaluation and making LLM benchmarks more ‘future-proof’ in general.

## B Additional Details about LiveBench Experiments

In this section, we detail further descriptions about the LiveBench benchmark itself and our experiments. For example, we include further depictions of the comparisons of LiveBench to ChatBot Arena and Arena-Hard in Figure 3. We display the full results table for LiveBench in Table 3. We display the list of all verifiable instructions in Table 4. Other details are below.

### B.1 Additional Experimental Results Discussion

**Instruction Following Results Discussion** To better understand why Opus achieves a higher ranking than GPT-4o in LiveBench, we manually examined some samples. We noticed that certain instructions were somewhat unusual and could be challenging for aligned models to perform correctly. For instance, here is a case where Claude-3-Opus succeeds while GPT-4o utterly fails:

Table 4: The list of 24 verifiable instructions used in our evaluation with brief descriptions. Although the original paper had 25 instructions, we removed the language instruction. The remaining rows are from [53].

Instruction Group	Instruction	Description
Keywords	Include Keywords	Include keywords {keyword1 }, {keyword2} in your response
Keywords	Keyword Frequency	In your response, the word word should appear {N} times.
Keywords	Forbidden Words	Do not include keywords {forbidden words} in the response.
Keywords	Letter Frequency	In your response, the letter {letter} should appear {N} times.
Length Constraints	Number Paragraphs	Your response should contain {N} paragraphs. You separate paragraphs using the markdown divider: * * *
Length Constraints	Number Words	Answer with at least / around / at most {N} words.
Length Constraints	Number Sentences	Answer with at least / around / at most {N} sentences.
Length Constraints	Number Paragraphs + First Word in i-th Paragraph	There should be {N} paragraphs. Paragraphs and only paragraphs are separated with each other by two line breaks. The {i}-th paragraph must start with word {first_word}.
Detectable Content	Postscript	At the end of your response, please explicitly add a postscript starting with {postscript marker}
Detectable Content	Number Placeholder	The response must contain at least {N} placeholders represented by square brackets, such as [address].
Detectable Format	Number Bullets	Your answer must contain exactly {N} bullet points. Use the markdown bullet points such as: * This is a point.
Detectable Format	Title	Your answer must contain a title, wrapped in double angular brackets, such as <<poem of joy>>.
Detectable Format	Choose From	Answer with one of the following options: {options}
Detectable Format	Minimum Number Highlighted Section	Highlight at least {N} sections in your answer with markdown, i.e. *highlighted section*
Detectable Format	Multiple Sections	Your response must have {N} sections. Mark the beginning of each section with {section_splitter} X.
Detectable Format	JSON Format	Entire output should be wrapped in JSON format.
Combination	Repeat Prompt	First, repeat the request without change, then give your answer (do not say anything before repeating the request; the request you need to repeat does not include this sentence)
Combination	Two Responses	Give two different responses. Responses and only responses should be separated by 6 asterisk symbols: *****.
Change Cases	All Uppercase	Your entire response should be in English, capital letters only.
Change Cases	All Lowercase	Your entire response should be in English, and in all lowercase letters. No capital letters are allowed.
Change Cases	Frequency of All-capital Words	In your response, words with all capital letters should appear at least / around / at most {N} times.
Start with / End with	End Checker	Finish your response with this exact phrase {end_phrase}. No other words should follow this phrase.
Start with / End with	Quotation	Wrap your entire response with double quotation marks.
Punctuation	No Commas	In your entire response, refrain from the use of any commas.

**Question:** {{Prompt}} In your response, the letter i should appear less than 1 times.

**Claude-3-Opus:** «The Strateg«c» Battleground of D«g»tal Campa«gn»ng»\n\nAs the UK gears up for the upcom«ng» general elect«on», pol«t»cal part«es» are d«v»ng headf«rst» «nto» the realm of d«g»tal campa«gn»ng.

**GPT-4o:** «The Impact of Digital Advertising on Election Timing»\n\nWhy schedule an election for November? This query was posed by digital marketing expert Mike Harris in a message to his friend, Labour\u2019s campaign manager, Morgan McSweeney, earlier this year.

