
ZEROSUMEVAL: SCALING LLM EVALUATION WITH INTER-MODEL COMPETITION

Anonymous authors

Paper under double-blind review

ABSTRACT

Evaluating the capabilities of Foundation Models has traditionally relied on static benchmark datasets, human assessments, or model-based evaluations — methods that often suffer from overfitting, high costs, and biases. We introduce ZeroSumEval, a novel competition-based evaluation protocol that leverages zero-sum games to assess LLMs with dynamic benchmarks that resist saturation. ZeroSumEval encompasses a diverse suite of games, including security challenges (Capture the Flag), classic board games (chess), and knowledge tests (MathQuiz). These games are designed to evaluate a range of AI capabilities such as strategic reasoning, planning, knowledge application, safety, and adaptability. A key novelty is integrating automatic prompt optimization to ensure fair comparisons by eliminating biases from human prompt engineering and support arbitrary prompting strategies. Furthermore, ZeroSumEval measures AI models’ abilities to self-improve from limited observations and assesses their robustness against adversarial or misleading examples during prompt optimization. Building upon recent studies that highlight the effectiveness of game-based evaluations for LLMs, ZeroSumEval enhances these approaches by providing a standardized and extensible framework for rigorous assessment. We find ZeroSumEval correlates strongly with expensive human evaluations (Chatbot Arena) and disagrees with benchmarks with known overfitting and saturation issues. Inspecting match traces reveals models that allocate more tokens to thought processes perform strongly in games involving planning capabilities.

1 INTRODUCTION

Large Language Models (LLMs) are being developed at an unprecedented pace (Zhao et al., 2024), requiring significant investment for their training and refinement (Kevin Lee, 2024; Miller, 2022; Kimball, 2024). As the performance and complexity of these models continue to grow (Chen et al., 2024b), selecting the most appropriate model for a specific application has become an increasingly challenging and costly decision (Kaplan et al., 2020; Hoffmann et al., 2022). Benchmarking emerges as a critical tool in this context (Laskar et al., 2023; Qin et al., 2023), providing standardized metrics and evaluations to guide these choices.

With the rapid growth of generative technologies built on top of Large Language Models (OpenAI, 2022; Google, 2024; Anthropic, 2024b; Ormazabal et al., 2024; Mistral, 2024; Dubey et al., 2024a; Yang et al., 2024), it has been increasingly difficult to evaluate these models comprehensively (Guo et al., 2023). Current benchmarking practices face several significant issues. Many benchmarks suffer from data contamination (Yang et al., 2023), where models inadvertently train on portions of the test data (Dubey et al., 2024a; Groeneveld et al., 2024), leading to inflated performance metrics. Sensitivity to prompt variations (Alzahrani et al., 2024b) and a lack of diversity in evaluation tasks (Laskar et al., 2024) further undermine the reliability and robustness of these benchmarks. Additionally, the high cost and effort required to develop new benchmarks often result in outdated evaluation methods that do not keep pace with the rapid development of LLMs (Kiela et al., 2021; Vu et al., 2023).

An observed disparity exists between the computational resources measured in floating-point operations per second, or FLOPs used to train LLMs and those allocated for their evaluation. Training these models involves massive computational efforts (Hoffmann et al., 2022), yet the evaluation

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

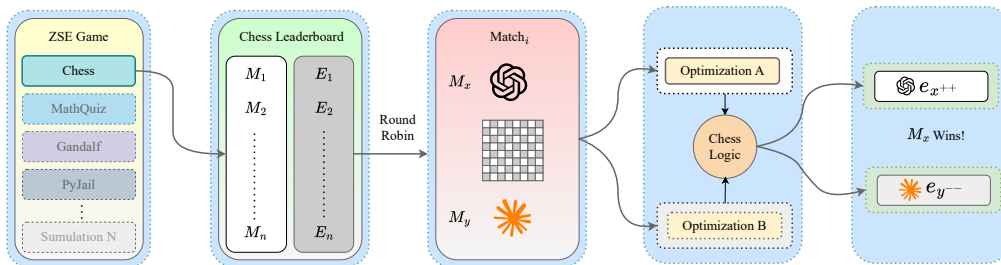


Figure 1: The ZEROSUMEVAL suite of benchmarks provides dynamic simulations with head to head model competition to create robust and scalable model evaluations and leaderboards. Integrated automatic prompt optimization minimizes biases introduced by prompting and hand-engineering.

phase typically utilizes a negligible fraction of this capacity (Laskar et al., 2024). Scaling up evaluation by increasing the number of evaluation tokens is essential for a more thorough understanding of model capabilities. Traditionally, this scaling involves incorporating human-crafted independent and identically distributed (i.i.d.) data (Holland et al., 2018), which is resource-intensive (Hutchinson et al., 2021) and may not adequately capture the complexities of language (Mehrabi et al., 2021) and reasoning required to challenge advanced LLMs (Gudibande et al., 2023) or even LLM generated (Karpinska et al., 2021).

Previous work has proposed the use of games as benchmarks (Topsakal et al., 2024), offering a promising avenue for evaluating complex reasoning (Wong et al., 2023) and decision-making abilities of LLMs (Warstadt et al., 2023; Park et al., 2023; Wang et al., 2023). Games provide interactive and dynamic environments that can test models beyond static datasets. However, existing game-based benchmarks are often (i) inflexible and limited in scope, (ii) not easily extendable, (iii) restricted in their effectiveness for comprehensive model evaluation, and (iv) depend on predefined prompts.

Scaling evaluation is fundamental not only for assessing performance but also for uncovering hidden dynamics within LLMs, such as potential backdoors or biases (Schuster et al., 2020), and for evaluating their emerging reasoning capabilities (Brown et al., 2020; Sanh et al., 2022; Wei et al., 2023b;a). Implementing environments for simulations or games offers a scalable solution to these challenges (OpenAI et al., 2019; OpenAI, 2019; Silver et al., 2016; 2017; Zheng et al., 2021).

Existing evaluation protocols possess several key issues:

(i) **Prompt Sensitivity:** Previous work (Zheng et al., 2024; Pezeshkpour & Hruschka, 2023; Lu et al., 2022; Alzahrani et al., 2024a; Wang et al., 2024a) has shown that models are sensitive to benchmark formats. By sheer chance, a model could be presented with a prompting method that’s either favorable or detrimental. These prompt modifications are shown to result in substantially different relative performance between models (Alzahrani et al., 2024a). By testing models in varied scenarios within a controlled environment, we can assess and improve their robustness to different prompts. Crucially, different models are not optimized for the same prompts due to variations in data mixtures and algorithmic implementations. Using identical prompts across all models may therefore lead to unfair comparisons.

(ii) **Limited Diversity:** Traditional evaluation methods often rely on static datasets, which are inherently limited by their dependency on human curation and annotation. This makes it challenging to continuously introduce new, diverse test data. An extensible simulated environment, however, allows for a wide array of dynamically generated games and scenarios, enhancing the diversity and scalability of evaluation tasks.

(iii) **Extensibility:** Once established, the environment can be easily expanded to include new games, rules, and scenarios, facilitating continuous evaluation improvements.

(iv) **Crowd and Annotator Bias:** LLM evaluations conducted by large crowds often tend to be susceptible to social hacking, and it can depend on geographic, temporal, and narrative factors Gururan-

gan et al. (2018). Controlled and interpretable environments can mitigate these biases by providing consistent, objective evaluation criteria.

