# InfiBench: Evaluating the Question-Answering Capabilities of Code Large Language Models

# Linyi Li Simon Fraser University linyi\_li@sfu.ca

# Shijie Geng

ByteDance Inc & Rutgers University sg1309@rutgers.edu

Zhenwen Li\* Yibo He\* Hao Yu\* Ziyue Hua\* Peking University **Guanghan Ning**ByteDance Inc

**Siwei Wang**ByteDance Inc

Tao Xie
Key Lab of HCST (PKU), MOE
taoxie@pku.edu.cn

# Hongxia Yang

The Hong Kong Polytechnic University (PolyU) hongxia.yang@polyu.edu.hk

#### **Abstract**

Large Language Models for code (code LLMs) have witnessed tremendous progress in recent years. With the rapid development of code LLMs, many popular evaluation benchmarks, such as HumanEval, DS-1000, and MBPP, have emerged to measure the performance of code LLMs with a particular focus on code generation tasks. However, they are insufficient to cover the full range of expected capabilities of code LLMs, which span beyond code generation to answering diverse coding-related questions. To fill this gap, we propose **InfiBench**, the first large-scale freeform question-answering (QA) benchmark for code to our knowledge, comprising 234 carefully selected high-quality Stack Overflow questions that span across 15 programming languages. InfiBench uses four types of model-free automatic metrics to evaluate response correctness where domain experts carefully concretize the criterion for each question. We conduct a systematic evaluation for over 100 latest code LLMs on InfiBench, leading to a series of novel and insightful findings. Our detailed analyses showcase potential directions for further advancement of code LLMs. InfiBench is fully open source at https://infi-coder.github.io/infibench and continuously expanding to foster more scientific and systematic practices for code LLM evaluation.

#### 1 Introduction

In recent years, Large Language Models (LLMs) have been revolutionizing the software development landscape [17, 12], demonstrating exceedingly strong and comprehensive capabilities in comprehending, generating, debugging, and summarizing code [9, 24]. For example, code LLM-powered products like GitHub Copilot [14] reached millions of active users within just one year of their launch.

Alongside the huge success of proprietary LLMs such as GPT-3.5 / GPT-4 [36] and Gemini [13], the development of open-source code LLMs<sup>2</sup> [35, 43, 39, 30] has been advancing at an unprecedented fast pace. As of June 2024, the Hugging Face Open LLM Leaderboard [4] has cataloged over 3,300 submissions of such models.

38th Conference on Neural Information Processing Systems (NeurIPS 2024) Track on Datasets and Benchmarks.

<sup>\*</sup>Equal contribution.

<sup>&</sup>lt;sup>2</sup>We define code LLMs as LLMs that show decent capabilities in the code domain, no matter whether they are exclusively trained or finetuned with code data or not.

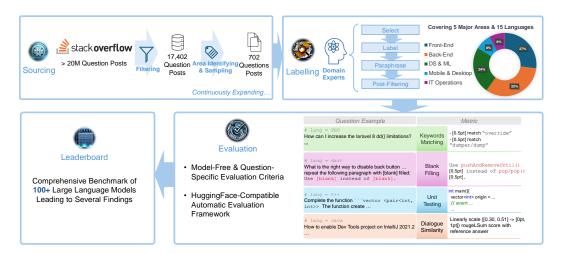


Figure 1: **InfiBench overview**. We construct the InfiBench benchmark by filtering high-quality and diverse question posts from Stack Overflow and annotating question-level evaluation criteria with domain experts. With an model-free automatic evaluation framework, we evaluate over 100 latest code LLMs (one of the most extensive evaluations for code LLMs to the best of our knowledge), leading to several insightful findings.

Table 1: Comparison between InfiBench and common existing benchmarks. Existing benchmarks weigh heavily on code generation, unit-test-based evaluation, and major programming languages. InfiBench processes a much higher diversity to reflect real-world code LLMs' usage scenarios. More discussion in Section 2.6.

Benchmark	Domain	# Question	Evaluation	Data Source	Highest LLM Score
HumanEval [9]	Python Programming	164	Test Cases	Hand-Written	90.2%
MBPP [3]	Python Programming	974	Test Cases	Hand-Written	81.1%
APPS [16]	Python Programming	10,000	Test Cases	Competitions	/ (no report yet)
DS-1000 [21]	Python Programming	1,000	Test Cases + Surface Form Constraints	StackOverflow	/ (no report yet)
HumanEval+ [26]	Python Programming	164	Augmented Test Cases	HumanEval	86.6%
HumanEvalPack [33]	Repair, Explain, Generation in 6 Languages	2,952	Test Cases	HumanEval	47.8%/52.1%/78.3%
LBPP [32]	Python Programming	161	Test Cases	Hand-Written	64%
SWE-bench [19]	Python Debugging / Repair	2,294	Test Cases	GitHub	22.06%
SWE-bench Verified [38]	Python Debugging / Repair	500	Test Cases	SWE-bench	45.20%
InfiBench	Free-Form Code Question Answering in 15 Languages	234	Keyword + Blank Filling + Test Cases + Text Similarity	Stack Overflow	70.64%

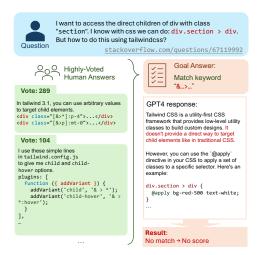


Figure 2: A challenging question paraphrased from Stack Overflow where GPT-4 fails to answer.

Given the plethora of code LLMs available, the development of reliable code benchmarks seems to lag in four aspects: (1) Benchmarks beyond code generation are relatively few. Benchmarks for code LLMs typically focus on a specific task or domain, often overly focus on code generation. For example, the widely-used HumanEval [9] and MBPP [3] purely focus on Python code generation, and DS-1000 [21] focuses on Python code generation in the field of data science. (2) Independent code benchmarks are relatively few. Recent efforts evolve existing benchmarks (e.g., HumanEval) to include more scenarios [33], languages [48], and tests [26]. However, these efforts lead to a series of benchmarks sharing the same source data (e.g., HumanEval Python problems), reducing score independence. (3) Existing code benchmarks are saturating. Strong LLMs are saturating existing benchmarks, e.g., GPT-4 has already achieved 90.2% Pass@1 score on HumanEval [37], while in real-world scenarios, GPT-4

can still fail as exemplified in Figure 2. (4) **Common benchmarks may be contaminated.** Some LLMs have unconventional high performance in common benchmarks and are suspected to have memorized benchmark-related data [10, 46, 32], obscuring the evaluation results. *Can we systematically and comprehensively evaluate code LLMs' abilities in challenging real-world usage scenarios?* 

To answer the question, we introduce InfiBench, a systematic benchmark for evaluating the free-form question-answering capabilities of code LLMs. As the first benchmark of its kind, the core principle of InfiBench aims to accurately represent how developers interact with and utilize such models in real-world scenarios. To achieve this, InfiBench comprises 234 questions that are carefully selected and proportionally filtered from the natural high-quality question distribution of Stack Overflow, without any constraints on topics, programming languages, question types, or answer forms. As a result, the curated 234 questions span 15 programming languages and 5 major areas: *front-end*, *back-end*, *DS&ML* (*data science and machine learning*), *mobile and desktop*, and *ITOps* (*information technology operations*).

Question diversity comes with evaluation challenges for two reasons. (1) Lack of metric. Unlike code generation or multiple-choice benchmarks, which can be evaluated through standardized methods like unit testing, there is no universal metric for response correctness for free-form questions. (2) Challenges with model-based evaluation. Model-based evaluations such as those involving GPT-4 are not only costly but also raise concerns about privacy and bias.

To mitigate the evaluation challenges, InfiBench includes an automatic evaluation framework that integrates four types of *model-free* metrics: keyword matching, blank filling, unit testing, and dialogue similarity. For each question, we invite industry domain experts to paraphrase the prompt, select the most appropriate metric, and write down the concrete criteria using domain-specific knowledge, with highly-voted answers from Stack Overflow as a reference. These questions and evaluation criteria are then cross-validated to ensure correctness and objectiveness and further calibrated to improve consistency across languages. Human experiments show that InfiBench evaluation aligns with humans better than LLM-based evaluation, achieving 85.1% agreement rate compared to 77.8% achieved by GPT-40-based evaluation.

As a novel and systematic benchmark disjoint with existing ones in terms of both forms and data sources, we believe that InfiBench is an ideal tool to measure existing code LLMs objectively. Hence, we conduct a systematic evaluation for **over 100 code LLMs** spanning both proprietary and opensource worlds using the InfiBench framework — the latest and most extensive evaluation for code LLMs to the best of our knowledge. Our evaluation leads to several insightful findings: (1) On InfiBench, GPT-4 achieves a score of 70.64%, being far from perfect but still far exceeding the most capable open-source models as of June 2024. On the other hand, GPT3.5 is surpassed by a few open-source models. (2) At similar model sizes, coding LLMs are usually visibly stronger than general LLMs; finetuning LLMs are usually visibly stronger than base LLMs. (3) The performance differences between different model families are huge, where one model could surpass another with less than 1/10 parameters, highlighting the importance of training data quality and techniques. (4) The scaling law is empirically verified for open-source models with fewer than 40B parameters, but not for those with more, where a turning point emerges. InfiBench is fully open source under CC BY-SA 4.0 license and continuously expanding<sup>3</sup>, including both the benchmark and Hugging-Face-compatible evaluation tools. All resources are available at https://infi-coder.github.io/infibench.

