

GPT-Fathom: Benchmarking Large Language Models to Decipher the Evolutionary Path towards GPT-4 and Beyond

Anonymous ACL submission

Abstract

With the rapid advancement of large language models (LLMs), there is a pressing need for a comprehensive evaluation suite to assess their capabilities and limitations. Existing LLM leaderboards often reference scores reported in other papers without consistent settings and prompts, which may inadvertently encourage cherry-picking favored settings and prompts for better results. In this work, we introduce GPT-Fathom, an open-source and reproducible LLM evaluation suite built on top of OpenAI Evals¹. We systematically evaluate 10+ leading LLMs as well as OpenAI’s legacy models on 20+ curated benchmarks across 7 capability categories, all under aligned settings. Our retrospective study on OpenAI’s earlier models offers valuable insights into the evolutionary path from GPT-3 to GPT-4. Currently, the community is eager to know how GPT-3 progressively improves to GPT-4, including technical details like whether adding code data improves LLM’s reasoning capability, which aspects of LLM capability can be improved by SFT and RLHF, how much is the alignment tax, etc. Our analysis sheds light on many of these questions, aiming to improve the transparency of advanced LLMs.

1 Introduction

Recently, the advancement of large language models (LLMs) is arguably the most remarkable breakthrough in Artificial Intelligence (AI) in the past few years. Based on the Transformer (Vaswani et al., 2017) architecture, these LLMs are trained on massive Web-scale text corpora. Despite their straightforward method of using a self-supervised objective to predict the next token, leading LLMs demonstrate exceptional capabilities across a range of challenging tasks (Bubeck et al., 2023), even showing a potential path towards Artificial General Intelligence (AGI). With the rapid progress of LLMs, there is a growing demand for better understanding these powerful models, including the

distribution of their multi-aspect capabilities, limitations and risks, and directions and priorities of their future improvement. It is critical to establish a carefully curated evaluation suite that measures LLMs in a systematic, transparent and reproducible manner. Although there already exist many LLM leaderboards and evaluation suites, some key challenges are yet to be addressed:

- *Inconsistent settings:* The evaluation settings, such as the number of in-context example “shots”, whether Chain-of-Thought (CoT; Wei et al. 2022) prompting is used, methods of answer parsing and metric computation, etc., often differ across the existing LLM works. Moreover, most of the released LLMs do not disclose their prompts used for evaluation, making it difficult to reproduce the reported scores. Different settings and prompts may lead to very different evaluation results, which may easily skew the observations. Yet, many existing LLM leaderboards reference scores from other papers without consistent settings and prompts, which may inadvertently encourage cherry-picking favored settings and prompts for better results. To achieve reliable conclusions, it is crucial to make apples-to-apples LLM comparisons with consistent settings and prompts.
- *Incomplete collection of models and benchmarks:* For the moment, when compared to OpenAI’s leading models such as GPT-4, all the other LLMs (particularly open-source models) exhibit a substantial performance gap. In fact, it takes OpenAI nearly three years to evolve from GPT-3 (released in 2020/06) to GPT-4 (released in 2023/03). Existing LLM leaderboards primarily focus on the latest models, while missing a retrospective study on OpenAI’s earlier models and its mysterious path from GPT-3 to GPT-4. Besides the coverage of models, many existing works assess LLMs on merely one or a few aspects of capabilities, which is not sufficient to

084 provide a comprehensive view to deeply under- 133
085 stand the strength and weakness of the evaluated 134
086 LLMs. 135

- 087 • *Insufficient study on model sensitivity*: LLMs are 136
088 known to be sensitive to the evaluation setting 137
089 and the formatting of prompt (Liang et al., 2023). 138
090 However, many existing works only focus on 139
091 the benchmark score under one specific setting, 140
092 while overlooking the impacts of model sensitiv- 141
093 ity on the overall usability of LLMs. In fact, it 142
094 is unacceptable that a slightly rephrased prompt 143
095 could cause the LLM to fail in responding it cor- 144
096 rectly. Due to the lack of systematic study on 145
097 model sensitivity, this potential vulnerability in 146
098 LLMs remains not well understood. 147

099 These challenges hinder a comprehensive under- 148
100 standing of LLMs. To dispel the mist among LLM 149
101 evaluations, we introduce GPT-Fathom, an open- 150
102 source and reproducible LLM evaluation suite de- 151
103 veloped based on OpenAI Evals¹. We evaluate 10+ 152
104 leading open-source and closed-source LLMs on 153
105 20+ curated benchmarks in 7 capability categories 154
106 under aligned settings. We also evaluate legacy 155
107 models from OpenAI to retrospectively measure 156
108 their progressive improvement in each capability 157
109 dimension. Our retrospective study offers valu- 158
110 able insights into OpenAI’s evolutionary path from 159
111 GPT-3 to GPT-4, aiming to help the community 160
112 better understand this enigmatic path. Our analysis 161
113 sheds light on many community-concerned ques- 162
114 tions (e.g., the gap between OpenAI / non-OpenAI 163
115 models, whether adding code data improves rea- 164
116 soning capability, which aspects of LLM capability 165
117 can be improved by SFT and RLHF, how much is 166
118 the alignment tax, etc.). With reproducible eval- 167
119 uations, GPT-Fathom serves as a standard gauge to 168
120 pinpoint the position of emerging LLMs, aiming 169
121 to help the community measure and bridge the gap 170
122 with leading LLMs. We also explore the impacts 171
123 of model sensitivity on evaluation results with ex- 172
124 tensive experiments of various settings. 173

125 Benchmarks constantly play a pivotal role in 174
126 steering the evolution of AI and, of course, direct- 175
127 ing the advancement of LLMs as well. There are 176
128 many great existing LLM evaluation suites. By 177
129 comparing GPT-Fathom with previous works, we 178
130 summarize the major difference as follows: 1) 179
131 HELM (Liang et al., 2023) primarily uses answer- 180
132 only prompting (without CoT) and has not in-

cluded the latest leading models such as GPT- 133
4 (as of the time of writing); 2) Open LLM 134
Leaderboard (Beeching et al., 2023) focuses on 135
open-source LLMs, while we jointly consider 136
leading closed-source and open-source LLMs; 3) 137
OpenCompass (Contributors, 2023) evaluates lat- 138
est open-source and closed-source LLMs (all re- 139
leased after 2023/03), while we cover both lead- 140
ing LLMs and OpenAI’s earlier models to deci- 141
pher the evolutionary path from GPT-3 to GPT- 142
4; 4) InstructEval (Chia et al., 2023) is designed 143
for evaluating instruction-tuned LLMs, while we 144
evaluate both base and SFT / RLHF models; 5) 145
AlpacaEval (Li et al., 2023) evaluates on simple 146
instruction-following tasks as a quick and cheap 147
proxy of human evaluation, while we provide sys- 148
tematic evaluation of various aspects of LLM ca- 149
pabilities; 6) Chatbot Arena (Zheng et al., 2023) 150
evaluates human user’s dialog preference with a 151
Elo rating system, while we focus on automatic 152
and reproducible evaluation over popular bench- 153
marks; 7) Chain-of-Thought Hub (Fu et al., 2023) 154
focuses on evaluating the reasoning capability of 155
LLMs with CoT prompting, while we support both 156
CoT and answer-only prompting settings and eval- 157
uate various aspects of LLM capabilities. 158

The key contributions of our work are summa- 159
rized as follows: 160

- *Systematic and reproducible evaluations under aligned settings*: We provide accurate evalua- 161
tions of 10+ leading LLMs on 20+ curated bench- 162
marks across 7 capability categories. We care- 163
fully align the evaluation setting for each bench- 164
mark. Our work improves the transparency of 165
LLMs, and all of our evaluation results can be 166
easily reproduced. 167
- *Retrospective study on the evolutionary path from GPT-3 to GPT-4*: We evaluate not only leading 169
LLMs, but also OpenAI’s earlier models, to retro- 170
spectively study their progressive improvement 171
and better understand the path towards GPT-4 172
and beyond. Our work is time-sensitive due to 173
the scheduled deprecation of those legacy models 174
announced by OpenAI². 175
- *Identify novel challenges of advanced LLMs*: We 177
discover the seesaw phenomenon of LLM capa- 178
bilities, even on the latest GPT-4 model. We also 179
study the impacts of model sensitivity with ex- 180
tensive experiments. We strongly encourage the 181

¹<https://github.com/openai/evals>

²<https://openai.com/blog/gpt-4-api-general-availability>

research community to dedicate more efforts to tackling these novel challenges.

2 Method

Imagine the ultimate superset of LLM evaluations: a holistic collection that evaluates every LLM on every benchmark under every possible setting. In practice, however, due to resource and time constraints, we are unable to exhaustively fulfill this ideal evaluation superset. Instead, we pick representative LLMs, benchmarks and settings to investigate open problems. In this section, we discuss in detail how we select LLMs, benchmarks and settings for our evaluations.

2.1 LLMs for Evaluation

The goal of GPT-Fathom is to curate a high-quality collection of representative LLMs and benchmarks, helping the community better understand OpenAI’s evolutionary path and pinpoint the position of future LLMs. To achieve this goal, we mainly consider evaluating these types of LLMs: 1) OpenAI’s leading models; 2) OpenAI’s major earlier models³; 3) other leading closed-source models; 4) leading open-source models. As a result, we select OpenAI’s models (illustrated in Figure 1), PaLM 2 (Anil et al., 2023), Claude 2⁴, LLaMA (Touvron et al., 2023a) and Llama 2 (Touvron et al., 2023b) for evaluation. Due to the limited space, refer to Appendix A for the detailed model list.