(v) **Saturation:** With the rapid improvement of LLMs, evaluation benchmarks quickly become obsolete and saturated, with frontier models achieving almost perfect scores, which necessitates the development of new benchmarks. On the opposite extreme, benchmarks that are too difficult would result in almost random scores. Both extremes result in a lack of granularity to distinguish models. Therefore, benchmarks posing moderate difficulty to frontier models will need to be continuously developed as models improve. For instance, GSM8K (Cobbe et al., 2021) tests models on grade school-level math, and most state-of-the-art models achieve scores above 90% (Dubey et al., 2024b; Anthropic, 2024a). Thus, the more difficult MATH (Hendrycks et al., 2021b) dataset, which consists of math competition questions, was developed and is now commonly used¹. A similar trend is observed in academic examination benchmarks with the migration from MMLU (Hendrycks et al., 2020) to MMLU-Pro (Wang et al., 2024b) and GPQA. (Rein et al., 2023)².

To address these challenges, we introduce ZEROSUMEEVAL, a flexible and extensible open-source framework designed to scale LLM evaluation through the simulation of two-player zero-sum games. Our framework allows for comprehensive and robust assessment by providing models with multiple opportunities to make legal moves, thereby accommodating occasional errors and offering a more nuanced understanding of their capabilities.

1. **Scaling Evaluation by Simulation:** We demonstrate how simulation environments can effectively scale the evaluation process.

2. **Flexible and Extensible Framework:** ZEROSUMEEVAL is designed to be adaptable, allowing researchers and practitioners to customize and extend the evaluation environment to suit diverse needs.

3. **Robustness to Prompt Sensitivity:** By incorporating automatic prompt optimization, our framework mitigates issues related to prompt sensitivity, leading to more reliable evaluation outcomes.

4. **Enhanced Interpretability:** The structured environment facilitates easier interpretation of model behaviors, aiding in the identification of strengths and weaknesses.

5. **Error Accommodation:** Models are given multiple chances to make legal moves, ensuring that occasional missteps due to inherent stochasticity do not disproportionately affect the overall evaluation.

2 RELATED WORK

2.1 STATIC LLM BENCHMARKS

Until recently, LLMs were evaluated on Natural Language Understanding (NLU) tasks from benchmark collections like GLUE (Wang, 2018) and SuperGLUE (Wang et al., 2019), which included tasks like paraphrase classification and sentiment analysis. As LLMs developed, they acquired emergent capabilities beyond generating plausible text, such as reasoning, generating code, and instruction following (Brown et al., 2020; Wei et al., 2022). With these newly found capabilities, new benchmarks were developed to quantify these abilities. As models improve, more difficult benchmarks are created. For example:

- **Reasoning:** undergraduate level academic questions are tested via MMLU (Hendrycks et al., 2020), while GPQA (Rein et al., 2023) tests models with graduate level questions. All aforementioned benchmarks score models based on the likelihood of specific tokens for the answer keys in a multiple-choice setting.

- **Mathematics:** GSM8K (Cobbe et al., 2021) evaluates models on elementary level arithmetic, while MATH (Hendrycks et al., 2021b) tests on competition level mathematics. Both benchmarks evaluate the model in a few-shot setting by encouraging models to output chains of thought followed by the numeric answer in a specific format.

¹HuggingFace’s Open LLM Leaderboard (Beeching et al., 2023; Fourrier et al., 2024) migrated from GSM8K in v1 to MATH in v2.

²Similar to 1, the leaderboard transitioned from MMLU in v1 to MMLU-Pro and GPQA in v2.

162 • **Coding:** HumanEval (Chen et al., 2021) test models on basic coding, while APPS (Hendrycks
163 et al., 2021a) uses coding competition questions. These benchmarks generate Python code by
164 prompting LLMs with function docstrings or written specifications, and run input/output test-cases
165 on the generated code.

166
167 Criticism of these types of static benchmarks are outlined in Section 1.

168 169 2.2 COMPARATIVE LLM BENCHMARKS

170
171 **LLM Game Evaluations** To address the static benchmark issues highlighted in Section 1, the
172 paradigm of evaluating agentic capabilities through simulations has been applied successfully in
173 multiple prior works. Evaluation frameworks comprising multiple games include: (i) *ChatArena*
174 (Wu et al., 2023), which includes Chess, Tic-Tac-Toe, Rock-Paper-Scissors, and others, (ii)
175 *GridGames* (Topsakal et al., 2024), implementing Tic-Tac-Toe, Connect Four, and Gomoku, and
176 (iii) *GameBench* (Costarelli et al., 2024), which is the most diverse, as they developed 9 games,
177 include non-deterministic and imperfect information games.

178
179 **Limitations of LLM Game Evaluations** All the aforementioned benchmark frameworks are im-
180 plemented with manually written prompts for all models, and sometime suggest a strategy within
181 the prompt, such as ChatArena prompting models to output a random move in Rock-Paper-Scissors.
182 GameBench tries to optimize model results by utilizing two prompting strategies: (i) Chain of
183 Thought (CoT), and (ii) Reasoning via Planning (RAP), but the issue of static prompt still persists.
184 This could explain the poor performances they observed, such as GPT-4 achieving almost random
185 results on some tasks.

186
187 **Comparative Human Evaluations** A popular head-to-head LLM evaluation framework is Chat-
188 bot Arena³ (Chiang et al., 2024), which allows users to prompt two anonymous LLMs with arbitrary
189 prompts and to vote for the better response. This creates a diverse evaluation that effectively ranks
190 all models in a leaderboard. However, it suffers from two issues: (i) human evaluations are slow
191 and laborious, and adding new models requires prolonged evaluation periods until sufficient votes
192 are acquired for a confident placement, and (ii) human evaluations contain human biases, such as
193 prompt over-representation (Dunlap et al., 2024) and bias to verbose and “pretty” responses (Chen
194 et al., 2024a; Park et al., 2024; Li et al., 2024).

195 3 METHODOLOGY

196
197 In this section, we describe the technical details of ZEROSUMEEVAL including design choices, the
198 importance of automatic prompt optimization, and game selection/categorization. At its core, ZE-
199 ROSUMEEVAL provides controlled environments to observe models competing against each other to
200 win competitive games. In particular, ZSE controls (i) the role and information each model has ac-
201 cess to at any point in the simulation and (ii) the data models can use to optimize/modify their own
202 prompts.

203 204 3.1 CAPABILITIES

205
206 The games within ZSE are designed to evaluate specific capabilities in a controlled environment:

207
208 **Reasoning** Board games and cybersecurity scenarios require models to perform complex, multi-
209 step reasoning. They test the models’ ability to process information, predict outcomes, and formulate
210 strategies in dynamically changing environments.

211
212 **Planning** Board games also involve long-term strategy, requiring models to anticipate the conse-
213 quences of their actions several moves ahead. This assesses the model’s foresight, adaptability, and
214 capacity for nuanced decision-making.

215
³formerly LMSYS, not to be confused with ChatArena.

Knowledge Application Models must recall and apply mathematical knowledge to solve problems in question answering type games. This setup provides a direct assessment of the models’ ability to retrieve, interpret, and implement factual information in structured problem-solving.

Creativity Models successful at cybersecurity type games must exhibit creativity to successfully create secure environments and break them.

3.2 GAME DESIGN

ZEROSUM-EVAL supports an expanding suite of game types designed to test the aspects of LLM performance described above. The mix we showcase includes both well-known and established games, such as chess, as well as more special-purpose games (e.g. MathQuiz). For completeness and reproducibility, we describe the implementations of MathQuiz and PyJail. The following set of games are selected to encompass a range of cognitive capabilities, including strategic reasoning, planning, knowledge application, and creativity:

Board Games (Chess) Classic board games like chess serve as a benchmark for strategic reasoning and long-term planning. They require models to engage in multi-step thinking, manage trade-offs, and foresee opponent moves. This category is instrumental in evaluating a model’s ability to plan several moves ahead, adapt its strategies, and make complex decisions under uncertainty⁴.