# 2 Benchmark Creation

InfiBench is created from a high-quality subset of Stack Overflow questions up until June 14, 2023. In this section, we describe the data curation process and the evaluation framework in detail.

#### 2.1 Data Curation

Stack Overflow is a question-and-answer website for developers with more than 24 million registered users as of June 2024 [41]. Since the website is a large collection of natural and diverse coding questions from real-world developers, we believe that questions from Stack Overflow can effectively evaluate code LLM's capabilities in real-world usage scenarios.

The full Stack Overflow dataset contains 23.54 million question posts and 34.68 million answer posts. Each question post has a total view count. Each answer post is attached to a question and has a vote count. The question creator can choose one answer as officially accepted.

<sup>&</sup>lt;sup>3</sup>In other words, **infi**nitely expanding, after which the **bench**mark is named.

Table 2: **InfiBench data statistics by area and language**. We uniformly sample a subset from the initial seed set (see Section 2.1) according to the area quota (see Section 2.2) for domain experts to select questions and annotate the correctness criterion to construct the benchmark.

		Initial Sec	ed Set	Tentative	Fir	nal InfiBench B	enchmark	
Area	Language	# Questions	% Area	# Questions	# Questions	% Questions	# Area	% Area
		" Questions	Quota	Quota	Quota	Quota	Quota	Quota
	Javascript	4912		44	44	18.80%		
Front-End	CSS	87	40.41%	10	10	4.27%	63	26.92%
	HTML	600		10	9	3.85%		
	Java	930		18	17	7.26%		
	C#	629		12	12	5.13%		
	PHP	462		10	9	3.85%		
Back-End	Go	117	18.71%	10	9	3.85%	77	32.91%
	Ruby	71		10	10	4.27%		
	Rust	96		10	10	4.27%		
	C/C++	287		10	10	4.27%		
DS & ML	Python	2779	21.39%	47	47	20.09%	56	23.93%
DS & ML	R	184	21.3970	10	9	3.85%	30	23.93%
	Dart	1562		19	19	8.12%	19	8.12%
Mobile &	Kotlin	383	18.13%	10				
Desktop	Swift	551	10.13 /0	10	Removed du	ring Post-Filteri	ng (see Se	ection 2.3)
	VBA	16		9				
IT Ops.	Bash	188	1.36%	21	19	8.12%	19	8.12%
To	tal	13854	100.0%	270	234	100.00 %	234	100.00%

As we aim to create a benchmark where the correctness evaluation criteria are clear, we view the positively voted answers as an important reference source. Hence, we choose to keep only the questions that have at least three positively voted answers and an officially accepted answer, which turn out to be 1,090,238 questions. For these one million questions, we choose to keep questions that are frequently viewed and relatively new. To fulfill this criterion, we draw a scatter plot of these  $\approx 1$  million questions, plotting the number of days since their creation until June 14, 2023 (data collection end-date) on the x-axis against the logarithm of their view counts on the y-axis. As shown in Figure 3, we empirically determine to keep questions that lie above the line connecting (0,5) and (3000, 15.5), resulting in a subset of 17,402 questions.

Utilizing the mandatory question tags of these questions, we then manually construct a tag tree that covers the 200 most frequent tags, enabling us to identify the top programming languages and areas for 14,330 out of these 17,402 questions. These questions are from 24 programming languages, with each language being categorized into one primary area among the five (front-end, back-end, DS&ML, mobile and desktop, and ITOps). Lastly, we exclude 6 programming languages that either describe data or are domain-specific: JSON, regex, Markdown, YAML, CSV, and SQL. As a result, we compile 13,854 questions that serve as the *initial seed set*.

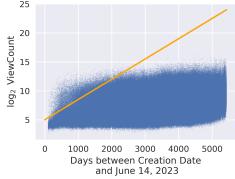


Figure 3: Scatter plot of filtered Stack Overflow questions. Questions above the orange line kept.

#### 2.2 Sampling

Based on a user study of developers' demand from our organization, we allocate the tentative area quota to be 25%, 25%, 25%, 15%, and 10% for front-end, back-end, DS&ML, mobile and desktop, and IT Ops, respectively. Inspired by HumanEval size and considering the labelling labor cost, we set 200 questions as the target benchmark size. Hence, the tentative size quotas by area are 50, 50, 50, 30, and 20 respectively. We then proportionally distribute the area quotas to language quotas based on the frequency of each language in the initial seed set. However, we observe that following this rule, certain languages such as CSS and C/C++ end up with fewer than 10 questions, which may yield unreliable language-level sub-score, so, for these languages, we set their quotas to 10.

As a result, we derive the *tentative* question quota for each language as shown in Table 2, which sums up to 270 questions. After determining the tentative question quota, we uniformly sample from the initial seed set a roughly two times larger pool for the domain expects to select and annotate.

#### 2.3 Human Annotation

We recruited five domain experts inside our company to create the benchmark, each in charge of one area. The annotation process is composed of three steps:

- Step 1: Question Selection and Type Annotation. Domain experts select high-quality questions from the inspecting set and annotate the question type to be one of the four: code completion, code debugging, config and environment debugging, and knowledge question-answering.
- Step 2: Prompt Paraphrasing. Domain experts paraphrase and simplify the original question body into succinct and explicit instructions. We include this step for two main purposes: (1) Reduce domain gap. From user-shared conversations collected from ShareGPT, we observe that when interacting with code LLMs, users tend to provide short and direct instructions like "Fix problem..." and "Debug code...". However, when posting Stack Overflow questions, users tend to be lengthy with courtesy words. We ask the domain experts to paraphrase the question to code LLM user's style without changing the semantics. (2) Reduce the impact of memorization and data contamination. Some code LLMs may be trained or finetuned with Stack Overflow data. Paraphrasing the questions can help to mitigate the result advantages of these models. Benchmark results in Table 4 reveal the effectiveness of this step where copying Stack Overflow answers only achieves a 65.18% score. We defer further discussion in Section 2.5.
- Step 3: Correctness Criterion Annotation. Domain experts choose one or multiple evaluation metrics from our supported ones (see Section 2.4) and annotate the concrete criterion following a YAML schema. External files can be attached if needed, e.g., unit tests and reference answers.

Calibration and Post-Filtering. To improve annotation consistency and objectiveness, we introduce a few checkpoints for domain experts to read others' annotated cases, discuss them, and reach consensus for controversial cases. After the 270 tentative questions were annotated, we then ran an initial evaluation of all these questions on over 30 code LLMs. This initial evaluation helps us to identify questions whose criteria are incorrect or out of distribution. We filter out these questions and then remove all questions from Kotlin, Swift, and VBA languages since the questions in these languages are too few after filtering. After this calibration and post-filtering process, the final benchmark includes 234 questions spanning over 15 languages. Their statistics are shown in Table 2. As we can observe, compared to the population area distribution of high-quality Stack Overview questions (see "% Area Quota" column under "Initial Seed Set"), the area distribution of final benchmark questions (see "% Area Quota" column under "Final InfiBench Benchmark") is more balanced and less biased towards front-end, mobile, and desktop topics.

#### 2.4 Evaluation Criteria and Evaluation Framework

In response to the diversified questions, InfiBench evaluation framework integrates four types of model-free and automatic metrics as below. Domain experts choose one or multiple metric types along with their weights and concretize.

- **Keywords Matching.** Though the responses can be in diverse forms, for a significant portion of benchmark questions, we find that the existence of some keywords strongly determines the quality of the response. Domain experts can write rules that match keywords and regular expressions or construct recursive logical expressions on top of keyword-matching results. When multiple keywords exist, each matching result can have its weight in the final score.
- Blank Filling. For some questions, it is challenging to measure the correctness given the response uncertainty. In this case, domain experts can instruct the model to answer the question by following a given template and filling in the blanks in the template. The blanks can correspond to either natural language or code snippet. Then, similar to keywords matching, each blank can match potential keywords, regular expressions, or recursive logic expressions built upon matching results. This metric type tests not only the model's QA ability but also its instruction-following ability.
- **Unit Testing.** For code-intensive questions, we can follow existing benchmarks to evaluate response correctness by unit tests. For this type, domain experts may add more specifications in the prompt to allow for unit-test-based evaluation, such as specifications on function name, input arguments, and output format. Domain experts can further import the context setup and cleanup script.
- **Dialogue Similarity.** For natural-language-intensive questions, domain experts can extract and shorten the reference answers from Stack Overflow, and then use the ROUGE score [25] to evaluate the response similarity with reference answers. The ROUGE score was initially proposed and

widely used in evaluating the quality of text summarization and machine translation. To map the ROUGE score back to our benchmark scale, we allow domain experts to tune the mapping interval and scores within the interval are then linearly mapped to our score scale.

The example questions and corresponding criteria are illustrated in Figure 1. Detail statistics of metric type ratios, question type ratios, and prompt length are shown in Table 3.

Score Computation. We treat each question equally with one point each. Given 234 questions in the benchmark, the full score is 234, and we by default report the percentage score (achieved score divided by 234) unless otherwise noted. The one point for each question can be further decomposed into a few scoring points within each question. For example, a question may contain four keywords with weights 2, 1, 1, and 1 each. Then, matching each keyword can contribute to 0.4, 0.2, 0.2, and 0.2 points respectively to the final score.