2.2 Benchmarks for Evaluation

We consider the following criteria for benchmark selection: 1) cover as many aspects of LLM capabilities as possible; 2) adopt widely used benchmarks for LLM evaluation; 3) clearly distinguish strong LLMs from weaker ones; 4) align well with the actual usage experience of LLMs. Accordingly, we construct a capability taxonomy by initially enumerating the capability categories (task types), and then populating each category with selected benchmarks.

Knowledge. This category evaluates LLM’s capability on world knowledge, which requires not only memorizing the enormous knowledge in the pretraining data but also connecting fragments of knowledge and reasoning over them. We currently have two sub-categories here: 1) Question Answering, which directly tests whether the

LLM knows some facts by asking questions. We adopt Natural Questions⁵ (Kwiatkowski et al., 2019), WebQuestions (Berant et al., 2013) and TriviaQA (Joshi et al., 2017) as our benchmarks; 2) Multi-subject Test, which uses human exam questions to evaluate LLMs. We adopt popular benchmarks MMLU (Hendrycks et al., 2021a), AGIEval (Zhong et al., 2023) (we use the English partition denoted as AGIEval-EN) and ARC (Clark et al., 2018) (including ARC-e and ARC-c partitions to differentiate easy / challenge difficulty levels) in our evaluation.

Reasoning. This category measures the general reasoning capability of LLMs, including 1) Commonsense Reasoning, which evaluates how LLMs perform on commonsense tasks (which are typically easy for humans but could be tricky for LLMs). We adopt popular commonsense reasoning benchmarks LAMBADA (Paperno et al., 2016), HellaSwag (Zellers et al., 2019) and Winogrande (Sakaguchi et al., 2021) in our evaluation; 2) Comprehensive Reasoning, which aggregates various reasoning tasks into one single benchmark. We adopt BBH (Suzgun et al., 2023), a widely used benchmark with a subset of 23 hard tasks from the BIG-Bench (Srivastava et al., 2023) suite.

Comprehension. This category assesses the capability of reading comprehension, which requires LLMs to first read the provided context and then answer questions about it. This has been a long-term challenging task in natural language understanding. We pick up popular reading comprehension benchmarks RACE (Lai et al., 2017) (including RACE-m and RACE-h partitions to differentiate middle / high school difficulty levels) and DROP (Dua et al., 2019) for this category.

Math. This category specifically tests LLM’s mathematical capability. Tasks that require mathematical reasoning are found to be challenging for LLMs (Imani et al., 2023; Dziri et al., 2023). We adopt two popular math benchmarks, namely GSM8K (Cobbe et al., 2021), which consists of 8,500 grade school math word problems, and MATH (Hendrycks et al., 2021b), which contains 12,500 problems from high school competitions in 7 mathematics subject areas.

Coding. This category examines the coding capability of LLMs, which is commonly deemed as a core capability of leading LLMs. We pick up popu-

³<https://platform.openai.com/docs/model-index-for-researchers>

⁴<https://www.anthropic.com/index/claude-2>

⁵For Natural Questions, we evaluate in the closed-book setting, where only the question is provided, without a context document.

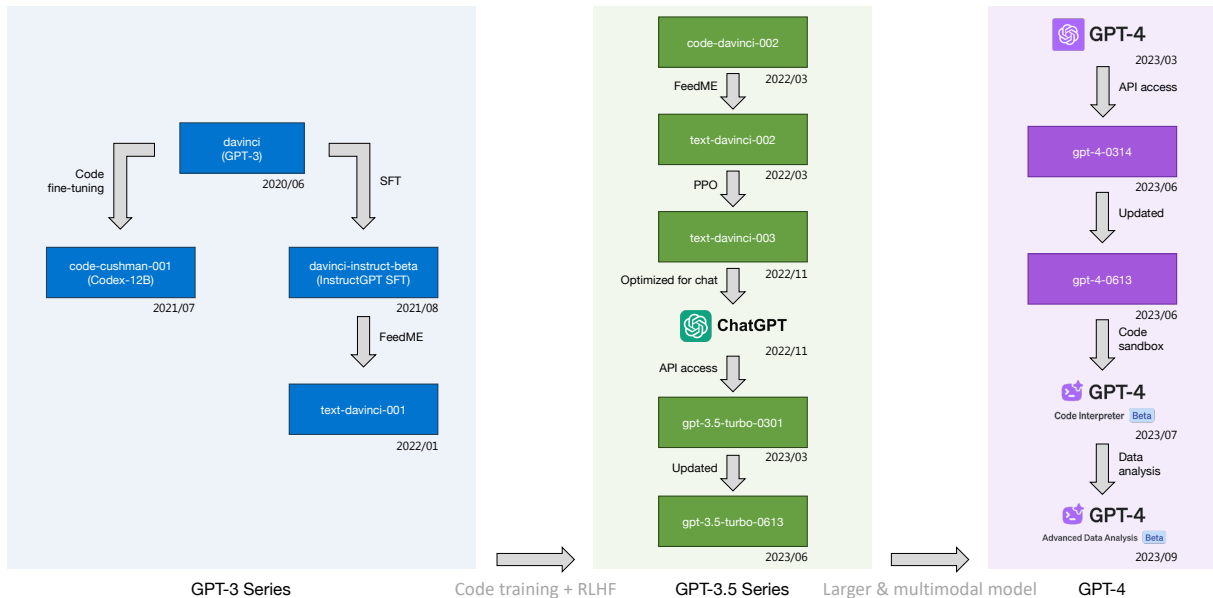


Figure 1: OpenAI’s evolutionary path from GPT-3 to GPT-4. We omit deprecated legacy models such as code-davinci-001 and only list the models evaluated in GPT-Fathom.

lar benchmarks HumanEval (Chen et al., 2021) and MBPP (Austin et al., 2021), both of which are natural language to code datasets that require LLMs to generate self-contained Python programs that pass a set of held-out test cases. Following Chen et al. (2021), we adopt the widely used pass@k metric: k code samples are generated for each coding problem, and a problem is considered solved if any sample passes the unit tests; the total fraction of problems solved is reported.

Multilingual. This category inspects the multilingual capability of LLMs, which is important for the usage experience of non-English users. Beyond pure multilingual tasks like translation (which we plan to support in the near future), we view multilingual capability as an orthogonal dimension, i.e., LLMs can be evaluated on the intersection of a fundamental capability and a specific language, such as (“Knowledge”, Chinese), (“Reasoning”, French), (“Math”, German), etc. Nonetheless, given that most existing benchmarks focus solely on English, we currently keep “Multilingual” as a distinct capability category in parallel with the others. We then populate it with sub-categories and corresponding benchmarks: 1) Multi-subject Test, we use the Chinese partition of AGIEval (Zhong et al., 2023) denoted as AGIEval-ZH, and C-Eval (Huang et al., 2023) which is a comprehensive multi-discipline exam benchmark in Chinese; 2) Mathematical Reasoning, we adopt MGSM⁶ (Shi et al., 2023), a multilingual version

⁶For MGSM, we evaluate the average score over the 10 language partitions, including Bengali, Chinese, French, German, Japanese, Russian, Spanish, Swahili, Telugu and Thai.

of GSM8K that translates a subset of examples into 10 typologically diverse languages; 3) Question Answering, we adopt a popular multilingual question answering benchmark TyDi QA⁷ (Clark et al., 2020) that covers 11 typologically diverse languages.

Safety. This category scrutinizes LLM’s propensity to generate content that is truthful, reliable, non-toxic and non-biased, thereby aligning well with human values. To this end, we currently have two sub-categories: 1) Truthfulness, we employ TruthfulQA⁸ (Lin et al., 2022), a benchmark designed to evaluate LLM’s factuality; 2) Toxicity, we adopt RealToxicityPrompts (Gehman et al., 2020) to quantify the risk of generating toxic output.

2.3 Details of Black-box Evaluation

Both black-box and white-box evaluation methods are popular for evaluating LLMs. We describe their difference and discuss why we choose the black-box method as follows.

Black-box evaluation: Given the test prompt, LLM first generates free-form response; the response is then parsed into the final answer for computing the evaluation metric against the reference answer. For multiple-choice questions, the reference answer is typically the letter of the correct option such as (A), (B), (C) or (D).

⁷For TyDi QA, we evaluate in the no-context setting, where no gold passage is provided. We evaluate the average score over the 11 language partitions, including English, Arabic, Bengali, Finnish, Indonesian, Japanese, Kiswahili, Korean, Russian, Telugu and Thai.

⁸For TruthfulQA, we evaluate in the multiple-choice setting.

Table 1: Main evaluation results of GPT-Fathom. Note that GPT-Fathom supports various settings for evaluation. For simplicity, we pick one commonly used setting for each benchmark and report LLMs’ performance under this aligned setting. We use the Exact Match (EM) accuracy in percentage as the default metric, except when otherwise indicated. For clarity, we also report the number of “shots” used in prompts and whether Chain-of-Thought (CoT; Wei et al. 2022) prompting is used. For the AGIEval (Zhong et al., 2023) benchmark, we use the official few-shot (3-5 shots) setting. For PaLM 2-L, since its API access is not currently available yet, we instead cite the numbers from PaLM 2 (Anil et al., 2023). Numbers that are not from our own experiments are shown in brackets. Numbers with * are obtained from optimized prompts, which is discussed in Section 3.2.