Question-Answer Games (MathQuiz) These games are constructed to measure models’ knowledge recall and logical reasoning abilities. MathQuiz, for instance, challenges models to both create and answer arithmetic and mathematical questions, assessing their understanding of mathematical concepts, computational accuracy, and step-by-step problem-solving skills. Our implementation of MathQuiz tasks a teacher player to create a challenging math problem and prove that the problem is valid and solvable. A student player then attempts to answer the generated math problem. The student wins the game by answering the question correctly or if the teacher fails to create a valid question.

Cybersecurity Games (PyJail) PyJail involves python “capture the flag” cybersecurity challenges, targeting the model’s ability to create puzzles and interact with a restricted python environment to strategize solutions. The PyJail game is structured into three stages. The first statically parses a player generated PyJail program to provide feedback on the syntax and semantic structure. Given validity, the challenge code is inserted into the environment, and the same player model must commit a solution that is tested dynamically to prove the challenge’s feasibility. A unique flag is stored in the target variable at runtime, which prevents any trivial method to cheat the challenge. The second player will complete the same step, provided a restricted view of the environment and limited context. The game ends if first player is unable to create a valid challenge or the flag is retrieved by the attacker.

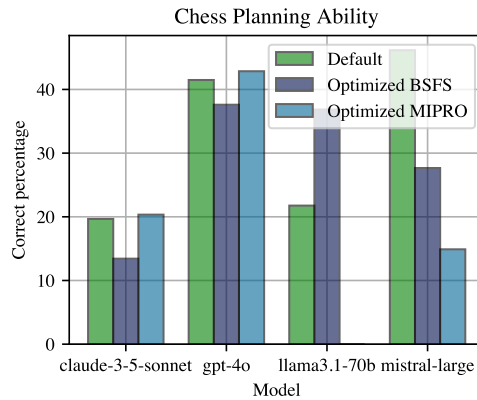


Figure 2: The effect of prompt optimization on the proportion of correct moves. Moves are classified as correct if the evaluation, as determined by Stockfish 17 (The Stockfish Developers, 2024) with depth 15, does not decrease by more than 0.3 points (pawn equivalent). Models react differently to prompts and have varying prompt optimization abilities.

⁴Chess has a rich history as a testbed for strategy and planning. See <https://github.com/carlini/chess-llm> and <https://huggingface.co/spaces/mlabonne/chessllm> for examples of LLMs playing chess.

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323

3.3 SCALABLE VERIFICATION

The MathQuiz and PyJail games require competing models to generate complex challenge environments and solutions. Since verification of the knowledge-based challenges by a human in the loop is not scalable, we design a method to verify model output using an automated manager in a two-fold generation and verification process. This is accomplished by defining a target outcome (e.g., the answer to a math question or a CTF flag) as the basis for verifying generated input, and regulating the model context at each stage.

The exact process (illustrated in Figure 3) is outlined as follows:

- (i) The generator model receives a target and attempts to output a valid challenge that resolves to the specific target.
- (ii) In the verification step, the manager restricts the model’s context to ensure no direct access to the target, and asks the generator model to solve the previously generated challenge.
- (iii) If the manager determines the verification is successful (by matching the target with the generator’s solution), the game proceeds. Otherwise, the generator model is deemed to have failed to generate a valid challenge.

This method ensures the generated challenge environment is valid and a solution is proven possible by the generator. The design also correctly penalizes models that directly generate memorized questions as it is likely to have been memorized by other models, thereby encouraging models to create challenging and novel questions. Finally, the scalability of the evaluation is preserved as the capabilities of models scale.

3.4 AUTOMATIC PROMPTING

Automatic prompting is an essential component of the ZEROSUMEVAL framework for several reasons. First, it allows models to learn to play new games through self-optimization, demonstrating their ability to adapt to different scenarios without human intervention. Second, it removes the human element in prompt engineering, thereby reducing biases and variations introduced by manual prompt construction (Zheng et al., 2024; Pezeshkpour & Hruschka, 2023; Alzahrani et al., 2024a). Third, automatic prompting serves as a measure of a model’s ability to self-improve at inference time, providing insight into its adaptability and strategic reasoning skills.

We leverage the DSPy (Khattab et al., 2023) approach to implement automatic prompt optimization in our framework. DSPy allows models to autonomously explore and select optimal prompts based on the current game context, dynamically adjusting strategies to maximize performance. We also make use of DSPy Assertions (Singhvi et al., 2024) to simulate interactivity between the models and the game environment by allowing a number of retries (with feedback from the game) when the model makes an invalid move. Although we find DSPy has the flexibility and generalizability to support various models and games, ZEROSUMEVAL supports alternative automatic prompt optimization techniques if required.

Through prompt optimization, models can develop improved strategies as they encounter diverse game scenarios. For example, in a chess game, models equipped with optimized prompting demonstrated a higher proportion of correct moves compared to their counterparts using default prompts Figure 2. This not only reveals the models’ enhanced strategic reasoning but also emphasizes the significance of prompt optimization in robust performance evaluation.

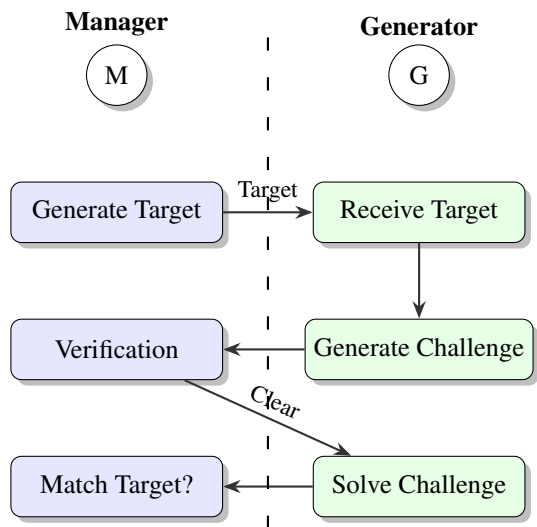


Figure 3: State diagram of the verification process involving the Game Manager and the Generator. Blue boxes indicate deterministic steps and green boxes indicate steps involving the model.

By incorporating automatic prompting, ZEROSUMEVAL addresses benchmark sensitivity. The prompt optimization integrates game validation mechanisms into the optimization process, allowing models to observe tangible outcomes and refine their prompt strategies. Consequently, this mitigates the variations in performance due to prompt sensitivity, leading to a more consistent and reliable evaluation of model capabilities.

Datasets and Optimizers To perform the automatic prompt optimization process, models require examples of gameplay (inputs and outputs) and prompt optimizers. We create standard datasets manually for each game available to all models for the optimization. The available datasets are described in Table 1. Through DSPy, ZEROSUMEVAL supports multiple types of optimizers. In this work, we focus on (i) BootstrapFewShot (ii) BootstrapFewShotRandomSearch (Khatab et al., 2023) and (iii) MIPROv2 (Opsahl-Ong et al., 2024).

Dataset	Source	Description
chess_stockfish	conacts/stockfish_dataset ⁵	stockfish vs stockfish games
chess_puzzles	(Schwarzschild et al., 2021)	chess puzzles.
mathquiz_gsm8k	(Cobbe et al., 2021)	grade school level math QA
mathquiz_hendrycks_math	(Hendrycks et al., 2021b)	advanced math QA
pyjail_ctf_llm	(Shao et al., 2024)	Pyjail style Capture The Flags (CTFs).