Table 3: InfiBench statistics.

(a) Question type. (b) Metric type.

Question Type	Ratio
Code Completion	30.37%
Knowledge Question-Answering	30.37% 27.04% 26.67%
Code Debugging	26.67%
Config & Environment Debugging	15.93%

Metric Type	Ratio
Keywords Matching	57.41%
Blank Filling	12.22%
Unit Testing	19.26%
Dialogue Similarity	11.85%

# (c) Prompt token length with Code Llama tokenizer.

min	25% quantile	median	mean	75% quantile	max
43	145.75	223	338.46	359.50	5047

**Implementation.** We have implemented an automated evaluation framework with Python, publicly available at https://infi-coder.github.io/infibench. Specifically, for blank-filling evaluation, we use the longest common subsequence matching via dynamic programming to capture the filled blanks in the response. For unit-testing evaluation, we construct a runtime environment that supports the test execution for nine languages. We plan to integrate the framework into the Hugging Face Open LLM Leaderboard [4] to further ease the evaluation burden.

How does InfiBench Evaluation Align with Human? To evaluate the alignment between InfiBench evaluation and human expert evaluation, we randomly sample 100 questions without replacement from the benchmark and select three strong LLMs to generate responses: GPT-4-0613, GPT-3.5turbo, and Mistral Codestral-22b. For each question, we randomly choose two out of these three model responses to construct response pairs, resulting in 100 response pairs  $\mathcal{R} = \{(A_i, B_i) : 1 \le i \le 100\}$ . For each response pair  $(A, B) \in \mathcal{R}$ , we use InfiBench, GPT-40, and human expert to evaluate into four outcomes: A is more correct than B(A > B); B is more correct then A(B > A); both A and B are correct  $(A \approx B \uparrow)$ ; both A and B are incorrect  $(A \approx B \downarrow)$ . Our purpose is to evaluate how InfiBench evaluation aligns with humans, specifically when compared to the widely-used LLM-as-a-judge (i.e., model-based evaluation) [47]. The concrete grading criteria is as below:

- InfiBench gives a score between [0,1] for each response in the pair. If the score difference in the pair is larger than 0.2, we label the outcome to be A > B or B > A respectively; otherwise, if the maximum score among the two is larger than 0.5, we label the outcome to be  $A \approx B \uparrow$ ; otherwise, we label the outcome to be  $A \approx B \downarrow$ .
- For GPT-40 evaluation, we deploy the prompting template from LLM-Blender [18, Appendix E] and trigger GPT-40 for grading the four outcomes. We enhance the reliability of the comparison by switching A and B and prompting GPT-40 twice. We record the preference only when a consistent preference exists.
- For human evaluation, we recruit human annotators who came up with the criteria to label the comparison preference since they are familiar with the questions and have strong expertise. Annotators have no access to the evaluation results of InfiBench and GPT-40, nor which source model generates the response. Annotators were instructed to directly label each pair with the four outcomes.

We defer the consensus matrices between InfiBench/GPT-40 and human annotators along with more findings in Appendix C. If we only count the cases where both InfiBench/GPT-40 and humans have clear preferences, the agreement rate between InfiBench and humans is 85.1%, and the agreement rate between GPT-40 and humans is 77.8%. Hence, the InfiBench evaluation aligns with human experts better than the GPT-40 evaluation (with >80% confidence). We observe that the advantage of InfiBench comes from the ability to detect deceptive answers. some model responses pretend to be helpful with lengthy wording and hallucinations. GPT-40 is more likely to be cheated than InfiBench, which looks for key concepts that should exist in a helpful answer.

#### 2.5 Mitigations on Memorization and Data Contamination

InfiBench is created from the publicly available Stack Overflow corpus to reflect real-world scenarios, and this corpus may already exist in the training set of some code LLMs (e.g., DeepSeek Coder [15] and StarCoder 2 [28]). Hence, some code LLMs may achieve a high score simply due to memorization. To mitigate this, we asked the domain experts to paraphrase every question as an essential step (see Section 2.3). Hence, copying either the highly voted answers or officially accepted answers of the original questions only achieves 65.18%, being far from perfect and inferior to GPT-4's 70.64%. Furthermore, code LLMs that use Stack Overflow data do not demonstrate significant advantages over those without. Hence, we deem the effect of contamination as small.

On the other hand, we release the post IDs of the source question posts of InfiBench. Hence, future LLM training could consider this benchmark to conduct deduplication and ablation studies on data contamination. Another usage of our benchmark is to evaluate retrieval-augmented (RAG) code LLMs where perfect retrieval from Stack Overflow and moderate adaptation should solve these questions, which we leave as future work.

# 2.6 Comparison with Existing Benchmarks

In Table 1, we compare InfiBench with several existing benchmarks for code LLMs. As reflected in the table, InfiBench strongly complements existing benchmarks for code LLMs by (1) extending them beyond code generation to a wide range of real-world tasks, (2) diversifying them since InfiBench does not share the same source as existing ones, and (3) increasing the differentiation as an unsaturated benchmark. Related benchmarks are further illustrated in Section 5. On the other hand, the benchmark is limited in size due to the high cost of correctness criteria labelling, and we are continuously expanding the benchmark.

# 3 Evaluation and Leaderboard

We systematically evaluated over 100 code LLMs spanning both proprietary and open-source worlds on InfiBench. To the best of our knowledge, this is the most extensive evaluation for code LLMs.

**Evaluation Protocol.** We adopt best@10 as the main evaluation metric: 10 responses are sampled and evaluated for each question, then the best score per question is recorded and summed up. Throughout the evaluation, we set sampling temperature T=0.2 and top p=0.9.

Furthermore, we swept sampling parameters with GPT-4 and the detailed results are in Appendix G. In a nutshell, for maximizing the performance under best@10, the best parameters are T=1.0 and p=0.9, leading to a score of  $76.15\%\pm0.21\%$  (in comparison to  $70.64\%\pm0.82\%$  in our main setting T=0.2, p=0.9). In particular, the temperature T affects much and the effect of top p is minor. We decided to stick to the original parameters T=0.2 and p=0.9 in the main evaluation since this setting is more akin to the real-world scenario where user generates once with low temperature.

We design two system prompts (shown in Appendix H), one for normal questions and the other for open-ended questions with an additional sentence to encourage succinct responses. For generic models, we generate the prompt with "{system prompt}\n{content prompt}" format; for instruction-finetuned or chat models, we generate the prompt with their prompt templates.

For proprietary models, we evaluate the latest models from OpenAI (GPT-4, GPT-40, etc), Anthropic (Claude 3), and Mistral AI (Mistral Small/Medium/Large) with API calling. When budget permits, we repeat each evaluation three times and report standard deviation. For open-source models, we download models from Hugging Face and evaluate them on an 8xA100 server with bigcode-evaluation-harness [5]. When the model size is within 30B parameters, we repeat each evaluation three times and report the standard deviation. All raw model responses are available at https://figshare.com/articles/dataset/InfiBench\_Detail\_Evaluation\_Data/26104864. More details on the evaluation protocol are in Appendix E.

**Leaderboard.** In Table 4, we present aggregated InfiBench leaderboards by model family, model type, and model size. The full leaderboard is deferred to Appendix E due to space limit. The table includes scores from using the original Stack Overflow answer posts as reference. Results are also presented as a scatter plot in Figure 4, where normal models are shown as scatters with error bars,

Table 4: **Aggregated InfiBench leaderboards (best viewed zoomed in and in color).** "Size" column records number of parameters. For MoE models, "total params. / params. activated during inference" is recorded. Bar colors stand for General Base, General Finetuned, Code Base, and Code Finetuned models respectively. Icon "a" stands for proprietary models otherwise open-source. Full leaderboard in Appendix E.

(a) InfiBench leaderboard by model family, where best (b) InfiBench leaderboard by model type, where top five model within each model family is shown.