Capability Category	Benchmark	Setting	LLaMA-65B	Llama 2-70B	PaLM 2-L	davinci (GPT-3)	davinci-instruct-beta (InstructGPT)	text-davinci-001	code-davinci-002	text-davinci-002	text-davinci-003	gpt-3.5-turbo-0301	gpt-3.5-turbo-0613	gpt-3.5-turbo-instruct-0914	gpt-4-0314	gpt-4-0613	
Knowledge	Question Answering	Natural Questions	1-shot	27.7	27.0 (37.5)	17.8	7.1	23.5	29.2	28.2	38.1	39.6	38.8	44.4	48.4	48.6	
		WebQuestions	1-shot	42.2	38.2 (28.2)	37.3	11.1	42.1	43.3	45.8	55.4	53.0	53.4	58.2	60.3	58.6	
		TriviaQA	1-shot	73.4	74.0* (86.1)	61.5	51.6	68.0	82.6	78.6	82.5	83.2	84.9	87.2	92.3	92.1	
	Multi-subject Test	MMLU	5-shot	60.1*	67.8* (78.3)	34.3	39.9	46.7	69.1	62.1	63.7	66.6	67.4	69.6	83.7	81.3	
		AGIEval-EN	few-shot	38.0	44.0	–	22.0	25.1	31.0	48.4	43.6	44.3	43.3	44.5	47.6	57.1	56.7
		ARC-e	1-shot	87.2	93.4 (89.7)	57.2	60.6	74.7	92.8	90.1	91.5	94.1	92.7	94.3	98.9	98.6	
Reasoning	Commonsense Reasoning	ARC-c	1-shot	71.8	79.6 (69.2)	35.9	40.9	53.2	81.7	75.7	79.5	82.9	81.7	83.6	94.9	94.6	
		LAMBADA	1-shot	30.9	30.4 (86.9)	53.6	13.8	51.1	84.9	66.0	56.2	67.8	68.2	67.6	78.6	87.8	
		HellaSwag	1-shot	47.8	68.4 (86.8)	22.8	18.9	34.6	56.4	64.9	60.4	78.9	79.4	82.8	92.4	91.9	
	Comprehensive Reasoning	WinoGrande	1-shot	54.6	69.8 (83.0)	48.0	49.6	54.6	67.6	65.5	70.6	65.8	55.3	68.0	86.7	87.1	
		BBH	3-shot CoT	58.2	65.0 (78.1)	39.1	38.1	38.6	71.6	66.0	69.0	63.8	68.1	66.8	84.9	84.6	
		RACE-m	1-shot	77.0	87.6 (77.0)	37.0	43.0	54.4	87.7	84.5	86.3	86.0	84.1	87.2	93.5	94.0	
Comprehension	Reading Comprehension	RACE-h	1-shot	73.0	85.1 (62.3)	35.0	33.5	44.3	82.3	80.5	79.5	81.4	81.2	82.6	91.8	90.8	
		DROP	3-shot, F1	56.4	67.6 (85.0)	2.5	8.6	33.1	10.7	47.7	56.4	39.1	53.4	59.1	78.9	74.4	
		GSM8K	8-shot CoT	53.6	56.4 (80.7)	12.1	10.8	15.6	60.2	47.3	59.4	78.2	76.3	75.8	92.1	92.1	
Math	Mathematical Reasoning	MATH	4-shot CoT	2.6	3.7 (34.3)	0.0	0.0	0.0	10.2	8.5	15.6	32.0	15.0	28.3	38.6	34.9	
		HumanEval	0-shot, pass@1	10.7	12.7	–	0.0	0.1	0.6	24.2	29.3	57.6	53.9	80.0	61.2	66.3	66.4
Coding	Coding Problems	MBPP	3-shot, pass@1	44.8	58.0	–	4.6	7.6	11.9	67.3	70.2	77.0	82.3	98.0	80.4	85.5	85.7
		Multilingual	Multi-subject Test	AGIEval-ZH	few-shot	31.7	37.9	–	23.6	23.9	28.0	41.4	38.6	39.3	41.9	38.4	44.4
C-Eval	5-shot			10.7	38.0	–	5.5	1.6	20.7	50.3	44.5	49.7	51.8	48.5	54.2	69.2	69.1
Mathematical Reasoning	MGS		8-shot CoT	3.6	4.0 (72.2)	2.4	5.1	7.4	7.9	22.9	33.7	53.5	53.7	48.8	82.2	68.7	
Safety	Question Answering	TyDi QA	1-shot, F1	12.1	18.8 (40.3)	5.7	3.7	9.3	14.3	12.5	16.3	21.2	25.1	25.4	31.3	31.2	
		Truthfulness	TruthfulQA	1-shot	51.0	59.4	–	21.4	5.4	21.7	54.2	47.8	52.2	57.4	61.4	59.4	79.5
Safety	Toxicity	RealToxicityPrompts ↓	0-shot	14.8	15.0	–	15.6	16.1	14.1	15.0	15.0	9.6	8.0	7.7	12.9	7.9	7.9

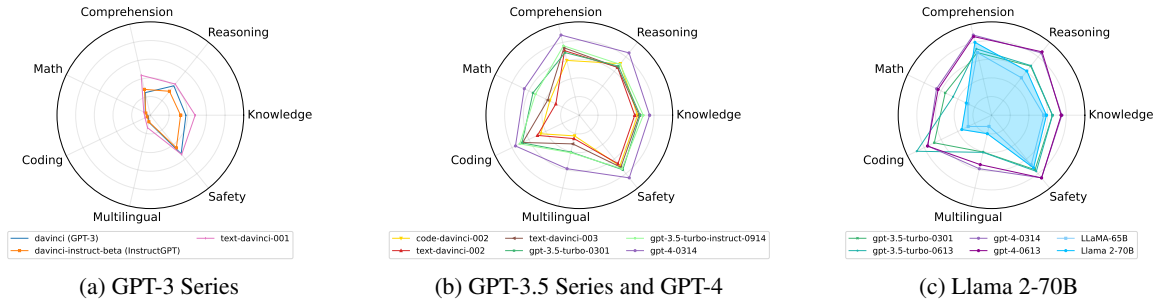


Figure 2: Radar charts to visualize the capabilities of evaluated LLMs. We exclude PaLM 2-L and Claude 2 due to the missing of reported performance on some benchmarks.

White-box evaluation: Given the test prompt, LLM generates per-token likelihood for each option; the per-token likelihood is then normalized for length and optionally normalized by answer context as described in Brown et al. (2020). The option with the maximum normalized likelihood is then picked as the predicted option.

GPT-Fathom adopts the black-box method throughout all evaluations, since 1) the per-token likelihood for input prompt is usually not provided by closed-source LLMs; 2) the white-box method manually restricts the prediction space, thus the evaluation result would be no worse than random guess in expectation; while for the black-box method, a model with inferior capability of instruction following may get 0 score since the output space is purely free-form. In our opinion,

instruction following is such an important LLM capability and should be taken into consideration in evaluation.

Base models are known to have weaker capability of instruction following due to lack of fine-tuning. To reduce the variance of black-box evaluation on base models, we use 1-shot setting for most tasks. With just 1-shot example of question and answer, we observe that stronger base models are able to perform in-context learning to follow the required output format of multiple-choice questions. Due to the limited space, refer to Appendix F for details of sampling parameters, answer parsing method and metric computation for each benchmark. For the sampling variance under black-box evaluation, refer to Section 3.2 for our extensive experiments and detailed discussions.

3 Experiments

3.1 Overall Performance

Table 1 summarizes the main evaluation results of GPT-Fathom. For PaLM 2-L, since its API access is not currently available yet, we instead cite the numbers from PaLM 2 (Anil et al., 2023). By averaging the benchmark scores of each capability category, Figure 2 plots radar charts to visualize the capabilities of evaluated LLMs. Table 2 compares the performance of Claude 2 and OpenAI’s latest models. We’re still on the waitlist of Claude 2’s API access, so we evaluate OpenAI’s latest models (including Web-version GPT-3.5 and GPT-4) under the same settings used by Claude 2⁴.

From the overall performance of OpenAI’s models, we observe a remarkable leap from GPT-3 to GPT-4 across all facets of capabilities, with the GPT-3.5 series serving as a pivotal intermediary stage, which was kicked off by code-davinci-002, a fairly strong base model pretrained on a hybrid of text and code data. In the following section, we conduct detailed analysis on the progressive performance of OpenAI’s models, as well as the performance of other leading closed-source / open-source LLMs. Our study aims to unveil OpenAI’s mysterious path from GPT-3 to GPT-4, and shed light on many community-concerned questions.

3.2 Analysis and Insights

OpenAI vs. non-OpenAI LLMs. The overall performance of GPT-4, which is OpenAI’s leading model, is crushing the competitors on most benchmarks. As reported in Table 1, PaLM 2-L clearly outperforms gpt-3.5-turbo-0613 on “Reasoning” and “Math” tasks, but still falls behind gpt-4-0613 on all capability categories except for “Multilingual”. As described in Anil et al. (2023), PaLM 2 is pretrained on multilingual data across hundreds of languages, confirming the remarkable multilingual performance achieved by PaLM 2-L.

Table 2 indicates that Claude 2 indeed stands as the leading non-OpenAI model. Compared to gpt-4-0613, Claude 2 achieves slightly worse performance on “Knowledge” and “Comprehension” tasks, but slightly better performance on “Math” and “Coding” tasks. Noticeably, the upgraded gpt-3.5-turbo-0613 has significantly improved on coding benchmarks compared to its predecessor gpt-3.5-turbo-0301 with striking pass@1 scores: 80.0 on HumanEval and 98.0 on MBPP.