Table 1: Overview of datasets used in the evaluation framework.

An interesting direction out of the scope of this work is enabling models to learn games via self-play. This would reduce manual effort needed to create new games for ZEROSUMEVAL and measure a model’s ability to effectively explore a space without supervision.

3.5 RATINGS

ZEROSUMEVAL utilizes an easily computable rating system derived from the outcomes of competitive games between models. Each model receives a rating based on its win-loss record over multiple games, allowing for a rapid and scalable oversight of model capabilities. This framework seamlessly incorporates new games, providing continuous and dynamic evaluation as models improve.

Following recent suggestions for LLM rating systems by Boubdir et al. (2023); Chiang et al. (2023), we employ the Bradley-Terry (BT) rating system, an alternative to the Elo system, to rate models. The BT model is permutation-invariant and assumes a fixed win rate for each model pair, maximizing the likelihood of observed outcomes (Bradley & Terry, 1952). This choice is more suitable than the traditional Elo system, which was designed for human chess players with varying skill levels, whereas LLMs have fixed skill levels defined by their weights (Elo, 1967).

ZEROSUMEVAL’s rating system facilitates analysis of model behaviors. It allows us to observe not only the relative strategic planning capabilities of models but also their capacity for self-improvement through prompt optimization. For instance, analysis of models’ gameplay strategies in chess revealed that prompt-optimized models allocate more reasoning words in their decision-making process, suggesting a deeper level of planning (Figure 4).

4 EXPERIMENTS

In this section, we describe the experiments to demonstrate the effectiveness of the ZEROSUMEVAL as a dynamic leaderboard. We also design experiments to evaluate the effect of prompt optimization on the performance of various large language models (LLMs) under various simulations.

4.1 MODEL SELECTION AND EXPERIMENTAL SETUP

We select four models of varying sizes and capabilities for this study: GPT-4o, Claude 3.5 Sonnet, LLaMA 3.1-70B-Instruct, and Mistral-Large. These models represent a range of architectures and training scales, providing a diverse set for evaluating the generalizability of the ZEROSUMEVAL framework.

378 The experiments involve running a multiple round-robin tournaments to simulate competitive game-
 379 play among the model (50-100 games per experiment). In addition to measuring model perfor-
 380 mance on the games in the ZEROSUMEVAL suite, we also examine how the models’ performance
 381 changes with different prompt optimization techniques. Each tournament round involves all possi-
 382 ble match permutations between model variants, after which the models’ ratings are calculated using
 383 the Bradley-Terry model (Bradley & Terry, 1952). The primary goal of this ablation study is to as-
 384 sess each model’s responsiveness to the optimization process and to identify resulting behavioral
 385 changes.

386 For the automatic prompt optimization, we utilize three optimizers commonly used in DSPy: Boot-
 387 strapFewshot (BSFS), BootstrapFewshotWithRandomSearch (BSFSRS), and MIPROv2, targeting
 388 the ChainOfThought module in DSPy.
 389

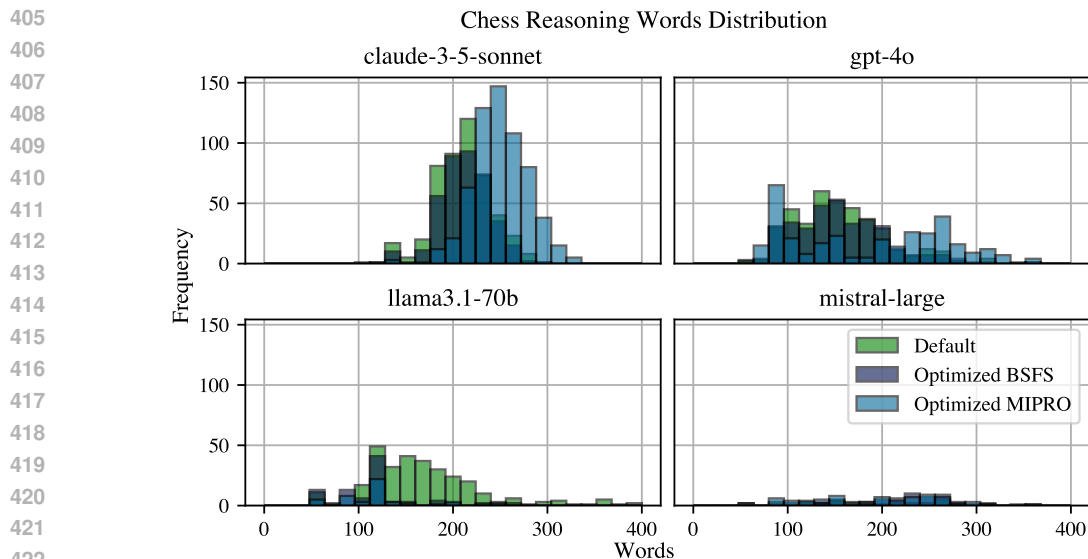
390 4.2 GAMES FOR ANALYSIS

391 Although ZEROSUMEVAL supports a range of games for assessing different capabilities, our de-
 392 tailed set of experiments focus primarily on Chess to analyze the models’ planning abilities. This
 393 decision is motivated by the interpretability of Chess gameplay and its complexity, which provides
 394 an ideal testbed for assessing strategic reasoning and decision-making.
 395

397 5 RESULTS

399 5.1 RATINGS AND PERFORMANCE TRENDS

401 Table 2 provides the ratings for each model variant across the games. The results indicate GPT-
 402 4o and Claude 3.5 Sonnet typically perform best with GPT-4o slightly ahead. This agrees with
 403 leaderboards based on human ratings, such as ChatbotArena (Chiang et al., 2024).
 404



424 Figure 4: The distribution of CoT words used for each model and prompt optimization technique.
 425 In general, prompt optimized models spend more words reasoning than their non-optimized coun-
 426 terparts, especially with MIPROv2 optimization.
 427

429 5.2 IMPACT OF PROMPT OPTIMIZATION ON PERFORMANCE

430 The experimental results (Table 3) reveal significant variations in model performance as a result of
 431 prompt optimization. Prompt optimization can even flip ranking as is the case with MIPROv2 -

highlighting the significant effect of prompt sensitivity. Prompt-optimized models typically exhibit improved strategic reasoning, as evidenced by an increased number of correct moves and a more favorable distribution of move evaluations (Figure 2). ZEROSUM EVAL provides the capability to compare models across prompt optimization strategies, leading to fairer evaluations and more robust leaderboards.

Figure 4 illustrates the distribution of CoT words for each model with different prompt optimization techniques. Notably, models optimized using MIPROv2 demonstrate a tendency to allocate more words to their reasoning process compared to their default counterparts, suggesting deeper planning and strategic consideration.

Model	Chess (MIPRO)	MathQuiz (Default)	PyJail (Default)
GPT-4o	1202.97	1048.12	1025.58
Claude 3.5 Sonnet	1000.00	962.51	1017.17
Mistral-Large	940.88	982.85	1000.00
LLaMA 3.1 70B	856.15	1006.52	953.15

Table 2: Performance ratings of various models across different tasks. The ratings are computed using the MIPRO-optimized approach for the Chess task and default settings for MathQuiz and PyJail tasks.