Family	Best Model Name	Size	InfiBench Score	Type	Rank	Model Family / Model Name	Size	InfiBench Score
⊕GPT-4	GPT-4-0613	?	$70.64\% \pm 0.82\%$	Pro-	1	aGPT-4/GPT-4-0613	?	$70.64\% \pm 0.82\%$
DeepSeek Coder	deepSeek-coder-V2-instruct	236B / 21B	65.49%	prie-	2	GPT-4/GPT-4-turbo-1106	?	$68.42\% \pm 0.38\%$
Claude 3	Claude 3 Opus	?	63.89%	tary		GPT-4/GPT-4o-2024-05-13		66.19%
Mistral Open	Codestral-22b	22B	$62.98\% \pm 0.56\%$	Model				63.89%
Phind	Phind-CodeLlama-34B-v2	34B	59.00%					58.22%
⊕Mistral	mistral-large	?	58.22%	Code	1			$65.49\%$ $62.98\% \pm 0.56\%$
DeepSeek LLM	deepseek-llm-67b-chat	67B	57.41%		2			62.98% ± 0.30%
GPT-3.5	GPT-3 5-turbo-0613	2	$56.47\% \pm 1.34\%$					59.00%
Owen	Owen-72B	72B	55.34%		-			58.47%
Magicoder	Magicoder-S-CL-7B	7B	$52.71\% \pm 0.72\%$			***************************************		47.36%
WizardLM	WizardCoder-Python-34B-V1.0	34B	52.59%	Code	2		34B	43.13%
Code Llama	CodeLlama-34b-Instruct	34B	50.45%	Base	3	StarCoder2/15B	15B	42.52% ± 1.24%
01.AI	Yi-34B-Chat	34B	49.58%	Model	4	Code Llama/CodeLlama-13b	13B	$41.66\% \pm 0.84\%$
Zephyr	Zephyr 7B beta	7B	$46.31\% \pm 1.11\%$		5	Code Llama/CodeLlama-13b-Python	13B	$41.31\% \pm 0.90\%$
StarCoder2	15B-Instruct	15B	$45.89\% \pm 0.95\%$	C1	1	DeepSeek LLM/deepseek-llm-67b-chat	67B	57.41%
DeepSeek MoE	deepseek-moe-16b-chat	16B / 2.8B	$45.18\% \pm 1.65\%$	Fine-	2			
OctoPack	OctoCoder	15.5B	$44.55\% \pm 0.79\%$	tuned	3			52.97%
gemma	gemma-7b-it	7B	$40.68\% \pm 1.23\%$	Model	-			49.58% 46.31% ± 1.11%
	Llama2-70B-Chat	70B	39.30%					
InternLM	InternLM-Chat-20B	20B	$37.41\% \pm 0.75\%$	C1	•			55.34% 43.69% ± 1.09%
Baichuan2	Baichuan2-13B-Chat	13B	$34.40\% \pm 1.34\%$	Base	-			39.87%
StarCoder	StarCode+	15.5B	30.67% ± 1.57%	Model	3			37.69%
CodeGen2.5	CodeGen2.5-7B-Instruct	7B	$29.57\% \pm 1.53\%$		5			$31.69\% \pm 0.29\%$
ChatGLM	ChatGLM3-6B	6B	$28.23\% \pm 0.58\%$			Quenquen 12	7.5	
⊕davinci	davinci-002	?	$21.25\% \pm 1.17\%$					
Phi	Phi1.5	1.5B	20.56% ± 0.09%	a) In	GD.	anah landarbaard bu	modal	siza where has
CodeGeeX	CodeGeeX2-6B	6B	$19.88\% \pm 0.36\%$	(c) III	шь	encir leaderboard by	modei	size, where bes
CodeGen2	CodeGen2-16B	16B	16.97% ± 1.15%	node	1 w/i	thin the threshold is sh	own	
IEITYuan	Yuan2-51B-hf	51B	15.25%		. ,,,	thin the threshold is si	0 11 11.	
CodeGen	CodeGen-16B-multi	16B	$13.62\% \pm 1.18\%$	Size Thr	eshold	Model Family / Model Name	Size In	nfiBench Score
						⊕GPT-4/GPT-4-0613	?	$70.64\% \pm 0.82\%$
				<100B		Mistral Open/Codestral-22b	22B	$62.98\% \pm 0.56\%$
Human				<20B		DeepSeek Coder/deepseek-coder-6.7b-instru	ct 6.7B	$53.25\% \pm 0.40\%$
	Officially-Accepted Answer Post		52.90%	<5B		DeepSeek Coder/deepseek-coder-1.3b-instru		$41.32\% \pm 1.12\%$
	aGPT-4 beepSeek Coder cClaude 3 Mistral Open Phind aMistral Open Phind aMistral Open Phind cOPT-3.5 Qwen Magicoder WizardLM Code Llama 01.A1 Zephyn StartCoder2 DeepSeek MoE OctoPack gemma Llama 2 Llama 2 InternLM Baichuan2 StartCoder CodeGen2.5 CodeGen2.5 CodeGen2.5 CodeGeeX CodeGe	aGPT-4 DeepSeck Coder Glade 3 Mistral Open Phind Mistral DeepSeck LLM GPT-3-5 GWen Mistral DeepSeck LLM GPT-3-5 GWen Magicoder WizardLM Code Llama-34b-nstruct GPT-3-5 GWen Magicoder WizardLM Code Llama-34b-nstruct GPT-3-5 Wagicoder WizardLM Code Llama-34b-nstruct GPT-3-5 Wagicoder WizardLM Code Llama-34b-nstruct Code Llama UnAl Yi-34B-Chat Zephyr Beta LSB-Instruct DeepSeck MoE Geepseck-moe-16b-chat CottoPack Gemma Llama 2 Llama-2 70B-Chat InternLM InternLM-Chat-20B Batichuna-1-37b-Instruct CodeGen-2.5 CodeGen-2.5 CodeGen-2.5-7B-Instruct CodeGen-2.5 CodeGen-2.5-7B-Instruct CodeGen-2.5-7B-Instruct CodeGen-2.5 CodeGen-2.5-7B-Instruct CodeGen-2.5 CodeGen-2.5-7B-Instruct CodeGen-2.5 CodeGen-2.5-7B-Instruct CodeGen-2.5-7B-	GPT-4   GPT-4-0613   2   7	### GPT-4   GPT-4-0613   70.64% ± 0.829   65.89%	AGPT-4	Code   Code	AGPT-4   GPT-4-0613	CodeCear   CodeCear

MoE models are shown as horizontal segments with error ranges connecting the activated parameters during inference and total parameters, and strong proprietary models are shown as horizontal lines.

In both tables and the figure, we classify LLMs by general/code and base/finetuned. The general LLMs are claimed to have strong capabilities beyond code, e.g., in various natural language tasks, while the code LLMs are exclusively optimized for the code domain. The base LLMs only went through the pretraining phase, while the finetuned LLMs are claimed to have instruction-following capabilities or are finetuned on instruction or human preference datasets.

# 4 Analysis and Discussion

The best model so far, GPT-4, is still far from perfect, and open-source models are competitive but still far from GPT-4. GPT-4 achieves the highest score 70.64% (interestingly, achieved by GPT-4-0613 instead of the more recent GPT-4o), then Claude 3 Opus with a score 63.89%, and then Codestral-22b [1] with a score 62.98% and deepseek-coder-33b-instruct [15] with a score 62.96%. The result implies that: (1) Noting that the full score is 100%, even the powerful GPT-4 is still far from perfect, which is in contrast to its ≈90% HumanEval score. We inspect the score breakdown. For the two most frequent metric types, keywords matching and unit testing, GPT-4 achieves similar scores 66.61% and 76.00% respectively. For blank filling, the score is relatively lower at 58.08%. These scores imply that GPT-4 may still lack generic ability in answering diversified real-world questions related to code. When instructed to follow a given template to answer (blank filling), due to the more strict requirement and narrower solution space, its lack of capability is more pronounced. (2) There is still a visible gap between open-source models and GPT-4. The gap between the most powerful open-source model, Codestral-22b, and GPT-4 is roughly 8 points. On the other hand, noticing that GPT-3.5-turbo achieves 56.47%, the open-source model, Codestral-22b, is now reliably better than GPT-3.5-turbo with merely 22B parameters which is promising.

**Among open-source models, different models have various performances.** Figure 4 systematically visualizes the performance of different open-source models at diverse scales. Although there is a general tendency that larger models achieve higher scores, the scores among different models at a similar scale differ largely. For example, on scale 7B, the best-performing model is at around

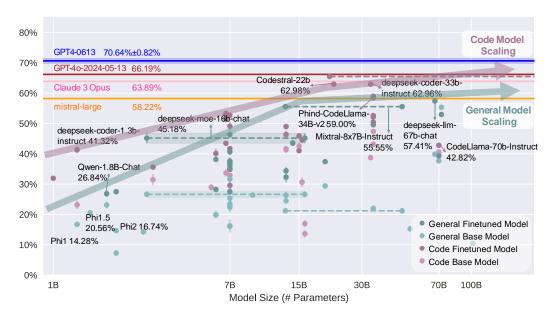


Figure 4: Scatter plot for all evaluated LLMs on InfiBench. x-axis is the model size in terms of number of parameters and y-axis is InfiBench score. Projected empirical scaling laws for both general and code models are drawn. Detail discussion in Section 4.

55%, pretty close to GPT-3.5, while the low-performing model stays at around 15%. Moreover, deepseek-coder-1.3b-instruct achieves 41.32% at 1.3B and surpasses a few models at scale 70B or 100B. Hence, though scaling matters, the training techniques and training data are equally important or even more, helping to reduce the required scale for achieving a certain score by more than  $10\times$ .

Hard problems generalize their difficulties. We rate the benchmark problem difficulty with five levels by how well GPT-4 and GPT-3.5-turbo answer them, as detailed in Appendix D. Example questions from each level are shown in Appendix I. We present the detail result table including the sub-score for each difficulty level in Appendix E. Interestingly, the trend is *highly consistent* that sub-scores decrease along with the increase of problem level. Specifically, hard problems for the most powerful model yet, GPT-4, are also generally hard for open-source models. These hard problems usually correspond to code generation with long and domain-specific context or challenging blank-filling questions since blank-filling is a specific task that rarely appears in training data before.