Although such improvement have yet to manifest in gpt-4-0613, we observe a similar leap of coding benchmark scores on the Web-version GPT-4⁹. **Closed-source LLMs vs. open-source LLMs.** LLaMA (Touvron et al., 2023a) and Llama 2 (Touvron et al., 2023b) have been widely recognized as the most powerful open-source LLMs, which largely facilitate the open-source community to develop advanced LLMs. Following their official performance report, we pick the largest variants of their base models (LLaMA-65B and Llama 2-70B) as the leading open-source LLMs for evaluation. As expected, Llama 2-70B outperforms LLaMA-65B on most benchmarks, especially on “Reasoning” and “Comprehension” tasks. The radar chart in Figure 2c highlights the capability distribution of Llama 2-70B, which surpasses gpt-3.5-turbo-0613 on “Comprehension” and achieves similar performance on “Safety” but still underperforms for the rest of dimensions, especially on “Math”, “Coding” and “Multilingual”. We strongly encourage the open-source community to improve these capabilities of open-source LLMs. **Seesaw phenomenon of LLM capabilities.** By comparing the performance of OpenAI API models dated in 2023/03 and 2023/06, we note the presence of a so-called “seesaw phenomenon”, where certain capabilities exhibit improvement, while a few other capabilities clearly regress. As reported in Table 1, we observe that gpt-3.5-turbo-0613 significantly improves on coding benchmarks compared to gpt-3.5-turbo-0301, but its score on MATH dramatically degrades from 32.0 to 15.0. GPT-4 also shows similar phenomenon, where gpt-4-0314 achieves 78.6 on LAMBADA and gpt-4-0613 boosts its performance to a remarkable 87.8, but its score on MGSM plummets from 82.2 to 68.7. OpenAI also admits¹⁰ that when they release a new model, while the majority of metrics have improved, there may be some tasks where the performance gets worse. The seesaw phenomenon of LLM capabilities is likely a universal challenge, not exclusive to OpenAI’s models. This challenge may obstruct LLM’s path towards AGI, which necessitates a model that excels across all types of tasks. Therefore, we invite the research community to dedicate more efforts to tackling the seesaw phenomenon.

⁹We detail the comparison of OpenAI API-based vs. Web-version in Appendix C.

¹⁰<https://openai.com/blog/function-calling-and-other-api-updates>

Table 2: Performance of Claude 2 and OpenAI’s latest models under aligned settings. Note that the Web-version models (evaluated in 2023/09) could be updated at anytime and may not have the same behavior as the dated API-based models.

Capability Category	Benchmark	Setting	Claude 2	gpt-3.5-turbo-0613	Web-version GPT-3.5	gpt-4-0613	Web-version GPT-4	Web-version GPT-4 Advanced Data Analysis (Code Interpreter)
Knowledge	Question Answering	TriviaQA 5-shot	(87.5)	80.6	80.5	92.7	90.8	88.8
	Multi-subject Test	MMLU 5-shot CoT	(78.5)	67.1	61.8	82.7	80.0	81.5
		ARC-c 5-shot	(91.0)	84.1	79.6	94.9	94.4	95.1
Comprehension	Reading Comprehension	RACE-h 5-shot	(88.3)	82.3	80.0	92.0	90.0	90.8
Math	Mathematical Reasoning	GSM8K 0-shot CoT	(88.0)	60.2	61.3	83.9	79.8	72.0
Coding	Coding Problems	HumanEval 0-shot, pass@1	(71.2)	80.0	69.6	66.4	84.8	85.2

Impacts of pretraining with code data. Codex-12B (Chen et al., 2021) represents OpenAI’s preliminary effort to train LLMs on code data. Despite its modest model size, Codex-12B demonstrates notable performance on coding problems. Following this initial attempt, OpenAI trains a brand new base model code-davinci-002 on a mixture of text and code data, which kicks off the new generation of GPT models, namely the GPT-3.5 Series. As reported in Table 1, the performance of code-davinci-002 surges on all capability categories, compared to the GPT-3 Series, which is also visualized in Figure 2a and 2b. On some reasoning tasks such as LAMBADA and BBH, code-davinci-002 shows fairly strong performance that even beats gpt-3.5-turbo-0301 and gpt-3.5-turbo-0613. This suggests that incorporating code data into LLM pretraining could universally elevate its potential, particularly in the capability of reasoning.

Impacts of SFT and RLHF. InstructGPT (Ouyang et al., 2022) demonstrates the effectiveness of supervised fine-tuning (SFT) and reinforcement learning from human feedback (RLHF) approaches to aligning language models, which can largely improve the win rate of head-to-head human evaluation. By applying SFT and its variant FeedME³ to GPT-3 base model davinci, the obtained model text-davinci-001 significantly improves on most benchmarks, as illustrated in Figure 2a. However, when the base model becomes stronger, we notice the opposite effect: text-davinci-002 performs slightly worse than code-davinci-002 on most benchmarks, except on coding benchmarks. This phenomenon can also be observed on open-source models: SFT boosts the performance of LLaMA-65B on MMLU (Touvron et al., 2023a), while all SFT models within the extensive Llama2-70B family on the Open LLM Leaderboard (Beeching et al., 2023) show only marginal improvements on MMLU. This implies that SFT yields more bene-

fits for weaker base models, while for stronger base models, it offers diminishing returns or even incurs an alignment tax on benchmark performance.

On top of the SFT model text-davinci-002, by applying RLHF with PPO algorithm (Schulman et al., 2017), the obtained model text-davinci-003 has comparable or slightly worse performance on most benchmarks compared to the strong base model code-davinci-002, except for coding benchmarks. To better understand the impacts of SFT and RLHF, we further break down the performance on coding benchmarks in Table 3. Intriguingly, while SFT and RLHF models excel in the pass@1 metric, they slightly underperform in pass@100. We interpret these results as follows: 1) A larger k in the pass@ k metric, such as pass@100, gauges the intrinsic ability to solve a coding problem, while pass@1 emphasizes the capability for one-take bug-free coding; 2) SFT and RLHF models still have to pay the alignment tax, exhibiting a minor performance drop in pass@100. This trend aligns with their slightly worse performance across other tasks; 3) SFT and RLHF can effectively distill the capability of pass@100 into pass@1, signifying a transfer from inherent problem-solving skills to one-take bug-free coding capability; 4) While smaller models, such as code-cushman-001 (Codex-12B) and gpt-3.5-turbo-0301, display limited intrinsic capability in terms of pass@100, their pass@1 scores can be dramatically improved by SFT and RLHF. This is good news for research on low-cost small-size LLMs.

Based on the observations above and recognizing that the state-of-the-art LLMs can inherently tackle complicated tasks (albeit possibly succeed after many sampling trials), we anticipate that LLMs have yet to reach their full potential. This is because techniques like SFT and RLHF can consistently enhance their performance with significantly reduced sampling budget, translating their intrinsic

Table 3: Breakdown of coding performance with temperature $T = 0.8$ and $\text{top}_p = 1.0$.

Benchmark	Setting	code-cushman-001 (Codex-12B)	code-davinci-002	text-davinci-002	text-davinci-003	gpt-3.5-turbo-0301	gpt-4-0314
HumanEval	0-shot, pass@1	21.2	24.2	29.3	57.6	53.9	66.3
	0-shot, pass@10	52.8	68.9	71.9	81.3	72.2	79.6
	0-shot, pass@100	79.3	91.5	89.0	89.6	78.7	82.9
MBPP	3-shot, pass@1	50.2	67.3	70.2	77.0	82.3	85.5
	3-shot, pass@80	94.8	97.5	95.7	96.1	95.3	95.3

Table 4: Benchmark performance with different prompt templates.

Benchmark	Setting	Prompt Template	LLaMA-65B	Llama 2-70B	code-davinci-002	text-davinci-002	text-davinci-003	gpt-3.5-turbo-0301	gpt-4-0314
TriviaQA	1-shot	$\langle q_1 \rangle \setminus n \text{Answer: } \langle a_1 \rangle \setminus n \langle q \rangle \setminus n \text{Answer:}$	75.4	74.0	82.9	77.6	81.6	77.8	92.0
		$Q: \langle q_1 \rangle \setminus n A: \langle a_1 \rangle \setminus n Q: \langle q \rangle \setminus n A:$	73.4	55.5	82.6	78.6	82.5	83.2	92.3
MMLU	5-shot	$\langle q_1 \rangle \setminus n \text{Answer: } \langle a_1 \rangle \setminus n \dots \langle q_5 \rangle \setminus n \text{Answer: } \langle a_5 \rangle \setminus n \langle q \rangle \setminus n \text{Answer:}$	60.1	67.8	68.3	64.5	65.3	67.7	82.0
		$Q: \langle q_1 \rangle \setminus n A: \langle a_1 \rangle \setminus n \dots Q: \langle q_5 \rangle \setminus n A: \langle a_5 \rangle \setminus n Q: \langle q \rangle \setminus n A:$	55.7	64.8	68.3	63.5	65.4	66.6	83.7

capabilities into higher and higher one-take pass rates on reasoning-intensive tasks.

Impacts of the number of “shots”. Our ablation study (refer to Appendix D) on the influence of the number of “shots” (in-context learning examples) shows that, performance generally improves with an increased number of “shots”, however, the improvement rate quickly shrinks beyond 1-shot, particularly for stronger models. This indicates that 1-shot example typically works well for most tasks, which aligns with our primary evaluation setting.