Model	Default Rating (CI)	BSFS Rating (CI)	BSFSRS Rating (CI)	MIPRO Rating (CI)
Claude 3.5 Sonnet	1028 (890-1153)	1000 (871-1126)	1000 (862-1147)	984 (837-1060)
Mistral-Large	942 (889-1005)	1016 (963-1073)	1014 (952-1069)	1023 (965-1090)
LLaMA 3.1 70B	978 (918-1054)	951 (888-1039)	1030 (962-1089)	1035 (967-1107)
GPT-4o	962 (880-1034)	1016 (909-1101)	966 (874-1044)	1055 (987-1133)

Table 3: Results of engaging each model in competition against itself optimized by our choices of optimizers. We can see that for Mistral, LLaMA, and GPT-4o, MIPRO outperforms all other optimizers. It is interestingly not the case with Claude. Ratings are shown with their 95% confidence intervals (CI). The highest rating for each model is in bold.

6 CONCLUSION

The dynamic, competitive nature of ZEROSUM EVAL’s evaluation provides a more robust and trustworthy measurement of AI model capabilities, advancing the state of benchmarking in large language models. By leveraging zero-sum games, we ensure that models are consistently challenged with diverse, evolving tasks, minimizing the risk of overfitting and saturation commonly observed in static benchmarks. Additionally, the integration of automatic prompt optimization offers a more holistic evaluation framework that captures not only a model’s static performance but also its dynamic capacity for self-improvement.

REFERENCES

- Norah Alzahrani, Hisham Abdullah Alyahya, Yazeed Alnumay, Sultan Alrashed, Shaykhah Alsubaie, Yusef Almushaykeh, Faisal Mirza, Nouf Alotaibi, Nora Altwairesh, Areeb Alowisheq, M Saiful Bari, and Haidar Khan. When benchmarks are targets: Revealing the sensitivity of large language model leaderboards, 2024a. URL <https://arxiv.org/abs/2402.01781>.
- Norah Alzahrani, Hisham Abdullah Alyahya, Yazeed Alnumay, Sultan Alrashed, Shaykhah Alsubaie, Yusef Almushaykeh, Faisal Mirza, Nouf Alotaibi, Nora Altwairesh, Areeb Alowisheq, et al. When benchmarks are targets: Revealing the sensitivity of large language model leaderboards. *arXiv preprint arXiv:2402.01781*, 2024b.
- Anthropic. Claude 3.5 Sonnet. <https://www.anthropic.com/news/claude-3-5-sonnet>, 2024a. Accessed: 2024-09-17.

486 AI Anthropic. Claude 3.5 sonnet model card addendum. *Claude-3.5 Model Card*, 2024b.
487

488 Edward Beeching, Clémentine Fourier, Nathan Habib, Sheon Han, Nathan Lambert, Nazneen
489 Rajani, Omar Sanseviero, Lewis Tunstall, and Thomas Wolf. Open llm leaderboard
490 (2023-2024). [https://huggingface.co/spaces/open-llm-leaderboard-old/
491 open_llm_leaderboard](https://huggingface.co/spaces/open-llm-leaderboard-old/open_llm_leaderboard), 2023.

492 Meriem Boubdir, Edward Kim, Beyza Ermis, Sara Hooker, and Marzieh Fadaee. Elo uncovered:
493 Robustness and best practices in language model evaluation, 2023. URL [https://arxiv.
494 org/abs/2311.17295](https://arxiv.org/abs/2311.17295).

495

496 Ralph Allan Bradley and Milton E Terry. Rank analysis of incomplete block designs: I. the method
497 of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952.

498

499 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhari-
500 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal,
501 Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M.
502 Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz
503 Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec
504 Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020. URL
505 <https://arxiv.org/abs/2005.14165>.

506

507 Guiming Hardy Chen, Shunian Chen, Ziche Liu, Feng Jiang, and Benyou Wang. Humans or llms
508 as the judge? a study on judgement biases, 2024a. URL [https://arxiv.org/abs/2402.
509 10669](https://arxiv.org/abs/2402.10669).

509

510 Hailin Chen, Fangkai Jiao, Xingxuan Li, Chengwei Qin, Mathieu Ravaut, Ruochen Zhao, Caiming
511 Xiong, and Shafiq Joty. Chatgpt’s one-year anniversary: Are open-source large language models
512 catching up?, 2024b. URL <https://arxiv.org/abs/2311.16989>.

513

514 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared
515 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri,
516 Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan,
517 Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian,
518 Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fo-
519 tios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex
520 Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders,
521 Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec
522 Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob Mc-
523 Grew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large
524 language models trained on code, 2021. URL <https://arxiv.org/abs/2107.03374>.

525

526 Wei-Lin Chiang, Tim Li, Joseph E. Gonzalez, and Ion Stoica. Chatbot arena: New
527 models & elo system update, Dec 2023. URL [https://lmsys.org/blog/
528 2023-12-07-leaderboard/](https://lmsys.org/blog/2023-12-07-leaderboard/).

529

530 Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li,
531 Dacheng Li, Hao Zhang, Banghua Zhu, Michael Jordan, Joseph E. Gonzalez, and Ion Sto-
532 ica. Chatbot arena: An open platform for evaluating llms by human preference, 2024. URL
533 <https://arxiv.org/abs/2403.04132>.

534

535 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
536 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John
537 Schulman. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*,
538 2021.

539

535 Anthony Costarelli, Mat Allen, Roman Hauksson, Grace Sodunke, Suhas Hariharan, Carlson Cheng,
536 Wenjie Li, Joshua Clymer, and Arjun Yadav. Gamebench: Evaluating strategic reasoning abilities
537 of llm agents, 2024. URL <https://arxiv.org/abs/2406.06613>.

538

539 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony

540 Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark,
541 Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere,
542 Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris
543 Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong,
544 Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny
545 Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino,
546 Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael
547 Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Ander-
548 son, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah
549 Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan
550 Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Ma-
551 hadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy
552 Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak,
553 Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Al-
554 wala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini,
555 Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yearly, Laurens van der
556 Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo,
557 Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Man-
558 nat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova,
559 Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal,
560 Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur
561 Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhar-
562 gava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong,
563 Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic,
564 Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sum-
565 baly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa,
566 Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang,
567 Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende,
568 Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney
569 Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom,
570 Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta,
571 Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vladan Petro-
572 vic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang,
573 Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuwei Wang, Yaelle Goldschlag, Yashesh Gaur,
574 Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre
575 Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aaron Grattafiori, Abha
576 Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay
577 Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein, Amanda
578 Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew
579 Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita
580 Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh
581 Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De
582 Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Bran-
583 don Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina
584 Mejia, Changan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai,
585 Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li,
586 Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana
587 Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil,
588 Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Ar-
589 caute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco
590 Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella
591 Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory
592 Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang,
593 Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Gold-
man, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman,
James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer
Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe
Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie
Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun

594 Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kun Huang, Kunal
595 Chawla, Kushal Lakhota, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva,
596 Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian
597 Khabsa, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson,
598 Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Ke-
599 neally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel
600 Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mo-
601 hammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navy-
602 ata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong,
603 Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli,
604 Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux,
605 Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao,
606 Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li,
607 Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott,
608 Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Sa-
609 tadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lind-
610 say, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Shuqiang
611 Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen
612 Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin Cho,
613 Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser,
614 Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Tim-
615 othy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan,
616 Vinay Satish Kumar, Vishal Mangla, Vitor Albiero, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu
617 Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Con-
618 stable, Xiaocheng Tang, Xiaofang Wang, Xiaojian Wu, Xiaolan Wang, Xide Xia, Xilun Wu,
619 Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi,
620 Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, Zachary DeVito, Zef
621 Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. The llama 3 herd of models, 2024a.
622 URL <https://arxiv.org/abs/2407.21783>.