Instruction finetuning is important for QA. Among models of similar scales and the same family, we find that the best-performing ones almost always include an instruction-finetuning phase, such as deepseek-llm-67b-chat, deepseek-coder-33b-instruct, CodeLlama-34B-Instruct, and Qwen-18B-Chat. In contrast, the pretraining models, such as davinci-002 and phi models, usually perform poorly despite their strong performances in code generation benchmarks. Instruction-finetuning is also critical for other code domain tasks such as code generation. As shown in Appendix F.1 where we plot model scores in QA (measured by InfiBench) and code generation (measured by HumanEval pass@1), instruction-tuning generally improves both QA and code generation, but the improvement is usually more significantly on code generation but more moderately on QA. As a result, we suggest generalizing the instruction-finetuning data beyond simple coding problems to improve code LLMs. Indeed, our preliminary experiments show that, after fine-tuning with the decontaminated and sanitized Stack Overflow data, we improved InfiBench scores for Codellama-13b-Instruct from 46.37% to 60.74% and for mixtral-8x7B-Instruct from 55.55% to 62.61%.

Some models may focus too much on code generation, especially the small ones. As detailed in Appendix F.1, we observe that for large models (>30B) and top entries, InfiBench and HumanEval pass@1 scores coincide well. However, for smaller models, the score tendencies start to diverge, where some models are relatively stronger on InfiBench (Mixtral-8x7B-Instruct) and more are relatively stronger on HumanEval (Phi1, Phi2, gemma-7b, ...). This phenomenon implies that a few models may be optimized too heavily on code generation benchmarks while ignoring the performance in generic code scenarios as represented by InfiBench, which in turn highlights the significance of free-form QA benchmarks like InfiBench in detecting capability imbalance in code LLMs.

Code LLama models have unique characteristics. We evaluated all Code Llama models [39]. As shown in Table 5, we found finetuning on Python data improves on HumanEval but hurts InfiBench scores, while instruction finetuning usually improves InfiBench scores but may hurt HumanEval. As a side product, we found CodeLlama-70B may be overly safeguarded and denies answering some

Table 5: Evaluation on eight models from the Code Llama [39] family showcases intense Python finetuning may hurt free-form QA ability, despite achieving higher HumanEval scores.

	Benchmark	Base	Python	Instruct
7B	HumanEval	33.5%	38.4% (+4.9%)	34.8% (+1.3%)
/ <b>D</b>	InfiBench	$37.62\%_{\pm 1.28\%}$	$32.89\%_{\pm0.45\%}$ (-4.73%)	$35.15\%_{\pm 1.28\%}$ (-2.47%)
13B	HumanEval	36.0%	43.3% (+7.3%)	42.7% (+6.7%)
130	InfiBench	$41.66\%_{\pm0.84\%}$	$41.31\%_{\pm 0.90\%} (-0.35\%)$	$46.37\%_{\pm 1.26\%}$ (+4.71%)
34B	HumanEval	48.8%	53.7% (+4.9%)	41.5% (-7.3%)
34D	InfiBench	47.36%	43.13% (-4.23%)	50.45% (+3.09%)
70B	HumanEval	53.0%	57.3% (+4.3%)	67.8% (+14.8%)
/0B	InfiBench	40.60%	40.29% (-0.31%)	42.82% (+2.22%)

safe questions in InfiBench. More model-specific findings are presented in Appendix F.

Code models and general models may exhibit different scaling laws, and open-source models scale well only within 40B yet. In Figure 4, we use the top-performing code and general models at each scale respectively to regress and extrapolate model performance at larger scales. As shown, code models tend to have higher capabilities compared to general models of the same scale, though the gap shrinks for larger models. Hence, when the compute budget is heavily limited, training exclusively in the code domain could be more efficient for building strong code LLMs.

In Figure 4, both predicting curves are split into two segments, steep in the first segment and much flat in the second. Following the first segment, open-source models catch up with GPT-4 at around 50B scale. However, following the second segment, they may need to be at >300B scale to catch up. The finding contradicts the common scaling law [20, 34, 7] where a strong linear relationship between model scale and capability exists. The contradiction implies that very large open-source models (>40B) may fail to achieve the expected performance at their scales, or there is some non-trivial barrier when scaling the model beyond 40B, or the scaling law may change at such a large scale. We leave further investigation as the future work. Notably, after the release of InfiBench, Deepseek-coder-v2 [49] was released as the largest code LLM to our knowledge in an MoE architecture with 236B total and 21B active parameters. On InfiBench, Deepseek-coder-v2 achieves 65.49%, setting the new baseline for open-source LLMs but still being inferior to GPT-4. More importantly, the score is within the predicted range of our empirical scaling law.

We defer dataset card and data accessibility details, discussion on limitations and societal impact, full leaderboard, additional findings, ablation studies, and data examples in appendices.

# 5 Related Work

Large language models [44, 11, 8] are transforming people's lives. In the coding domain, LLMs [9, 24] are shown to be capable of completing a wide range of tasks such as code generation, debugging, and question-answering. Recently, code LLMs are booming. New models, including both proprietary [14, 36] and open-source ones [4, 35, 42, 43, 22, 30, 39, 49], emerge almost every month.

Benchmarks for code LLMs are developing, though at a relatively slower pace. Common benchmarks, e.g., APPS [16], MBPP [3], and HumanEval [9], focus on code generation and unit-test-based evaluation. Some efforts augment these benchmarks by language translation (e.g., Multilingual HumanEval [2], HumanEval-X [48]), test augmentation (e.g., HumanEval+ [26]), task generalization (e.g., HumanEvalPack [33]), and human rewriting (e.g., LBPP [32]). To systematically evaluate real-world problem solving, recently, SWE-bench [19], its filtered version SWE-bench Verified [38], and RepoBench [27] are proposed but they still primarily focus on code generation. Some general-purpose benchmarks, e.g., Arena-Hard [23], contain code-related questions, but rely on LLM to judge and do not provide domain-specific scores. CodeXGLUE [29] considers multiple coding capabilities beyond code generation, but replies on existing data sources. In contrast to these benchmarks, InfiBench benchmark is built for evaluating free-form question-answering ability in the code domain beyond code generation in an automated and model-independent way.

# 6 Conclusion

We proposed InfiBench, a systematic benchmark for evaluating the question-answering ability of code LLMs in real-world scenarios, to facilitate development and scientific evaluation of LLMs. InfiBench comprises 234 high-quality questions from Stack Overflow and supports automatic model-free evaluation. A comprehensive evaluation of over 100 code LLMs reveals several findings and takeaways. The benchmark is publicly available and continuously expanding.

# Acknowledgement

We thank ByteDance Inc. for the support on computing resources, anonymous reviewers for their constructive feedback, and Kaixin Li (National University of Singapore) for contributing the Docker image after the initial release of InfiBench. This work was partially supported by National Natural Science Foundation of China under Grant No. 62161146003, and the Tencent Foundation/XPLORER PRIZE. Tao Xie is also affiliated with the School of Computer Science, Peking University, China. The corresponding authors are Linyi Li and Tao Xie.

# References

- [1] Mistral AI. Codestral: Hello, world! | mistral ai | frontier ai in your hands. https://mistral.ai/news/codestral/, 2024.
- [2] Ben Athiwaratkun, Sanjay Krishna Gouda, Zijian Wang, Xiaopeng Li, Yuchen Tian, Ming Tan, Wasi Uddin Ahmad, Shiqi Wang, Qing Sun, Mingyue Shang, Sujan Kumar Gonugondla, Hantian Ding, Varun Kumar, Nathan Fulton, Arash Farahani, Siddhartha Jain, Robert Giaquinto, Haifeng Qian, Murali Krishna Ramanathan, Ramesh Nallapati, Baishakhi Ray, Parminder Bhatia, Sudipta Sengupta, Dan Roth, and Bing Xiang. Multi-lingual evaluation of code generation models. In *Eleventh International Conference on Learning Representations*, 2023.
- [3] Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*, 2021.
- [4] Edward Beeching, Clémentine Fourrier, Nathan Habib, Sheon Han, Nathan Lambert, Nazneen Rajani, Omar Sanseviero, Lewis Tunstall, and Thomas Wolf. Open Ilm leaderboard. https://huggingface.co/spaces/HuggingFaceH4/open\_llm\_leaderboard, 2023.
- [5] Loubna Ben Allal, Niklas Muennighoff, Logesh Kumar Umapathi, Ben Lipkin, and Leandro von Werra. A framework for the evaluation of code generation models. https://github.com/bigcode-project/bigcode-evaluation-harness, 2022.
- [6] Manish Bhatt, Sahana Chennabasappa, Yue Li, Cyrus Nikolaidis, Daniel Song, Shengye Wan, Faizan Ahmad, Cornelius Aschermann, Yaohui Chen, Dhaval Kapil, et al. Cyberseceval 2: A wide-ranging cybersecurity evaluation suite for large language models. *arXiv preprint arXiv:2404.13161*, 2024.
- [7] Xiao Bi, Deli Chen, Guanting Chen, Shanhuang Chen, Damai Dai, Chengqi Deng, Honghui Ding, Kai Dong, Qiushi Du, Zhe Fu, et al. Deepseek llm: Scaling open-source language models with longtermism. *arXiv* preprint arXiv:2401.02954, 2024.
- [8] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. In *Thirty-fourth International Conference on Neural Information Processing Systems*, 2020.
- [9] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- [10] Jasper Dekoninck, Mark Niklas Müller, and Martin Vechev. Constat: Performance-based contamination detection in large language models. arXiv preprint arXiv:2405.16281, 2024.
- [11] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805, 2018.
- [12] Angela Fan, Beliz Gokkaya, Mark Harman, Mitya Lyubarskiy, Shubho Sengupta, Shin Yoo, and Jie M Zhang. Large language models for software engineering: Survey and open problems. *arXiv preprint arXiv:2310.03533*, 2023.