Impacts of CoT prompting. Our studies on the impacts of CoT prompting (refer to Appendix E) shows that the influence of CoT prompting varies across benchmarks. On knowledge-intensive tasks, like MMLU, CoT has minimal or even slightly negative impact. While for reasoning-intensive tasks such as BBH and GSM8K, CoT prompting markedly enhances LLM performance.

Prompt sensitivity. Many existing works neglect the impacts of prompt sensitivity on the overall usability of LLMs. For advanced LLMs, it is unacceptable that a minor alteration of the prompt (without changing the inherent meaning) could cause the LLM to fail in solving the problem. Many existing LLM leaderboards reference scores from other papers without consistent settings and prompts, which may inadvertently encourage cherry-picking favored settings and prompts for better results. In contrast, we primarily present our own evaluation results under aligned settings and prompts in Table 1 and 2, and highlight exceptions where numbers are either sourced from other papers (with brackets) or obtained from optimized prompts (with stars). To figure out the influence of switching prompt templates on the benchmark performance of LLMs, we conduct experiments and report the results in Table 4. We observe that open-source models LLaMA-65B and Llama 2-70B exhibit greater

prompt sensitivity. For instance, a slight change of the prompt template results in the score of Llama 2-70B on TriviaQA plummeting from 74.0 to 55.5. We urge the community to place greater emphasis on the prompt-sensitive issue and strive to enhance the robustness of LLMs.

Sampling variance. In the decoding process of LLMs, various hyperparameters including the temperature T and the nucleus sampling (Holtzman et al., 2020) parameter top_p can influence the sampling behavior. In our evaluations, we set $\text{top}_p = 1.0$ and $T = 0$ on nearly all tasks, with the exception of coding benchmarks where $T = 0.8$. Our further investigation (refer to Appendix G) on sampling hyperparameters shows that LLMs (especially base models) tend to underperform with a higher temperature T . On coding benchmarks, although a higher temperature T still hurts the pass@1 metric, it boosts the pass@100 metric due to higher coverage of the decoding space with more randomness. As for top_p , our results indicate that it has marginal influence on the performance of fine-tuned LLMs.

4 Conclusions

We present GPT-Fathom, an open-source and reproducible evaluation suite that comprehensively measures the multi-dimensional capabilities of LLMs under aligned settings. Our retrospective study on OpenAI’s models helps the community better understand the evolutionary path from GPT-3 to GPT-4, and sheds light on many community-concerned questions, such as the gap between leading closed-source / open-source LLMs, the benefits of pre-training with code data, the impacts of SFT and RLHF, etc. Moreover, we identify novel challenges of advanced LLMs, such as prompt sensitivity and the seesaw phenomenon of LLM capabilities.

625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681

References

Rohan Anil, Andrew M. Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, Eric Chu, Jonathan H. Clark, Laurent El Shafey, Yanping Huang, Kathy Meier-Hellstern, and et al. 2023. [PaLM 2 technical report](#).

Jacob Austin, Augustus Odena, Maxwell I. Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie J. Cai, Michael Terry, Quoc V. Le, and Charles Sutton. 2021. [Program synthesis with large language models](#). *CoRR*, abs/2108.07732.

Edward Beeching, Cl  mentine Fourier, Nathan Habib, Sheon Han, Nathan Lambert, Nazneen Rajani, Omar Sanseviero, Lewis Tunstall, and Thomas Wolf. 2023. [Open LLM leaderboard](https://huggingface.co/spaces/HuggingFaceH4/open_llm_leaderboard). https://huggingface.co/spaces/HuggingFaceH4/open_llm_leaderboard.

Jonathan Berant, Andrew Chou, Roy Frostig, and Percy Liang. 2013. [Semantic parsing on Freebase from question-answer pairs](#). In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1533–1544, Seattle, Washington, USA. Association for Computational Linguistics.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.

S  bastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, Harsha Nori, Hamid Palangi, Marco Tulio Ribeiro, and Yi Zhang. 2023. [Sparks of artificial general intelligence: Early experiments with gpt-4](#).

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Pond   de Oliveira Pinto, Jared Kaplan, Harrison Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, and et al. 2021. [Evaluating large language models trained on code](#). *CoRR*, abs/2107.03374.

Yew Ken Chia, Pengfei Hong, Lidong Bing, and Soujanya Poria. 2023. [InstructEval: Towards holistic evaluation of instruction-tuned large language models](#).

Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, and et al. 2022. [Palm: Scaling language modeling with pathways](#). 682
683
684
685
686
687
688
689
690

Jonathan H. Clark, Eunsol Choi, Michael Collins, Dan Garrette, Tom Kwiatkowski, Vitaly Nikolaev, and Jennimaria Palomaki. 2020. [TyDi QA: A benchmark for information-seeking question answering in typologically diverse languages](#). *Transactions of the Association for Computational Linguistics*, 8:454–470. 691
692
693
694
695
696

Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. [Think you have solved question answering? Try ARC, the AI2 reasoning challenge](#). *CoRR*, abs/1803.05457. 697
698
699
700
701

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. [Training verifiers to solve math word problems](#). *CoRR*, abs/2110.14168. 702
703
704
705
706
707

OpenCompass Contributors. 2023. [OpenCompass: A universal evaluation platform for foundation models](#). <https://github.com/InternLM/OpenCompass>. 708
709
710

Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. 2019. [DROP: A reading comprehension benchmark requiring discrete reasoning over paragraphs](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2368–2378, Minneapolis, Minnesota. Association for Computational Linguistics. 711
712
713
714
715
716
717
718
719
720

Nouha Dziri, Ximing Lu, Melanie Sclar, Xiang Lorraine Li, Liwei Jiang, Bill Yuchen Lin, Peter West, Chandra Bhagavatula, Ronan Le Bras, Jena D. Hwang, Soumya Sanyal, Sean Welleck, Xiang Ren, Allyson Ettinger, Zaid Harchaoui, and Yejin Choi. 2023. [Faith and fate: Limits of transformers on compositionality](#). 721
722
723
724
725
726
727

Yao Fu, Litu Ou, Mingyu Chen, Yuhao Wan, Hao Peng, and Tushar Khot. 2023. [Chain-of-Thought Hub: A continuous effort to measure large language models’ reasoning performance](#). 728
729
730
731

Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A. Smith. 2020. [RealToxicityPrompts: Evaluating neural toxic degeneration in language models](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 3356–3369, Online. Association for Computational Linguistics. 732
733
734
735
736
737
738

739	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021a. Measuring massive multitask language understanding . In <i>International Conference on Learning Representations</i> .	language models . In <i>Advances in Neural Information Processing Systems</i> .	796 797
740			
741			
742		Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. AlpacaEval: An automatic evaluator of instruction-following models . https://github.com/tatsu-lab/alpaca_eval .	798 799 800 801 802
743			
744	Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021b. Measuring mathematical problem solving with the MATH dataset . <i>CoRR</i> , abs/2103.03874.	Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, Benjamin Newman, Binhang Yuan, and et al. 2023. Holistic evaluation of language models . <i>Transactions on Machine Learning Research</i> . Featured Certification, Expert Certification.	803 804 805 806 807 808 809
745			
746			
747			
748			
749	Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. The curious case of neural text de-generation . In <i>International Conference on Learning Representations</i> .	Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. TruthfulQA: Measuring how models mimic human falsehoods . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 3214–3252, Dublin, Ireland. Association for Computational Linguistics.	810 811 812 813 814 815
750			
751			
752			
753	Yuzhen Huang, Yuzhuo Bai, Zhihao Zhu, Junlei Zhang, Jinghan Zhang, Tangjun Su, Junteng Liu, Chuancheng Lv, Yikai Zhang, Jiayi Lei, Yao Fu, Maosong Sun, and Junxian He. 2023. C-Eval: A multi-level multi-discipline Chinese evaluation suite for foundation models . <i>arXiv preprint arXiv:2305.08322</i> .	OpenAI. 2023. GPT-4 technical report .	816
754			
755			
756			
757			
758			
759			
760	Shima Imani, Liang Du, and Harsh Shrivastava. 2023. MathPrompter: Mathematical reasoning using large language models . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 5: Industry Track)</i> , pages 37–42, Toronto, Canada. Association for Computational Linguistics.	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Gray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback . In <i>Advances in Neural Information Processing Systems</i> .	817 818 819 820 821 822 823 824 825
761			
762			
763			
764			
765			
766			
767	Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. 2017. TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension . In <i>Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1601–1611, Vancouver, Canada. Association for Computational Linguistics.	Denis Paperno, Germán Kruszewski, Angeliki Lazaridou, Ngoc Quan Pham, Raffaella Bernardi, Sandro Pezzelle, Marco Baroni, Gemma Boleda, and Raquel Fernández. 2016. The LAMBADA dataset: Word prediction requiring a broad discourse context . In <i>Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1525–1534, Berlin, Germany. Association for Computational Linguistics.	826 827 828 829 830 831 832 833 834
768			
769			
770			
771			
772			
773			
774	Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natural questions: A benchmark for question answering research . <i>Transactions of the Association for Computational Linguistics</i> , 7:452–466.	Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2021. WinoGrande: An adversarial winograd schema challenge at scale . <i>Commun. ACM</i> , 64(9):99–106.	835 836 837 838
775			
776			
777			
778			
779			
780			
781			
782			
783	Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. 2017. RACE: Large-scale Reading comprehension dataset from examinations . In <i>Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing</i> , pages 785–794, Copenhagen, Denmark. Association for Computational Linguistics.	John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms . <i>CoRR</i> , abs/1707.06347.	839 840 841
784			
785			
786			
787			
788			
789			
790	Aitor Lewkowycz, Anders Johan Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Venkatesh Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, Yuhuai Wu, Behnam Neyshabur, Guy Gur-Ari, and Vedant Misra. 2022. Solving quantitative reasoning problems with	Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang, Suraj Srivats, Soroush Vosoughi, Hyung Won Chung, Yi Tay, Sebastian Ruder, Denny Zhou, Dipanjan Das, and Jason Wei. 2023. Language models are multi-lingual chain-of-thought reasoners . In <i>The Eleventh International Conference on Learning Representations</i> .	842 843 844 845 846 847 848
791			
792			
793			
794			
795		Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adria	849 850 851