622 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
623 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony
624 Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark,
625 Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere,
626 Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris
627 Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong,
628 Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny
629 Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino,
630 Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael
631 Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Ander-
632 son, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah
633 Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan
634 Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Ma-
635 hadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy
636 Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak,
637 Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Al-
638 wala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini,
639 Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yearly, Laurens van der
640 Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo,
641 Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Man-
642 nat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova,
643 Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal,
644 Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur
645 Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhar-
646 gava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong,
647 Rogavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic,
648 Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sum-
649 baly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa,
650 Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang,

648 Sharath Rapparthi, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende,
649 Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney
650 Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom,
651 Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta,
652 Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vladan Petro-
653 vic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang,
654 Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuwei Wang, Yaelle Goldschlag, Yashesh Gaur,
655 Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre
656 Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aaron Grattafori, Abha
657 Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay
658 Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein, Amanda
659 Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew
660 Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita
661 Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh
662 Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De
663 Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Bran-
664 don Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina
665 Mejia, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai,
666 Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li,
667 Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana
668 Liskovich, Didem Foss, Dingkan Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil,
669 Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Ar-
670 caute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco
671 Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella
672 Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory
673 Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang,
674 Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Gold-
675 man, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman,
676 James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer
677 Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe
678 Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie
679 Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun
680 Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kun Huang, Kunal
681 Chawla, Kushal Lakhotia, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva,
682 Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian
683 Khabsa, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson,
684 Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Ke-
685 neally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel
686 Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mo-
687 hammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navy-
688 ata Bawa, Nayan Singhal, Nick Gebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong,
689 Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli,
690 Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux,
691 Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao,
692 Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li,
693 Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott,
694 Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Sa-
695 tadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lind-
696 say, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Shuqiang
697 Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen
698 Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin Cho,
699 Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser,
700 Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Tim-
701 othy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan,
Vinay Satish Kumar, Vishal Mangla, Vitor Albiero, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu
Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Con-
stable, Xiaocheng Tang, Xiaofang Wang, Xiaojuan Wu, Xiaolan Wang, Xide Xia, Xilun Wu,
Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi,
Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, Zachary DeVito, Zef

702 Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. The llama 3 herd of models, 2024b.
703 URL <https://arxiv.org/abs/2407.21783>.
704

705 Lisa Dunlap, Evan Frick, Tianle Li, Isaac Ong, Joseph E. Gonzalez, and Wei-Lin Chiang. What’s
706 up with llama 3? arena data analysis, May 2024. URL [https://lmsys.org/blog/
707 2024-05-08-llama3/](https://lmsys.org/blog/2024-05-08-llama3/).

708 Arpad E Elo. The proposed uscf rating system, its development, theory, and applications. *Chess*
709 *life*, 22(8):242–247, 1967.

710 Clémentine Fourier, Nathan Habib, Alina Lozovskaya, Konrad Szafer, and Thomas Wolf. Open
711 llm leaderboard v2. [https://huggingface.co/spaces/open-llm-leaderboard/
712 open_llm_leaderboard](https://huggingface.co/spaces/open-llm-leaderboard/open_llm_leaderboard), 2024.

713

714 Google. Gemini: A family of highly capable multimodal models, 2024.

715 Dirk Groeneveld, Iz Beltagy, Pete Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord,
716 Ananya Harsh Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, Shane Arora, David Atkin-
717 son, Russell Authur, Khyathi Raghavi Chandu, Arman Cohan, Jennifer Dumas, Yanai Elazar,
718 Yuling Gu, Jack Hessel, Tushar Khot, William Merrill, Jacob Morrison, Niklas Muennighoff,
719 Aakanksha Naik, Crystal Nam, Matthew E. Peters, Valentina Pyatkin, Abhilasha Ravichander,
720 Dustin Schwenk, Saurabh Shah, Will Smith, Emma Strubell, Nishant Subramani, Mitchell Worts-
721 man, Pradeep Dasigi, Nathan Lambert, Kyle Richardson, Luke Zettlemoyer, Jesse Dodge, Kyle
722 Lo, Luca Soldaini, Noah A. Smith, and Hannaneh Hajishirzi. Olmo: Accelerating the science of
723 language models, 2024. URL <https://arxiv.org/abs/2402.00838>.

724 Arnav Gudibande, Eric Wallace, Charlie Snell, Xinyang Geng, Hao Liu, Pieter Abbeel, Sergey
725 Levine, and Dawn Song. The false promise of imitating proprietary llms, 2023. URL <https://arxiv.org/abs/2305.15717>.
726

727 Zishan Guo, Renren Jin, Chuang Liu, Yufei Huang, Dan Shi, Supryadi, Linhao Yu, Yan Liu, Jiaxuan
728 Li, Bojian Xiong, and Deyi Xiong. Evaluating large language models: A comprehensive survey,
729 2023. URL <https://arxiv.org/abs/2310.19736>.

730

731 Suchin Gururangan, Swabha Swayamdipta, Omer Levy, Roy Schwartz, Samuel Bowman, and
732 Noah A. Smith. Annotation artifacts in natural language inference data. In Marilyn Walker,
733 Heng Ji, and Amanda Stent (eds.), *Proceedings of the 2018 Conference of the North American*
734 *Chapter of the Association for Computational Linguistics: Human Language Technologies, Vol-*
735 *ume 2 (Short Papers)*, pp. 107–112, New Orleans, Louisiana, June 2018. Association for Com-
736 putational Linguistics. doi: 10.18653/v1/N18-2017. URL [https://aclanthology.org/
737 N18-2017](https://aclanthology.org/N18-2017).

738 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and
739 Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint*
740 *arXiv:2009.03300*, 2020.

741 Dan Hendrycks, Steven Basart, Saurav Kadavath, Mantas Mazeika, Akul Arora, Ethan Guo, Collin
742 Burns, Samir Puranik, Horace He, Dawn Song, and Jacob Steinhardt. Measuring coding challenge
743 competence with apps. *NeurIPS*, 2021a.

744

745 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,
746 and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *NeurIPS*,
747 2021b.

748 Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza
749 Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hen-
750 nigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy,
751 Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre.
752 Training compute-optimal large language models, 2022. URL [https://arxiv.org/abs/
753 2203.15556](https://arxiv.org/abs/2203.15556).

754 Sarah Holland, Ahmed Hosny, Sarah Newman, Joshua Joseph, and Kasia Chmielinski. The dataset
755 nutrition label: A framework to drive higher data quality standards, 2018. URL [https://
arxiv.org/abs/1805.03677](https://arxiv.org/abs/1805.03677).

756 Ben Hutchinson, Andrew Smart, Alex Hanna, Emily Denton, Christina Greer, Oddur Kjartansson,
757 Parker Barnes, and Margaret Mitchell. Towards accountability for machine learning datasets:
758 Practices from software engineering and infrastructure, 2021. URL [https://arxiv.org/
759 abs/2010.13561](https://arxiv.org/abs/2010.13561).

760 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child,
761 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language
762 models, 2020. URL <https://arxiv.org/abs/2001.08361>.

763 Marzena Karpinska, Nader Akoury, and Mohit Iyyer. The perils of using mechanical turk to evaluate
764 open-ended text generation, 2021. URL <https://arxiv.org/abs/2109.06835>.

765 Mathew Oldham Kevin Lee, Adi Gangidi. Building meta’s genai infrastructure.
766 [https://engineering.fb.com/2024/03/12/data-center-engineering/
767 building-metas-genai-infrastructure/](https://engineering.fb.com/2024/03/12/data-center-engineering/building-metas-genai-infrastructure/), 2024. Accessed: September 28, 2024.