- [13] Google Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- [14] Github. Github copilot your ai pair programmer. https://github.com/features/copilot, 2023.
- [15] Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, Y Wu, YK Li, et al. Deepseek-coder: When the large language model meets programming—the rise of code intelligence. *arXiv preprint arXiv:2401.14196*, 2024.
- [16] Dan Hendrycks, Steven Basart, Saurav Kadavath, Mantas Mazeika, Akul Arora, Ethan Guo, Collin Burns, Samir Puranik, Horace He, Dawn Song, et al. Measuring coding challenge competence with apps. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021.
- [17] Xinyi Hou, Yanjie Zhao, Yue Liu, Zhou Yang, Kailong Wang, Li Li, Xiapu Luo, David Lo, John Grundy, and Haoyu Wang. Large language models for software engineering: A systematic literature review. *arXiv preprint arXiv:2308.10620*, 2023.
- [18] Dongfu Jiang, Xiang Ren, and Bill Yuchen Lin. Llm-blender: Ensembling large language models with pairwise ranking and generative fusion. In *Sixty-first Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2023.
- [19] Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R Narasimhan. Swe-bench: Can language models resolve real-world github issues? In *Twelfth International Conference on Learning Representations*, 2024.
- [20] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. arXiv preprint arXiv:2001.08361, 2020.
- [21] Yuhang Lai, Chengxi Li, Yiming Wang, Tianyi Zhang, Ruiqi Zhong, Luke Zettlemoyer, Wen-tau Yih, Daniel Fried, Sida Wang, and Tao Yu. Ds-1000: A natural and reliable benchmark for data science code generation. In *Fortieth International Conference on Machine Learning*, 2023.
- [22] Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou, Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, et al. Starcoder: may the source be with you! *arXiv preprint arXiv:2305.06161*, 2023.
- [23] Tianle Li, Wei-Lin Chiang, Evan Frick, Lisa Dunlap, Tianhao Wu, Banghua Zhu, Joseph E Gonzalez, and Ion Stoica. From crowdsourced data to high-quality benchmarks: Arena-hard and benchbuilder pipeline. *arXiv preprint arXiv:2406.11939*, 2024.
- [24] Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, et al. Competition-level code generation with alphacode. *Science*, 378(6624):1092–1097, 2022.
- [25] Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, 2004.
- [26] Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. Is your code generated by chatGPT really correct? rigorous evaluation of large language models for code generation. In *Thirty-seventh International Conference on Neural Information Processing Systems*, 2023.
- [27] Tianyang Liu, Canwen Xu, and Julian McAuley. Repobench: Benchmarking repository-level code auto-completion systems. *arXiv* preprint arXiv:2306.03091, 2023.
- [28] Anton Lozhkov, Raymond Li, Loubna Ben Allal, Federico Cassano, Joel Lamy-Poirier, Nouamane Tazi, Ao Tang, Dmytro Pykhtar, Jiawei Liu, Yuxiang Wei, et al. Starcoder 2 and the stack v2: The next generation. *arXiv preprint arXiv:2402.19173*, 2024.

- [29] Shuai Lu, Daya Guo, Shuo Ren, Junjie Huang, Alexey Svyatkovskiy, Ambrosio Blanco, Colin Clement, Dawn Drain, Daxin Jiang, Duyu Tang, et al. Codexglue: A machine learning benchmark dataset for code understanding and generation. In *Thirty-fifth International Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1)*, 2021.
- [30] Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xiubo Geng, Wenxiang Hu, Chongyang Tao, Jing Ma, Qingwei Lin, and Daxin Jiang. Wizardcoder: Empowering code large language models with evol-instruct. *arXiv preprint arXiv:2306.08568*, 2023.
- [31] Justus Mattern, Fatemehsadat Mireshghallah, Zhijing Jin, Bernhard Schoelkopf, Mrinmaya Sachan, and Taylor Berg-Kirkpatrick. Membership inference attacks against language models via neighbourhood comparison. In *Findings of the Association for Computational Linguistics: ACL* 2023, 2023.
- [32] Alexandre Matton, Tom Sherborne, Dennis Aumiller, Elena Tommasone, Milad Alizadeh, Jingyi He, Raymond Ma, Maxime Voisin, Ellen Gilsenan-McMahon, and Matthias Gallé. On leakage of code generation evaluation datasets. *arXiv preprint arXiv:2407.07565*, 2024.
- [33] Niklas Muennighoff, Qian Liu, Armel Zebaze, Qinkai Zheng, Binyuan Hui, Terry Yue Zhuo, Swayam Singh, Xiangru Tang, Leandro von Werra, and Shayne Longpre. Octopack: Instruction tuning code large language models. *arXiv preprint arXiv:2308.07124*, 2023.
- [34] Niklas Muennighoff, Alexander Rush, Boaz Barak, Teven Le Scao, Nouamane Tazi, Aleksandra Piktus, Sampo Pyysalo, Thomas Wolf, and Colin A Raffel. Scaling data-constrained language models. In *Thirty-seventh International Conference on Neural Information Processing Systems*, 2024.
- [35] Erik Nijkamp, Bo Pang, Hiroaki Hayashi, Lifu Tu, Huan Wang, Yingbo Zhou, Silvio Savarese, and Caiming Xiong. Codegen: An open large language model for code with multi-turn program synthesis. In *Eleventh International Conference on Learning Representations*, 2023.
- [36] OpenAI. Gpt-4 technical report. OpenAI, 2023.
- [37] OpenAI. Hello gpt-40 | openai. https://openai.com/index/hello-gpt-40/, 2024.
- [38] OpenAI. Introducing swe-bench verified | openai. https://openai.com/index/introducing-swe-bench-verified/, 2024.
- [39] Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, et al. Code llama: Open foundation models for code. *arXiv preprint arXiv:2308.12950*, 2023.
- [40] Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference attacks against machine learning models. In *Thirty-eight IEEE Symposium on Security and Privacy*, 2017.
- [41] StackExchange. All sites stackexchange. https://stackexchange.com/sites?view= list#users, 2024.
- [42] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- [43] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- [44] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Thirty-first International Conference on Neural Information Processing Systems*, 2017.
- [45] Boxin Wang, Weixin Chen, Hengzhi Pei, Chulin Xie, Mintong Kang, Chenhui Zhang, Chejian Xu, Zidi Xiong, Ritik Dutta, Rylan Schaeffer, et al. Decodingtrust: A comprehensive assessment of trustworthiness in gpt models. In *Thirty-seventh International Conference on Neural Information Processing Systems*, 2024.

- [46] Ruijie Xu, Zengzhi Wang, Run-Ze Fan, and Pengfei Liu. Benchmarking benchmark leakage in large language models. *arXiv preprint arXiv:2404.18824*, 2024.
- [47] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. In *Thirty-seventh International Conference on Neural Information Processing Systems*, 2024.
- [48] Qinkai Zheng, Xiao Xia, Xu Zou, Yuxiao Dong, Shan Wang, Yufei Xue, Zihan Wang, Lei Shen, Andi Wang, Yang Li, Teng Su, Zhilin Yang, and Jie Tang. Codegeex: A pre-trained model for code generation with multilingual evaluations on humaneval-x. In *Twenty-ninth ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2023.
- [49] Qihao Zhu, Daya Guo, Zhihong Shao, Dejian Yang, Peiyi Wang, Runxin Xu, Y Wu, Yukun Li, Huazuo Gao, Shirong Ma, et al. Deepseek-coder-v2: Breaking the barrier of closed-source models in code intelligence. *arXiv preprint arXiv:2406.11931*, 2024.

# Checklist

- 1. For all authors...
  - (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
  - (b) Did you describe the limitations of your work? [Yes] Limitations are discussed throughout Section 2 and specifically in Appendix B.
  - (c) Did you discuss any potential negative societal impacts of your work? [Yes] Societal impacts are discussed in Appendix B.
  - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
- 2. If you are including theoretical results...
  - (a) Did you state the full set of assumptions of all theoretical results? [N/A]
  - (b) Did you include complete proofs of all theoretical results? [N/A]
- 3. If you ran experiments (e.g. for benchmarks)...
  - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] All code and data are publicly available at https://infi-coder.github.io/infibench along with instructions needed to reproduce. The accessibility information is also available in detail in Appendix A.
  - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [N/A] This work does not involve model training. The inference hyperparameters are listed in Section 3 and ablation studies are presented in Appendix G.
  - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes] All experiments are repeated three times whenever budget and computing resource permit. Error bars and standard deviations are reported.
  - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] The information is included in Section 3.
- 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
  - (a) If your work uses existing assets, did you cite the creators? [Yes]
  - (b) Did you mention the license of the assets? [Yes] The main assets are from Stack Overflow which is open source under CC BY-SA 4.0 license. We inherit this license to release.
  - (c) Did you include any new assets either in the supplemental material or as a URL? [Yes] New assets (the benchmark and evaluation tool) is accessible through https://infi-coder.github.io/infibench.