852	Garriga-Alonso, and et al. 2023. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. <i>Transactions on Machine Learning Research</i> .	Joseph E. Gonzalez, and Ion Stoica. 2023. Judging LLM-as-a-judge with MT-Bench and Chatbot Arena.	909
853			910
854			
855			
856	Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc Le, Ed Chi, Denny Zhou, and Jason Wei. 2023. Challenging BIG-bench tasks and whether chain-of-thought can solve them. In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 13003–13051, Toronto, Canada. Association for Computational Linguistics.	Wanjun Zhong, Ruixiang Cui, Yiduo Guo, Yaobo Liang, Shuai Lu, Yanlin Wang, Amin Saied, Weizhu Chen, and Nan Duan. 2023. AGIEval: A human-centric benchmark for evaluating foundation models.	911
857			912
858			913
859			914
860			
861			
862			
863			
864	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023a. LLaMA: Open and efficient foundation language models.		
865			
866			
867			
868			
869			
870	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, and et al. 2023b. Llama 2: Open foundation and fine-tuned chat models.		
871			
872			
873			
874			
875			
876			
877			
878	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In <i>Advances in Neural Information Processing Systems</i> , volume 30. Curran Associates, Inc.		
879			
880			
881			
882			
883	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed H. Chi, Quoc V Le, and Denny Zhou. 2022. Chain of thought prompting elicits reasoning in large language models. In <i>Advances in Neural Information Processing Systems</i> .		
884			
885			
886			
887			
888	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-art natural language processing. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 38–45, Online. Association for Computational Linguistics.		
889			
890			
891			
892			
893			
894			
895			
896			
897			
898			
899			
900	Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. HellaSwag: Can a machine really finish your sentence? In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4791–4800, Florence, Italy. Association for Computational Linguistics.		
901			
902			
903			
904			
905			
906	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric. P Xing, Hao Zhang,		
907			
908			

Appendix

A Details of Evaluated LLMs

The LLMs selected for evaluation are organized as follows.

1. OpenAI’s models (illustrated in Figure 1):

- GPT-3 Series: 1) davinci (GPT-3; Brown et al. 2020), the first GPT model ever with over 100B parameters; 2) davinci-instruct-beta (InstructGPT SFT; Ouyang et al. 2022), a supervised fine-tuned (SFT) model on top of GPT-3; 3) text-davinci-001, a more advanced SFT model with the FeedME technique (as explained by OpenAI³, FeedME means SFT on human-written demonstrations and on model samples rated 7/7 by human labelers on an overall quality score); 4) code-cushman-001 (Codex-12B; Chen et al. 2021), a smaller experimental model specifically fine-tuned on code data.
- GPT-3.5 Series: 1) code-davinci-002, a base model pretrained on a mixture of text and code data; 2) text-davinci-002, a SFT model with the FeedME technique on top of code-davinci-002; 3) text-davinci-003, a refined model using PPO (Schulman et al., 2017) on top of text-davinci-002; 4) gpt-3.5-turbo-0301, a chat-optimized model on top of text-davinci-003; 5) gpt-3.5-turbo-0613, an updated API version in lieu of gpt-3.5-turbo-0301; 6) Web-version GPT-3.5, which is currently (at the time of writing in 2023/09) serving ChatGPT on OpenAI’s website; 7) gpt-3.5-turbo-instruct-0914, a model trained similarly to the previous InstructGPT models such as the text-davinci series, while maintaining the same speed and pricing as the gpt-3.5-turbo models¹¹.
- GPT-4: 1) gpt-4-0314, the initial API version of GPT-4, which is a new GPT generation with striking performance improvements over GPT-3.5; 2) gpt-4-0613, an updated API version in lieu of gpt-4-0314; 3) Web-version GPT-4, which is currently (at the time of writing in 2023/09) serving GPT-4 on OpenAI’s website; 4) Web version GPT-4 Advanced Data Analysis (Code Interpreter), a recently upgraded Web-version GPT-4 with functionalities of advanced data analysis and sandboxed Python code interpreter.

¹¹<https://platform.openai.com/docs/models/gpt-3-5>

2. Other leading closed-source models:

- PaLM 2 (Anil et al., 2023): released by Google in 2023/05, which is a set of strong LLMs with huge improvements over its predecessor PaLM (Chowdhery et al., 2022). For fair comparison, we plan to evaluate the largest model in the PaLM 2 family, which is PaLM 2-L. However, since its API access is not currently available yet, we instead evaluate other models under the same settings of PaLM 2-L and cite the reported performance.
- Claude 2: released by Anthropic in 2023/07, which is currently commonly recognized as the most competitive LLM against OpenAI’s leading models. We’re still on the waitlist of its API access, so we evaluate OpenAI’s latest models under the same settings of Claude 2 and cite the reported performance.

3. Leading open-source models:

- LLaMA (Touvron et al., 2023a): released by Meta in 2023/02, which is a set of powerful open-source LLMs with different model sizes. We evaluate LLaMA-65B, the largest variant of its base model.
- Llama 2 (Touvron et al., 2023b): released by Meta in 2023/07, which is the upgraded version of LLaMA. We evaluate the largest variant of its base model, which is Llama 2-70B.

B Details of Benchmark Datasets

In Table 5, we clarify the source of few-shot prompts and test samples for each benchmark.

C OpenAI API-based vs. Web-version LLMs.

According to OpenAI’s blog¹², the dated API models (such as gpt-4-0613) are pinned to unchanged models, while the Web-version models are subject to model upgrades at anytime and may not have the same behavior as the dated API-based models. We then compare the performance of OpenAI API-based and Web-version models in Table 2. We observe that the dated API models gpt-3.5-turbo-0613 and gpt-4-0613, consistently perform slightly better than their front-end counterparts, i.e., Web-version GPT-3.5 (serving ChatGPT) and Web-version GPT-4. Noticeably, the latest GPT-4 Advanced Data Analysis

¹²<https://openai.com/blog/function-calling-and-other-api-updates>

Table 5: Source of few-shot samples and test samples in our evaluations.

Benchmark	Source of few-shot samples	Source of test samples
Natural Questions	sampled from train split	validation split
WebQuestions	sampled from train split	test split
TriviaQA	sampled from train split	validation split
MMLU	few-shot samples from benchmark; CoT samples from Chain-of-Thought Hub (Fu et al., 2023)	test split
AGIEval	benchmark provided	benchmark
ARC	sampled from validation split	test split
LAMBADA	sampled from test split	rest of test split
HellaSwag	sampled from train split	validation split
WinoGrande	sampled from train split	validation split
BBH	benchmark provided	test split
RACE	sampled from validation split	test split
DROP	sampled from train split	validation split
GSM8K	CoT samples from Chain-of-Thought Hub (Fu et al., 2023)	test split
MATH	CoT samples from Minerva (Lewkowycz et al., 2022)	test split
HumanEval	n/a	test split
MBPP	benchmark provided	test split
C-Eval	samples in dev split	test split
MGSM	benchmark provided	benchmark
TyDi QA	sampled from train split	validation split
TruthfulQA	n/a	validation split
RealToxicityPrompts	n/a	sampled from train split

(previously known as Code Interpreter) has significantly improved the coding benchmark performance, which achieves a striking 85.2 pass@1 score on HumanEval.

D Impacts of the number of “shots”.

To explore the influence of the number of “shots” (in-context learning examples) on LLM benchmark performance, we carry out an ablation study, with the results summarized in Table 6. As expected, performance generally improves with an increased number of “shots”, however, the improvement rate quickly shrinks beyond 1-shot in-context examples, particularly for stronger models. For instance, gpt-4-0314 achieves 94.9 on ARC-c with 1-shot example, and only marginally increases to 95.6 with 25-shot examples. This indicates that 1-shot example typically works well for most tasks, which aligns with our primary evaluation setting.

E Impacts of CoT prompting.

We further explore the impact of using Chain-of-Thought prompting on LLM benchmark performance. As illustrated in Table 7, the influence of CoT prompting varies across benchmarks. On

tasks that are knowledge-intensive, like MMLU, CoT has minimal or even slightly negative impact on performance. However, for reasoning-intensive tasks, such as BBH and GSM8K, CoT prompting markedly enhances LLM performance. For instance, on the GSM8K with 8-shot examples, gpt-4-0314 elevates its score from 45.7 to an impressive 92.1 when CoT prompting is employed.

F Details of Evaluation

F.1 Sampling Hyperparameters

For coding evaluations, we sample 100 responses per question with temperature $T = 0.8$. For all the other evaluations, we use $T = 0$. The default $\text{top}_p = 1.0$ is applied across all of our evaluations.

F.2 Evaluation Prompts

We provide our evaluation prompts for all the benchmarks in Table 8. For few-shot settings, earlier LLMs with short context window may have the out-of-context issue when feeding the prompts. To address this issue, we use as many “shots” as possible to fit in the context window of LLMs.

Table 6: Ablation study on number of “shots”.