768 Omar Khattab, Arnav Singhvi, Paridhi Maheshwari, Zhiyuan Zhang, Keshav Santhanam, Sri Vard-
769 hamanan, Saiful Haq, Ashutosh Sharma, Thomas T. Joshi, Hanna Moazam, Heather Miller, Matei
770 Zaharia, and Christopher Potts. Dspy: Compiling declarative language model calls into self-
771 improving pipelines. *arXiv preprint arXiv:2310.03714*, 2023.

772 Douwe Kiela, Max Bartolo, Yixin Nie, Divyansh Kaushik, Atticus Geiger, Zhengxuan Wu, Bertie
773 Vidgen, Grusha Prasad, Amanpreet Singh, Pratik Ringshia, Zhiyi Ma, Tristan Thrush, Sebastian
774 Riedel, Zeerak Waseem, Pontus Stenetorp, Robin Jia, Mohit Bansal, Christopher Potts, and Adina
775 Williams. Dynabench: Rethinking benchmarking in nlp, 2021. URL [https://arxiv.org/
776 abs/2104.14337](https://arxiv.org/abs/2104.14337).

777 Spencer Kimball. Microsoft, brookfield to develop more than 10.5 gi-
778 gawatts of renewable energy. [https://www.cnbc.com/2024/05/01/
779 microsoft-brookfield-to-develop-more-than-10point5-gigawatts-of-renewable-energy.
780 html](https://www.cnbc.com/2024/05/01/microsoft-brookfield-to-develop-more-than-10point5-gigawatts-of-renewable-energy.html), 2024. Accessed: September 28, 2024.

781 Md Tahmid Rahman Laskar, M Saiful Bari, Mizanur Rahman, Md Amran Hossen Bhuiyan, Shafiq
782 Joty, and Jimmy Xiangji Huang. A systematic study and comprehensive evaluation of chatgpt on
783 benchmark datasets, 2023.

784 Md Tahmid Rahman Laskar, Sawsan Alqahtani, M Saiful Bari, Mizanur Rahman, Mohammad Ab-
785 dullah Matin Khan, Haidar Khan, Israt Jahan, Amran Bhuiyan, Chee Wei Tan, Md Rizwan Parvez,
786 et al. A systematic survey and critical review on evaluating large language models: Challenges,
787 limitations, and recommendations. *arXiv preprint arXiv:2407.04069*, 2024.

788 Tianle Li, Anastasios Angelopoulos, and Wei-Lin Chiang. Does style matter? disentangling
789 style and substance in chatbot arena, Aug 2024. URL [https://lmsys.org/blog/
790 2024-08-28-style-control/](https://lmsys.org/blog/2024-08-28-style-control/).

791 Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. Fantastically ordered
792 prompts and where to find them: Overcoming few-shot prompt order sensitivity. In Smaranda
793 Muresan, Preslav Nakov, and Aline Villavicencio (eds.), *Proceedings of the 60th Annual Meet-
794 ing of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 8086–8098,
795 Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.
796 acl-long.556. URL <https://aclanthology.org/2022.acl-long.556>.

797 Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. A survey
798 on bias and fairness in machine learning. *ACM Comput. Surv.*, 54(6), July 2021. ISSN 0360-0300.
799 doi: 10.1145/3457607. URL <https://doi.org/10.1145/3457607>.

800 Ron Miller. Google launches a 9 exaflop cluster of cloud tpu v4 pods
801 into public preview. [https://techcrunch.com/2022/05/11/
802 google-launches-a-9-exaflop-cluster-of-cloud-tpu-v4-pods-into-public-preview/
803 2022](https://techcrunch.com/2022/05/11/google-launches-a-9-exaflop-cluster-of-cloud-tpu-v4-pods-into-public-preview/). Accessed: [Your Access Date].

804 Mistral. Au large, 2024. URL <https://mistral.ai/news/mistral-large/>.

-
- 810 OpenAI. Openai five defeats dota 2 world champions. [https://openai.com/blog/
811 openai-five-defeats-dota-2-world-champions/](https://openai.com/blog/openai-five-defeats-dota-2-world-champions/), April 2019. Accessed: 2024-
812 09-28.
- 813 OpenAI. Chatgpt: Optimizing language models for dialogue, 2022. URL [https://openai.
814 com/blog/chatgpt/](https://openai.com/blog/chatgpt/).
- 815
- 816 OpenAI, Ilge Akkaya, Marcin Andrychowicz, Maciek Chociej, Mateusz Litwin, Bob McGrew,
817 Arthur Petron, Alex Paino, Matthias Plappert, Glenn Powell, Raphael Ribas, Jonas Schneider,
818 Nikolas Tezak, Jerry Tworek, Peter Welinder, Lilian Weng, Qiming Yuan, Wojciech Zaremba,
819 and Lei Zhang. Solving rubik’s cube with a robot hand, 2019. URL [https://arxiv.org/
820 abs/1910.07113](https://arxiv.org/abs/1910.07113).
- 821 Krista Opsahl-Ong, Michael J Ryan, Josh Purtell, David Broman, Christopher Potts, Matei Zaharia,
822 and Omar Khattab. Optimizing instructions and demonstrations for multi-stage language model
823 programs, 2024. URL <https://arxiv.org/abs/2406.11695>.
- 824
- 825 Aitor Ormazabal, Che Zheng, Cyprien de Masson d’Autume, Dani Yogatama, Deyu Fu, Donovan
826 Ong, Eric Chen, Eugenie Lamprecht, Hai Pham, Isaac Ong, Kaloyan Aleksiev, Lei Li, Matthew
827 Henderson, Max Bain, Mikel Artetxe, Nishant Relan, Piotr Padlewski, Qi Liu, Ren Chen, Samuel
828 Phua, Yazheng Yang, Yi Tay, Yuqi Wang, Zhongkai Zhu, and Zhihui Xie. Reka core, flash, and
829 edge: A series of powerful multimodal language models, 2024.
- 830 Joon Sung Park, Joseph O’Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and
831 Michael S Bernstein. Generative agents: Interactive simulacra of human behavior. In *Proceedings
832 of the 36th annual acm symposium on user interface software and technology*, pp. 1–22, 2023.
- 833 Ryan Park, Rafael Rafailov, Stefano Ermon, and Chelsea Finn. Disentangling length from quality
834 in direct preference optimization, 2024. URL <https://arxiv.org/abs/2403.19159>.
- 835
- 836 Pouya Pezeshkpour and Estevam Hruschka. Large language models sensitivity to the order of op-
837 tions in multiple-choice questions, 2023. URL <https://arxiv.org/abs/2308.11483>.
- 838 Chengwei Qin, Aston Zhang, Zhuosheng Zhang, Jiaao Chen, Michihiro Yasunaga, and Diyi
839 Yang. Is chatgpt a general-purpose natural language processing task solver? *arXiv preprint
840 arXiv:2302.06476*, 2023.
- 841
- 842 David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Di-
843 rani, Julian Michael, and Samuel R Bowman. Gpqa: A graduate-level google-proof q&a bench-
844 mark. *arXiv preprint arXiv:2311.12022*, 2023.
- 845 Victor Sanh, Albert Webson, Colin Raffel, Stephen H. Bach, Lintang Sutawika, Zaid Alyafeai, An-
846 toine Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja, Manan Dey, M Saiful Bari, Canwen
847 Xu, Urmish Thakker, Shanya Sharma Sharma, Eliza Szczechla, Taewoon Kim, Gunjan Chh-
848 ablani, Nihal Nayak, Debajyoti Datta, Jonathan Chang, Mike Tian-Jian Jiang, Han Wang, Mat-
849 teo Manica, Sheng Shen, Zheng Xin Yong, Harshit Pandey, Rachel Bawden, Thomas Wang,
850 Trishala Neeraj, Jos Rozen, Abheesh Sharma, Andrea Santilli, Thibault Fevry, Jason Alan
851 Fries, Ryan Teehan, Tali Bers, Stella Biderman, Leo Gao, Thomas Wolf, and Alexander M.
852 Rush. Multitask prompted training enables zero-shot task generalization, 2022. URL [https:
853 //arxiv.org/abs/2110.08207](https://arxiv.org/abs/2110.08207).
- 854 Tal Schuster, Roei Schuster, Darsh J. Shah, and Regina Barzilay. The limitations of stylometry for
855 detecting machine-generated fake news. *Computational Linguistics*, 46(2):499–510, June 2020.
856 doi: 10.1162/coli.a.00380. URL <https://aclanthology.org/2020.cl-2.8>.
- 857
- 858 Avi Schwarzschild, Eitan Borgnia, Arjun Gupta, Arpit Bansal, Zeyad Emam, Furong Huang, Micah
859 Goldblum, and Tom Goldstein. Datasets for studying generalization from easy to hard examples,
860 2021.
- 861 Minghao Shao, Sofija Jancheska, Meet Udeshi, Brendan Dolan-Gavitt, Haoran Xi, Kimberly Milner,
862 Boyuan Chen, Max Yin, Siddharth Garg, Prashanth Krishnamurthy, Farshad Khorrami, Ramesh
863 Karri, and Muhammad Shafique. Nyu ctf dataset: A scalable open-source benchmark dataset for
evaluating llms in offensive security, 2024. URL <https://arxiv.org/abs/2406.05590>.