- (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [Yes] We release the new asset inheriting the CC BY-SA 4.0 license as described in Section 1 and Appendix A.
- (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [Yes] Domain experts are required to remove such information by paraphrasing when constructing the benchmark.
- 5. If you used crowdsourcing or conducted research with human subjects...
  - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
  - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
  - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

# **Appendices**

In appendices, we present dataset card and data accessibility details in Appendix A, discussion on limitations and societal impact in Appendix B, agreement statistics between InfiBench/GPT-40 and human in Appendix C, question grouping by difficulty in Appendix D, full leaderboard in Appendix E, additional findings in Appendix F, study of sampling hyperparameters in Appendix G, prompts in Appendix H, and benchmark data examples in Appendix I.

# A Dataset Card and Accessibility Details

#### **Dataset Card**

- Name: InfiBench
- **Description**: Evaluation Dataset for the Question-Answering Capabilities of Code Large Language Models
- URL: https://infi-coder.github.io/infibench (all resources) / https://huggingface.co/datasets/llylly001/InfiBench (data part)
- Version: 2.1
- License: Creative Commons Attribution Share Alike 4.0
- Citation:

```
@misc{infibench,
    title={InfiBench: Evaluating the Question-Answering Capabilities
    of Code Large Language Models},
    howpublished = "\url{https://infi-coder.github.io/infibench}",
    author={InfiBench},
    year={2024}
}
```

- DOI: doi:10.57967/hf/2474
- Responsible AI Data Collection:

Data source is downloaded from the publicly available StackExchange archive (https://archive.org/download/stackexchange, https://ia904700.us.archive.org/view\_archive.php?archive=/6/items/stackexchange/stackoverflow.com-Posts.7z). Especially, we use the preprocessed version from https://huggingface.co/datasets/mikex86/stackoverflow-posts where all posts are formatted in Markdown text.

We choose to keep only the questions with at least three positively voted answers and an officially accepted answer, which turn out to be 1,090,238 questions. For these one million questions, we choose to keep frequently viewed and relatively new questions.

Utilizing the mandatory question tags of these questions, we then manually construct a tag tree that covers the 200 most frequent tags, enabling us to identify the top programming languages and areas for 14,330105 out of these 17,402 questions. We exclude 6 programming languages that either describe data or are domain-specific: JSON, regex, Markdown, YAML, CSV, and SQL. As a result, we compile 13,854 questions that serve as the initial seed set.

We randomly sample from the initial seed set. Then we recruited five domain experts inside our company to create the benchmark from the sampled initial seed set, each in charge of one area. The annotation process is composed of three steps: (1) Question Selection and Type Annotation; (2) Prompt Paraphrasing. (3) Correctness Criterion Annotation.

#### • Responsible AI — Data Biases:

The data essentially serves as an evaluation benchmark. We foresee data biases in the following aspects:

(1) Non-standard evaluation. Alongside the data is a comprehensive benchmark of existing code LLMs. The benchmark scores are evaluated under a specific set of hyperparameters

(e.g, temperature 0.2, top probability 0.9, best@10 at question level). Data usage under different evaluation conditions may result in misleading comparison results and conclusions. (2) Usage misinterpretation. The benchmark focuses on evaluating the response correctness of code LLMs for a set of real-world developers' questions. Our evaluation standard does not specifically take other aspects (naturalness, conciseness, fairness, politeness, etc) into consideration. Hence, this is risk of overinterpreting the evaluation results. When evaluating a code LLM, we recommend combining this benchmark score with other evaluations to be

(3) Potential data contamination. Though we have made our efforts to reduce the impact of data contamination, future code LLMs may train or fine-tune on this benchmark dataset to improve the score on InfiBench. This could be challenging to prevent as a cost of being fully public. On the other hand, as responsible LLM developers, we hope future practitioners would report how they use the benchmark data if beyond the original scope (for evaluation use).

a more comprehensive evaluation.

 Responsible AI — Personal Sensitive Information: During the data construction process, our domain experts paraphrased the question prompts to remove personal and sensitive information (PII) and a cross validation stage was introduced to further ensure the PII removal.

Croissant Dataset Description: https://huggingface.co/datasets/llylly001/InfiBench/blob/main/croissant-infibench.json. Note that the Croissant format is mainly designed for machine learning dataset description. However, InfiBench is more than a dataset; it is an evaluation benchmark including response evaluation standards, tools, and an accompanying leaderboard. Hence, the Croissant script records only the CSV file and covers question prompts and evaluation standards; whereas the open-source evaluation tool and leaderboard are not recorded which can be separately downloaded from https://infi-coder.github.io/infibench.

Data Accessibility. As briefly mentioned in the main text, all materials are made publicly available and accessible at the website: https://infi-coder.github.io/infibench without personal request. The materials include three parts: (1) Benchmark questions and evaluation metrics this part is additionally uploaded to Hugging Face (URL and DOI are in the above dataset card). (2) Automatic evaluation tool — this part is uploaded and maintained in a dedicated GitHub repo https://github.com/infi-coder/infibench-evaluator. In addition, we uploaded our extension of bigcode-evaluation-harness [5], namely infibench-evaluation-harness to a  $dedicated\ Git\ \bar{H}ub\ repo\ https://github.com/infi-coder/infibench-evaluation-harness.$ The extension includes the inference code on InfiBench for all evaluated LLMs. (3) Evaluation raw data and leaderboard — the leaderboard is displayed on the website https://infi-coder. github.io/infibench and the raw model responses are stored in the website repo https: //github.com/infi-coder/infibench. All materials are under the Creative Commons Attribution Share Alike 4.0 license. In the above dataset card and Appendix B, we anticipate potential inappropriate usage of the benchmark and we encourage the practitioners to document their usage of the benchmark if beyond model evaluation. In the future, we will continue the maintenance and expansion of the benchmark. Furthermore, we are developing an adaptor for automatic evaluation on Hugging Face so that InfiBench can be integrated into the Hugging Face Open LLM Leaderboard [4] to further ease the evaluation burden.

# **B** Limitations, Societal Impacts, and Future Work

In this appendix, we expand our discussion of limitations, potential societal impacts, and future work.

**Evaluation Metric.** In InfiBench, the expert-annotated evaluation metric is designed to mainly focus on response correctness, more specifically, whether the response contains key information that solves the given question. Concretely, the metric may evaluate whether the response passes a given set of unit tests, whether it suggests the right API or concept, whether it follows the instruction to provide relevant information, etc. Hence, the score comes with two limitations: (1) The score is subjective since the metric is annotated by human experts without an explicit and universal

Table 6: Confusion matrices between InfiBench/GPT-40 and human. Details in Appendix C. Bolded cells correspond to when both methods have clear preferences on one response.

#### (a) Between InfiBench and human.

#### (b) Betwen GPT-40 and human.

				Human		
		A > B	B > A	$A \approx B \uparrow$	$A \approx B \downarrow$	Tot.
	A > B	23	3	9	4	39
	B > A	4	17	12	2	35
InfiBench	$A \approx B \uparrow$	0	0	10	0	10
	$A \approx B \downarrow$	4	3	3	6	16
	Tot.	31	23	34	12	100

				Human		
		A > B	B > A	$A \approx B \uparrow$	$A \approx B \downarrow$	Tot.
	A > B	23	7	8	6	44
	B > A	3	12	9	3	27
GPT-40	$A \approx B \uparrow$	5	4	15	3	27
	$A \approx B \downarrow$	0	0	2	0	2
	Tot.	31	23	34	12	100

standard. Note that we did not aim to provide an objective metric since the developers' views of response correctness intrinsically vary and diverge for these diverse questions. On the other hand, we introduce a cross-validation and calibration stage to improve the metric representativeness of most developers' standards. We leave it as a future work to further quantitatively measure and improve the metric representativeness. (2) The score focuses mainly on correctness. Several other aspects define a model's usability, such as language naturalness (including conciseness, politeness, etc), trustworthiness (refusal of risky questions, fairness, unbiasedness, privacy, etc), and system-level metrics (latency, throughput, parallelism-friendliness, etc). Model evaluators and practitioners may keep in mind that InfiBench score is not a comprehensive usability measurement of code LLMs, and we strongly encourage them to combine InfiBench score with benchmarks on these other aspects (c.f. [6, 45]) to comprehensively evaluate LLMs.

**Data Contamination.** The limitations and mitigations on data contamination are discussed in Section 2.5. In addition, as a side effect of open source, future code LLMs may leverage the benchmark data to deliberately introduce data contamination to achieve a high score in InfiBench. To partly detect such data contamination, our evaluation of using the original stack Overflow answers might be a proxy. According to Table 4(a), even gold extraction from human answers cannot saturate the benchmark while strong LLMs like GPT-4 surpassed human answers. Hence, if a future model achieves scores close to human answers (between 50% and 65%) but cannot further improve beyond human along with scaling, data contamination may potentially happen. Detecting data contamination is itself a research topic where research on member inference attacks [40, 31] is involved. We did not integrate a detection module in the current release of InfiBench but we are planning to inspect this topic in the future.

Labelling Cost. InfiBench construction involves human labelling cost, where domain experts paraphrase the source question post and label the evaluation metric. Such a cost prevents the InfiBench from scaling up in terms of size, and the questions for less popular programming languages, such as Rust and Ruby, are relatively few. In an attempt to mitigate this limitation, we explored a few alternative evaluation metrics, such as dialogue similarity with officially accepted answers. However, these alternatives either require a language model which may induce bias and heavy computing cost, or deviate away from domain experts' correctness judgment. We leave the exploration of more scalable metrics and annotation procedures as future work and make the benchmark fully open source so community involvement may boost the expansion.