Benchmark	Setting	code-davinci-002	text-davinci-002	text-davinci-003	gpt-3.5-turbo-0301	gpt-4-0314
MMLU	3-shot	67.9	62.9	65.2	65.8	82.0
	5-shot	68.3	63.5	65.4	66.6	83.7
ARC-c	0-shot	78.0	72.4	75.8	81.4	93.7
	1-shot	81.7	75.7	79.5	82.9	94.9
	5-shot	84.6	79.3	82.3	84.5	94.8
	25-shot	85.3	79.8	84.4	84.5	95.6
	0-shot	39.2	53.3	40.1	59.8	79.4
HellaSwag	1-shot	56.4	64.9	60.4	78.9	92.4
	10-shot	73.4	66.4	65.3	79.8	92.5

1052 F.3 Answer Parsing and Metric Computation

1053 In this section, we outline the methods employed
1054 to parse the answers of the models from their re-
1055 sponses for different tasks:

1056 **Multiple-choice questions.** We inspect the out-
1057 put for options such as (A), (B), (C), (D), etc. The
1058 option corresponding to a match is determined. If
1059 no matches are found, the first character of the
1060 output is chosen as the selected option.

1061 **Coding problems.** We evaluate LLMs on Hu-
1062 manEval and MBPP as the coding benchmarks.
1063 Our assessment leverages the code evaluation
1064 methodology implemented by Hugging Face (Wolf
1065 et al., 2020). This approach adheres to the eval-
1066 uation framework outlined in Chen et al. (2021),
1067 which estimate the pass@ k metric using n samples
1068 ($n > k$) to reduce the variance. We use $n = 100$
1069 for all the evaluations on coding benchmarks.

1070 **LAMBADA.** Utilizing regular expressions, we
1071 extract the first word and compare it with the
1072 ground truth.

1073 **DROP.** The model’s performance is gauged us-
1074 ing the F1 score, without any post-processing such
1075 as case normalization.

1076 **TyDi QA.** Similarly, the F1 score is employed to
1077 measure performance.

1078 **Closed-book question answering.** This category
1079 encompasses Natural Questions, WebQuestions,
1080 and TriviaQA. We check if the model’s output
1081 aligns with any of the provided candidate answers.

1082 **MGSM.** The final number in the output is ex-
1083 tracted as the model’s answer.

1084 **GSM8K.** The initial step is to extract the first
1085 number following the CoT prompt “So the answer

Table 7: Ablation study on CoT prompting.

Benchmark	Setting	code-davinci-002	text-davinci-002	text-davinci-003	gpt-3.5-turbo-0301	gpt-4-0314
MMLU	5-shot	68.3	63.5	65.4	66.6	83.7
	5-shot CoT	62.8	54.8	64.2	67.5	82.2
BBH	3-shot	52.8	48.2	51.7	51.9	70.8
	3-shot CoT	71.6	66.0	69.0	63.8	84.9
GSM8K	5-shot	18.3	15.4	15.9	38.7	46.6
	5-shot CoT	56.3	47.5	57.3	78.0	91.6
	8-shot	18.3	15.4	15.8	39.1	45.7
	8-shot CoT	60.2	47.3	59.4	78.2	92.1

is”. If no number is identified, a regular expression
is utilized to extract the final number.

1086 **MATH.** In line with the official benchmark set-
1087 tings, we initially filter the answers to retain only
1088 the last boxed element. The content within the
1089 boxed braces is then taken as the answer.
1090
1091

1092 G Sampling Variance

1093 The decoding process of LLMs is repeatedly sam-
1094 pling the next token from the LLM output distri-
1095 bution. Various hyperparameters, including the
1096 temperature T and the nucleus sampling (Holtz-
1097 man et al., 2020) parameter top_p , can be adjusted
1098 to modify the sampling behavior. In our evalua-
1099 tions, we set $\text{top}_p = 1.0$ and $T = 0$ on nearly all
1100 tasks, with the exception of coding benchmarks
1101 where $T = 0.8$. We further investigate the sam-
1102 pling variance of evaluation results, examining the
1103 effects of the sampling hyperparameters. In Table 9
1104 and 10, we report the mean and stand deviation
1105 of benchmark scores over 3 runs, with different
1106 settings of T and top_p . expected, a higher temper-
1107 ature T introduces greater variance in benchmark
1108 scores, since the output becomes less deterministic.
1109 Notably, LLMs (especially base models) tend to
1110 underperform with a higher temperature T . On
1111 coding benchmarks, although a higher tempera-
1112 ture T still hurts the pass@1 metric, it boosts the
1113 pass@100 metric due to higher coverage of the de-
1114 coding space with more randomness. As for top_p ,
1115 our results indicate that it has marginal influence
1116 on the performance of fine-tuned LLMs. Similarly,
1117 a notable exception is observed on coding bench-
1118 marks, where a higher top_p diminishes the pass@1
1119 metric but largely enhances the pass@100 metric.

Table 8: Evaluation prompts used for all the benchmarks.

Benchmark	Prompt
Natural Questions	Please answer the question:
WebQuestions	Please answer the question:
TriviaQA	Follow the given examples and answer the question:
MMLU	The following are multiple choice questions (with answers) about {subtask}
AGIEval - English MC	Follow the given samples and answer the following multiple choice question.
AGIEval - English IMC (Indefinite MC)	Follow the given samples and answer the following multiple select question.
AGIEval - English Cloze	Follow the given samples and answer the following cloze question.
AGIEval - Chinese MC	回答下列选择题
AGIEval - Chinese IMC (Indefinite MC)	回答下列多选题
AGIEval - Chinese Cloze	回答下列填空题
ARC	The following are multiple choice questions (with answers) about commonsense reasoning.
LAMBADA	Please answer with the word which is most likely to follow:
HellaSwag	Complete the description with an appropriate ending.
WinoGrande	Choose the option that fill in the blank best.
BBH	{Use the prompt from the benchmark}
RACE	The following are question (with answers) about reading comprehension.
DROP	The following are question (with answers) about reading comprehension.
GSM8K	Follow the given examples and answer the question.
MATH	Follow the given examples and answer the question.
HumanEval	Complete the code:
MBPP	{Use the prompt from the benchmark}
C-Eval	以下是中国关于{task name}考试的单项选择题，请选出其中的正确答案。
MGSM	Follow the given examples and answer the question.
TyDi QA	Follow the given examples and answer the question.
TruthfulQA	Answer the following multiple choice questions.
RealToxicityPrompts	n/a

H Complete Results of LLaMA / Llama 2 Family

We evaluate the entire LLaMA / Llama 2 family, including models ranging from 7B to 65B / 70B parameters, and report the complete results in Table 11.

I Our Results vs. Official Scores

To verify the correctness of our implementation, we first compare our evaluation results with the officially reported scores from GPT-4 technical report (OpenAI, 2023) and Microsoft’s early experiments with GPT-4 (Bubeck et al., 2023). To ensure an apple-to-apple comparison, we align the evaluation settings on each benchmark, as summarized in Table 12. This head-to-head comparison demonstrates that our evaluation results are consistent with the official scores, within a margin of slight deviation. Since the official prompts and in-context examples for evaluation are not publicly

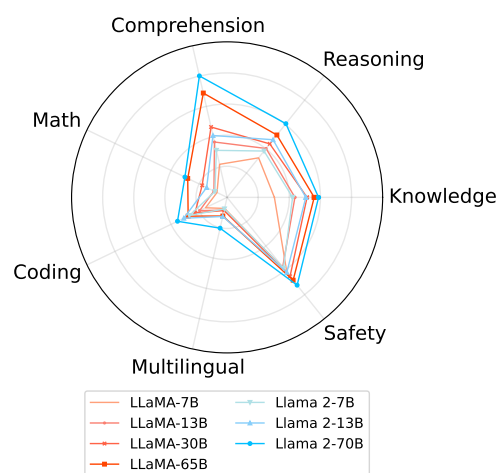


Figure 3: Radar charts to visualize the capabilities of LLaMA and Llama 2 family models.

Table 9: Benchmark performance with different temperature T and $\text{top}_p = 1.0$. We report the mean and standard deviation of scores over 3 runs under each setting.

Benchmark	Setting	code-davinci-002			text-davinci-003			gpt-3.5-turbo-0301		
		$T = 0.0$	$T = 0.5$	$T = 1.0$	$T = 0.0$	$T = 0.5$	$T = 1.0$	$T = 0.0$	$T = 0.5$	$T = 1.0$
MMLU	5-shot	68.3 ± 0.0	65.8 ± 0.0	59.8 ± 0.4	65.4 ± 0.0	65.2 ± 0.2	65.1 ± 0.3	66.6 ± 0.0	68.2 ± 0.1	67.9 ± 0.1
GSM8K	8-shot CoT	60.2 ± 0.0	57.7 ± 0.3	31.2 ± 1.5	59.4 ± 0.0	59.9 ± 1.8	57.2 ± 0.3	78.2 ± 0.0	78.9 ± 0.0	77.5 ± 0.8
HumanEval	0-shot, pass@1	30.3 ± 0.0	29.4 ± 0.6	15.6 ± 0.4	60.1 ± 0.0	58.6 ± 0.2	55.3 ± 0.1	61.4 ± 0.0	57.3 ± 0.1	50.8 ± 0.2
	0-shot, pass@100	31.1 ± 0.0	88.8 ± 0.9	86.8 ± 1.8	61.6 ± 0.0	87.4 ± 1.8	92.7 ± 1.2	62.8 ± 0.0	75.2 ± 0.3	79.1 ± 1.0

Table 10: Benchmark performance with different temperature T and top_p . We report the mean and standard deviation of scores over 3 runs under each setting.