864 David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche,
865 Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering
866 the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.
867

868 David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez,
869 Marc Lanctot, Laurent Sifre, Dhharshan Kumaran, Thore Graepel, Timothy Lillicrap, Karen Si-
870 monyan, and Demis Hassabis. Mastering chess and shogi by self-play with a general reinforce-
871 ment learning algorithm, 2017. URL <https://arxiv.org/abs/1712.01815>.

872 Arnav Singhvi, Manish Shetty, Shangyin Tan, Christopher Potts, Koushik Sen, Matei Zaharia, and
873 Omar Khattab. Dspy assertions: Computational constraints for self-refining language model
874 pipelines, 2024. URL <https://arxiv.org/abs/2312.13382>.

875

876 The Stockfish Developers. Stockfish, 2024. URL [https://github.com/](https://github.com/official-stockfish/Stockfish)
877 [official-stockfish/Stockfish](https://github.com/official-stockfish/Stockfish). Version 17.
878

879 Oguzhan Topsakal, Colby Jacob Edell, and Jackson Bailey Harper. Evaluating large language mod-
880 els with grid-based game competitions: An extensible llm benchmark and leaderboard, 2024.
881 URL <https://arxiv.org/abs/2407.07796>.

882 Tu Vu, Mohit Iyyer, Xuezhi Wang, Noah Constant, Jerry Wei, Jason Wei, Chris Tar, Yun-Hsuan
883 Sung, Denny Zhou, Quoc Le, et al. Freshllms: Refreshing large language models with search
884 engine augmentation. *arXiv preprint arXiv:2310.03214*, 2023.
885

886 Alex Wang. Glue: A multi-task benchmark and analysis platform for natural language understand-
887 ing. *arXiv preprint arXiv:1804.07461*, 2018.
888

889 Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer
890 Levy, and Samuel Bowman. Superglue: A stickier benchmark for general-purpose language
891 understanding systems. *Advances in neural information processing systems*, 32, 2019.

892 Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan,
893 and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models,
894 2023. URL <https://arxiv.org/abs/2305.16291>.
895

896 Haochun Wang, Sendong Zhao, Zewen Qiang, Nuwa Xi, Bing Qin, and Ting Liu. Beyond the
897 answers: Reviewing the rationality of multiple choice question answering for the evaluation of
898 large language models, 2024a. URL <https://arxiv.org/abs/2402.01349>.

899

900 Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming
901 Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, et al. Mmlu-pro: A more robust and challenging
902 multi-task language understanding benchmark. *arXiv preprint arXiv:2406.01574*, 2024b.

903 Alex Warstadt, Aaron Mueller, Leshem Choshen, Ethan Wilcox, Chengxu Zhuang, Juan Ciro,
904 Rafael Mosquera, Bhargavi Paranjabe, Adina Williams, Tal Linzen, and Ryan Cotterell (eds.).
905 *Proceedings of the BabyLM Challenge at the 27th Conference on Computational Natural Lan-
906 guage Learning*, Singapore, December 2023. Association for Computational Linguistics. URL
907 <https://aclanthology.org/2023.conll-babylm.0>.

908

909 Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yo-
910 gatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol
911 Vinyals, Percy Liang, Jeff Dean, and William Fedus. Emergent abilities of large language models,
912 2022. URL <https://arxiv.org/abs/2206.07682>.

913 Jason Wei, Najoung Kim, Yi Tay, and Quoc V. Le. Inverse scaling can become u-shaped, 2023a.
914 URL <https://arxiv.org/abs/2211.02011>.

915

916 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc
917 Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models,
2023b. URL <https://arxiv.org/abs/2201.11903>.

918 Lionel Wong, Gabriel Grand, Alexander K. Lew, Noah D. Goodman, Vikash K. Mansinghka, Jacob
919 Andreas, and Joshua B. Tenenbaum. From word models to world models: Translating from
920 natural language to the probabilistic language of thought, 2023. URL <https://arxiv.org/abs/2306.12672>.
921
922 Yuxiang Wu, Zhengyao Jiang, Akbir Khan, Yao Fu, Laura Ruis, Edward Grefenstette, and Tim
923 Rocktäschel. Chatarena: Multi-agent language game environments for large language models.
924 <https://github.com/chatarena/chatarena>, 2023.
925
926 An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li,
927 Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang,
928 Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jin Xu, Jingren Zhou, Jinze Bai,
929 Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng
930 Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai
931 Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan
932 Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Yang Fan, Yang Yao, Yichang
933 Zhang, Yu Wan, Yunfei Chu, Zeyu Cui, Zhenru Zhang, and Zhihao Fan. Qwen2 technical report,
934 2024. URL <https://arxiv.org/abs/2407.10671>.
935
936 Shuo Yang, Wei-Lin Chiang, Lianmin Zheng, Joseph E. Gonzalez, and Ion Stoica. Rethinking
937 benchmark and contamination for language models with rephrased samples, 2023.
938
939 Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min,
940 Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen,
941 Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and
942 Ji-Rong Wen. A survey of large language models, 2024. URL <https://arxiv.org/abs/2303.18223>.
943
944 Chujie Zheng, Hao Zhou, Fandong Meng, Jie Zhou, and Minlie Huang. Large language models are
945 not robust multiple choice selectors, 2024. URL <https://arxiv.org/abs/2309.03882>.
946
947 Stephan Zheng, Alexander Trott, Sunil Srinivasa, David C. Parkes, and Richard Socher. The ai
948 economist: Optimal economic policy design via two-level deep reinforcement learning, 2021.
949 URL <https://arxiv.org/abs/2108.02755>.

950 A APPENDIX

951
952
953 Table 4: Exact model versions used in our evaluations.

954 Model	954 Version
955 GPT-4o	955 gpt-4o-2024-08-06
956 Claude 3.5 Sonnet	956 claude-3-5-sonnet-20240620
957 Mistral-Large	957 mistralai/Mistral-Large-Instruct-2407
958 Llama 3.1 70B	958 meta-llama/Meta-Llama-3.1-70B-Instruct

959
960
961
962
963
964
965
966
967
968
969
970
971