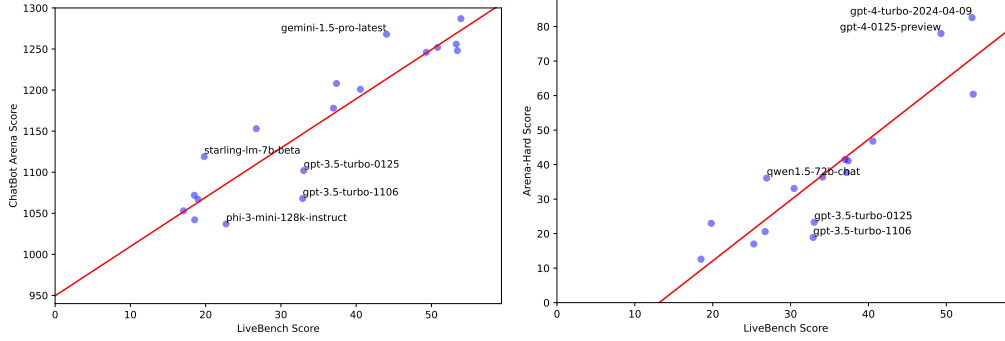


Figure 3: **The performance of models on different benchmarks, compared to a best-fit line.** We compare the different in relative performance of LLMs on LiveBench vs. ChatBot Arena, and LiveBench vs. Arena-Hard. We see that while many LLMs are near the best-fit lines, a few are notable outliers, providing evidence that their output style may be noticeably better or worse than their ability to answer questions.

Table 5: Aggregated scores on instruction following tasks over the four subtasks (Paraphrase, Story Generation, Simplify, and Summarize)

Task	Claude Opus	GPT-4o
change_case	0.954	0.819
combination	0.798	0.735
detectable_content	0.953	0.69
detectable_format	0.907	0.992
keywords	0.904	0.786
length_constraints	0.698	0.702
punctuation	1.0	0.904
startend	0.98	0.977

Additionally, we examined the aggregate scores for various constraints, such as letter frequency, word frequency, word existence, and the number of placeholders. We discovered that while Claude-3-Opus could follow these types of instructions, GPT-4o could not - see [Table 5](#).

We also analyze the performance of Claude-3-Opus and GPT-4o on the original IFeval benchmark ([Table 6](#)). We consider both strict accuracies, which require exact matching, and loose accuracies, where the evaluator reduces false negatives by giving partial points for partial matches (see [53]). While the performance is similar in the strict setting, Claude-3-Opus exhibits a higher upper bound in the loose accuracy setting. Moreover, our hypothesis about the larger performance gap observed in LiveBench compared to IFeval is that LiveBench is more challenging. In LiveBench, the models must understand the document, the task, and the instructions, which can involve complex compositions. In contrast, the original IFeval presents only the task and instructions. The increased difficulty in LiveBench comes from the fact that instead of a simple prompt like "write a blog about Japan," it requires summarizing or analyzing a Guardian article, necessitating comprehension of the text.

Table 6: Claude-3-Opus and GPT-4o performance on IFeval benchmark.

Model	Strict	Strict	Loose	Loose
	Prompt-Level	Instruction-Level	Prompt-Level	Instruction-Level
Claude-3-Opus	0.815	0.872	0.884	0.922
GPT-4o	0.811	0.867	0.852	0.896

Table 7: Model Performance on math and reasoning tasks with both ground-truth (GT) or LLM judging (LLM-Jdg.)

	AMC12 2024		AIME 2024		SMC 2023		Zebra Puzzles	
	GT	LLM-Jdg.	GT	LLM-Jdg.	GT	LLM-Jdg.	GT	LLM-Jdg.
GPT-4-Turbo	54	64.000	13.793	35.714	70.588	58.824	38	68
Claude-3-Opus	56	42.857	6.897	17.241	58.824	52.941	34	52

This increased difficulty in parsing and understanding the context likely contributes to the larger performance gap observed in LiveBench.

## B.2 Details from Ablation Studies

In this section, we give more details from [Section 3](#).

Recall that in [Section 3.3](#), we ran an ablation study by taking three math sub-tasks and one reasoning task, and scoring them by either matching with the ground truth answer, or by asking an LLM judge to score the answer as either correct or incorrect. We used a judge prompt similar to MT-Bench, which we duplicate below. Furthermore, to complement [Table 2](#), we give the model performance scores for ground-truth and LLM judging for the respective models and tasks, in [Table 7](#).

[Instruction] Please act as an impartial judge and evaluate the quality of the response provided by an AI assistant to the user question displayed below. Your evaluation should consider correctness alone. Identify and correct any mistakes. Be as objective as possible. After providing your explanation, you must rate the response on a scale of 1 to 10 by strictly following this format: "[[rating]]", for example: "Rating: [[5]]".  
[Question] question [The Start of Assistant’s Answer] answer [The End of Assistant’s Answer]

Our paper introduces a new benchmark for LLMs, which contains frequently-updated questions from new information sources, scores answers according to objective ground-truth values, and contains a wide variety of tasks. We do not see any inherently negative broader societal impacts of our work.

Our hope is that our work will have a positive impact for both practitioners and researchers: by providing a new benchmark with frequently-updated questions, our work has the potential to both accelerate future research and enable more comprehensive and rigorous evaluations of existing and future models. Furthermore, we hope that the general framework of our benchmark – frequently-updated questions with new information sources – will catch on, mitigating the negative effects of contamination in future LLM evaluation and making LLM benchmarks more ‘future-proof’ in general.