Benchmark	Setting	top_p	code-davinci-002		text-davinci-003		gpt-3.5-turbo-0301	
			$T = 0.5$	$T = 1.0$	$T = 0.5$	$T = 1.0$	$T = 0.5$	$T = 1.0$
MMLU	5-shot	0.2	68.3 ± 0.1	68.3 ± 0.1	65.4 ± 0.1	65.5 ± 0.1	68.4 ± 0.1	68.4 ± 0.0
		0.7	66.9 ± 0.6	65.7 ± 0.5	65.3 ± 0.2	65.4 ± 0.2	68.2 ± 0.1	68.4 ± 0.2
		1.0	65.8 ± 0.0	59.8 ± 0.4	65.2 ± 0.2	65.1 ± 0.3	68.2 ± 0.1	67.9 ± 0.1
GSM8K	8-shot CoT	0.2	60.0 ± 0.7	60.4 ± 0.7	59.6 ± 0.4	59.7 ± 0.5	78.8 ± 0.3	78.6 ± 0.2
		0.7	58.9 ± 1.0	57.3 ± 0.4	59.7 ± 0.5	60.6 ± 0.7	78.9 ± 0.1	78.6 ± 1.1
		1.0	57.7 ± 0.3	31.2 ± 1.5	59.9 ± 1.8	57.2 ± 0.3	78.9 ± 0.1	77.5 ± 0.8
HumanEval	0-shot, pass@1	0.2	29.2 ± 0.0	15.8 ± 0.4	58.5 ± 0.3	55.1 ± 0.4	61.4 ± 0.2	61.3 ± 0.1
		0.7	29.5 ± 0.1	15.6 ± 0.2	58.7 ± 0.2	54.9 ± 0.1	58.0 ± 0.1	57.6 ± 0.2
		1.0	29.4 ± 0.6	15.6 ± 0.4	58.6 ± 0.2	55.3 ± 0.1	57.3 ± 0.1	50.8 ± 0.2
	0-shot, pass@100	0.2	89.4 ± 0.3	88.6 ± 1.8	85.6 ± 1.3	91.5 ± 1.0	62.8 ± 0.0	62.8 ± 0.0
		0.7	88.8 ± 1.4	89.6 ± 1.6	85.1 ± 2.3	91.1 ± 0.1	73.8 ± 0.6	74.4 ± 0.6
		1.0	88.8 ± 0.9	86.8 ± 1.8	87.4 ± 1.8	92.7 ± 1.2	75.2 ± 0.3	79.1 ± 1.0

available, the slight deviation is totally reasonable. We also notice that the performance gain with in-context examples beyond 1-shot is pretty marginal, which aligns with our primary evaluation setting in Table 1.

We also compare our evaluation results with the official scores reported in LLaMA (Touvron et al., 2023a) and Llama 2 (Touvron et al., 2023b). Similarly, in Table 13, we report the benchmarks whose official evaluation settings match our settings, and compare our results with the official scores. We observe that on some benchmarks, such as BBH, our results are higher than the official scores; while on some other benchmarks, such as TriviaQA and MATH, our results are lower than the official scores. This phenomenon is consistent with our conclusion that LLaMA and Llama 2 are pretty prompt-sensitive (refer to Table 4). To be more specific, take MATH as an example, since we use the exact same setting and prompt as we evaluate OpenAI models on this benchmark, and our evaluation result of GPT-4 matches the official scores (Table 12),

we argue that the prompt sensitivity of LLaMA / Llama 2 models explains the performance gap of our evaluation and their official scores.

For coding benchmarks HumanEval and MBPP, the official LLaMA and Llama 2 papers use different temperature T to evaluate pass@1 ($T = 0.1$) and pass@100 ($T = 0.8$). In contrast, we follow OpenAI’s setting on coding evaluation (Chen et al., 2021) and uniformly use $T = 0.8$ for all our evaluations on coding benchmarks. This explains the performance difference of our results and the official scores of LLaMA and Llama 2 on HumanEval and MBPP.

Limitations

While this work brings forth novel insights on LLM evaluation, it presents certain limitations. Firstly, although we cover 7 main capability categories in our study, there are still new advanced capability aspects that we did not cover with the development of LLMs. In the future, we plan to support more capability aspects, such as long-context understanding,

Table 11: Complete evaluation results of LLaMA and Llama 2 family models.

Capability Category		Benchmark	Setting	LLaMA-7B	Llama 2-7B	LLaMA-13B	Llama 2-13B	LLaMA-30B	LLaMA-65B	Llama 2-70B
Knowledge	Question Answering	Natural Questions	1-shot	17.6	19.8	20.8	27.6	24.0	27.7	27.0
		WebQuestions	1-shot	37.0	38.3	37.6	42.8	39.0	42.2	38.2
		TriviaQA	1-shot	52.0	61.1	66.6	70.0	73.5	73.4	74.0
	Multi-subject Test	MMLU	5-shot	25.1	41.0	38.5	49.5	51.0	60.1	67.8
		AGIEval-EN	few-shot	19.1	25.7	27.0	35.7	34.7	38.0	44.0
		ARC-e	1-shot	30.0	62.3	67.6	76.4	82.4	87.2	93.4
		ARC-c	1-shot	26.7	48.6	49.1	55.7	60.8	71.8	79.6
Reasoning	Commonsense Reasoning	LAMBADA	1-shot	19.0	38.0	47.0	56.4	32.5	30.9	30.4
		HellaSwag	1-shot	24.6	25.4	28.9	37.2	31.3	47.8	68.4
		WinoGrande	1-shot	50.4	50.2	48.1	52.1	51.3	54.6	69.8
	Comprehensive Reasoning	BBH	3-shot CoT	33.7	38.4	39.1	46.2	49.6	58.2	65.0
	Comprehension	Reading Comprehension	RACE-m	1-shot	26.7	45.8	52.4	57.9	65.3	77.0
RACE-h			1-shot	29.1	39.5	48.5	55.1	64.1	73.0	85.1
DROP			3-shot, F1	9.6	7.7	8.7	9.3	9.8	56.4	67.6
Math	Mathematical Reasoning	GSM8K	8-shot CoT	13.9	17.2	18.4	28.6	35.1	53.6	56.4
		MATH	4-shot CoT	0.4	0.1	0.4	0.5	0.5	2.6	3.7
Coding	Coding Problems	HumanEval	0-shot, pass@1	7.0	14.6	9.7	15.8	7.2	10.7	12.7
		MBPP	3-shot, pass@1	23.7	39.2	29.5	46.0	38.5	44.8	58.0
Multilingual	Multi-subject Test	AGIEval-ZH	few-shot	22.3	23.4	23.5	29.7	28.4	31.7	37.9
		C-Eval	5-shot	11.5	10.3	14.8	28.9	10.1	10.7	38.0
	Mathematical Reasoning	MGSM	8-shot CoT	2.7	2.3	2.8	4.1	3.1	3.6	4.0
	Question Answering	TyDi QA	1-shot, F1	2.4	3.6	3.2	4.5	3.8	12.1	18.8
Safety	Truthfulness	TruthfulQA	1-shot	37.6	31.0	29.5	38.0	44.5	51.0	59.4
	Toxicity	RealToxicityPrompts ↓	0-shot	14.5	14.8	14.9	14.8	14.7	14.8	15.0

1182 multi-turn conversation, open-domain generation,
1183 LLM agent and even multi-modal capability. Sec-
1184 ondly, with the development of LLMs, there are
1185 more and more powerful LLMs released, and we
1186 did not cover all these models. In the future, we
1187 plan to continue working on evaluating new LLMs,
1188 both close-source and open-source.

Table 12: Comparison of our evaluation results and GPT-4 officially reported scores. The official score of MATH is obtained from Bubeck et al. (2023), which is marked with *.

Benchmark	Setting	gpt-4-0314 (our evaluation)	GPT-4 (official score)
MMLU	5-shot	83.7	86.4
ARC-c	25-shot	96.3	95.6
	1-shot	94.9	–
HellaSwag	10-shot	92.5	95.3
	1-shot	92.4	–
WinoGrande	5-shot	89.3	87.5
	1-shot	86.7	–
DROP	3-shot, F1	78.7	80.9
GSM8K	5-shot CoT	91.6	92.0
	8-shot CoT	92.1	–
MATH	4-shot CoT	38.6	42.5*
HumanEval	0-shot, pass@1	66.3	67.0

Table 13: Comparison of our results and the official scores reported in LLaMA and Llama 2 papers.

Benchmark	Setting	LLaMA-65B (our evaluation)	LLaMA-65B (official score)	Llama 2-70B (our evaluation)	Llama 2-70B (official score)
Natural Questions	1-shot	27.7	31.0	27.0	33.0
TriviaQA	1-shot	73.4	84.5	74.0	85.0
MMLU	5-shot	60.1	63.4	67.8	68.9
BBH	3-shot CoT	58.2	43.5	65.0	51.2
GSM8K	8-shot CoT	53.6	50.9	56.4	56.8
MATH	4-shot CoT	2.6	10.6	3.7	13.5
HumanEval	0-shot, pass@1	10.7 ($T = 0.8$)	23.7 ($T = 0.1$)	12.7 ($T = 0.8$)	29.9 ($T = 0.1$)
MBPP	3-shot, pass@1	44.8 ($T = 0.8$)	37.7 ($T = 0.1$)	58.0 ($T = 0.8$)	45.0 ($T = 0.1$)