# C Agreement Statistics between InfiBench/GPT-40 Evaluation and Human

In Section 2.4, we evaluated the alignment between InfiBench/GPT-40 evaluation and human evaluation by generating 100 response pairs for InfiBench questions and let InfiBench, GPT-40, and human annotators to grade into four outcomes.

Table 6 shows the confusion matrices between InfiBench/GPT-40 and human, where each cell corresponds to the frequency of each combination of outcomes among 100 pairs. The implication of each outcome is introduced in Section 2.4.

Learned from Table 6, if we only count the cases where both human and InfiBench have clear preferences, their agreement rate is  $\frac{40}{47} = 85.1\%$ ; if we only count the cases where both human and GPT-40 have clear preferences, their agreement rate is  $\frac{35}{45} = 77.8\%$ . Hence, the InfiBench evaluation aligns with human experts better than the GPT-40 evaluation (with >80% confidence). Furthermore,

we observe that GPT-40 has a stronger opinion and tends to choose one response more often, so it falls short when A and B are both bad responses, labelling none of them as "both bad". We also observe that InfiBench evaluation could be too strict due to pattern matching and fixed post-processing leading to over-differentiation—when a human believes A and B are both good responses, with only a 29.4% chance InfiBench labels them as "both good".

# **D** Difficulty Grouping

We systematically evaluated GPT-4 and GPT-3.5-turbo on the benchmark following the evaluation protocol in Section 3, based on which we classify the benchmark questions into five disjoint difficulty groups.

- Level 1 (93 questions, 39.7%): GPT-3.5-turbo can achieve a mean score  $\geq 0.5$ .
- Level 2 (55 questions, 23.5%): Among the rest questions, those where GPT-4's mean score  $\geq 0.5$ .
- Level 3 (44 questions, 18.8%): Among the rest questions, those where GPT-4 with sampling temperature 1.0 can achieve a maximum score ≥ 0.5 among 10 trials.
- Level 4 (18 questions, 7.7%): Among the rest questions, those GPT-4 with sampling temperature 0.2 can achieve a positive score among 100 trials.
- Level 5 (24 questions, 10.3%): The remaining questions, i.e., GPT-4 cannot get score among 100 trials.

Appendix E shows each code LLM's score in each difficulty group. The mean scores strictly decrease for higher difficulty levels, highlighting that the question difficulty is in general consistent across different code LLMs and our group assignment is reasonable. We hope that the grouping can help better reveal the strengths and weaknesses of a code LLM for different questions.

Question examples by difficulty groups are in Appendix I.

#### **E Evaluation Details and Full Benchmark Results**

**Evaluation Details of Code LLMs.** For proprietary model evaluation, we did not specify the max tokens to generate and found out that the longest response generated by GPT-4 has 662 tokens with Code Llama tokenizer.

For open-source model evaluation, for models with over 30B parameters, due to the GPU memory limit and efficiency concerns, we impose the longest context constraint of 4,096 tokens and experiment just once. Since there is only one question whose GPT-4 context (prompt + GPT-4 response) can exceed 4,096 tokens, we think this context constraint has little effect, reducing the score by 0.37% at most. For models within 30B parameters, since GPT-4 response has at most 662 tokens, we set the max number of tokens to generate to be  $\min\{1024, \text{context length - prompt length}\}$ , providing some wiggle room. Meanwhile, we repeat the evaluation three times for models within 30B parameters.

Evaluation Details of Original Stack Overflow Answers. As listed in Table 4(a) and Table 7, besides evaluating LLM responses, we evaluated the score of human-written original Stack Overflow answers since the question prompts are paraphrased from Stack Overflow. We consider three settings: (1) evaluating the officially-accepted answer post (note that we select only the Stack Overflow questions with an officially-accepted answer into the benchmark); (2) evaluating the highest-voted answer post (note that any registered user can equally vote for or against an answer); and (3) evaluating the highest-voted answer posts up to 10 and recording the highest score achieved by any post. For the last setting, we chose the number 10 because the main evaluation metric of model response is best@10. Moreover, we observe that all officially accepted answers for InfiBench questions are among the top 10 highest-voted answer posts. Note that there is no randomness of scores from Stack Overflow answers, so we do not repeat the evaluation nor report the standard deviation.

As expected, the last setting achieves the highest score 65.18% among the three settings. Due to its consistency with models' evaluation metric best@10, we deem this score most comparable with scores from LLMs. Interestingly, when considering only one answer post, the second setting, selecting the highest-voted answer, is better than the first setting, selecting the officially accepted answer.

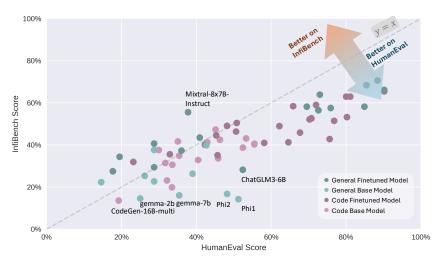


Figure 5: InfiBench and HumanEval scores as a scatter plot for LLMs. r=0.8058. Discussion in Appendix F.1.

**Full Benchmark Results.** We present the full leaderboard in Table 7 (by descending order of InfiBench scores) and Table 8 (by alphabetical order of model family names). These tables are expanded from the aggregated Table 4. In these tables, we show model properties including size and context length. We also present HumanEval [3] scores since HumanEval is one of the most widely used benchmarks for evaluating code LLMs (further discussion in Appendix F). Furthermore, we represent the score breakdown by difficulty levels, problem types, and evaluation metric types. The proportion of each difficulty level can be found in Appendix D, and the proportion of each problem type and evaluation metric type is shown in Table 3(a,b). InfiBench score can be computed by the weighted sum of breakdown subscores by proportions. We present the score of human-written original Stack Overflow answers in the last three rows.

In tables, the mean scores are computed from scores of all 106 code LLMs. We observe that the mean overall score, 37.82%, is still much inferior to human answers (which achieves over 50% even with just one attempt). The model performance is monotonically decreasing for higher difficulty levels; relatively equivalent across different problem types; and weaker under blank-filling and dialogue-similarity metrics than keyword-matching and unit-testing metrics.

# F Additional Findings and Discussion

In this appendix, we present additional findings and discussion that are omitted from Section 3.

#### F.1 Correlations between InfiBench and HumanEval Scores

We study the correlation between InfiBench and HumanEval pass@1 scores for different LLMs. In Figure 5, we plot LLMs with both InfiBench and HumanEval scores, in total 66 LLMs, in Table 7 as a scatter plot. The figure shows that scores on the two benchmarks are generally positively correlated, with a Pearson correlation coefficient r=0.8058. If conducting a linear regression, we would observe that different model types (i.e., general/code model, base/finetuned model) share almost the same linear relationship, indicating that both benchmarks can reflect the model capability in general. Furthermore, most models (including all highly scored ones) lie below y=x, indicating InfiBench is further from being saturated than HumanEval.

However, a few outlier models exist in Figure 5. Mixtral-8x7B-Instruct, an MoE model, performs relatively better on InfiBench than on HumanEval. Some other models, e.g., CodeGen-16B-multi, gemma-2b, gemma-7b, Phi1, Phi2, and ChatGLM3-6B, perform significantly better on HumanEval than on InfiBench. These models are relatively small or old-dated. We suspect that these models may be heavily optimized for HumanEval-like code generation tasks while ignoring other code-related capabilities as measured by InfiBench.

Table 7: Full leaderboard of all benchmarked LLMs ranked by InfiBench scores. Evaluation protocol in Section 3 and details explained in Appendix E. Icon "a" stands for proprietary models otherwise open-source. As a reference, HumanEval scores digested from [26] and each model's report are shown. Bar colors stand for General Base, General Finetuned, Code Base, and Code Finetuned models respectively. Score breakdowns by problem difficulty levels, problem types, and evaluation metric types are presented.

Rank	Model Family	Model Name	Size (# Param.)	Context Length	InfiBench Score	HumanEval	Level 1	Di Level 2	ifficulty Lev Level 3	els Level 4	Level 5	Code Completion	Proble Code Debugging	m Type Knowledge QA	Config & Env Debugging	Keyword Matching	Evaluation Unit Testing	Metric Type Blank Filling	Dialogue Similarity
1	⊕GPT-4	GPT-4-0613	?	8192	$70.64\% \pm 0.82\%$	88.4	92.31%	92.48%	51.90%	31.91%	0.00%	75.23%	69.74%	68.55%	66.63%	66.61%	76.00%	58.08%	84.27%
2	GPT-4 GPT-4	GPT-4-turbo-1106	?	8192	68.42% ± 0.38% 66.19%	85.4	89.90%	78.57%	54.16%	30.93% 28.05%	16.20%	74.82% 75.00%	65.36%	67.47% 65.65%	62.98%	64.98%	76.40%	53.91% 49.74%	52.85% 70.73%
4	DeepSeek Coder	GPT-4o-2024-05-13 deepSeek-coder-V2-instruct	236B / 21B	8192 128000	65.49%	90.2	88.77%	78.66% 76.97%	46.43% 50.58%	17.31%	5.21% 12.50%	74.77%	59.32% 63.89%	59.57%	61.70% 59.30%	61.59% 58.91%	76.00% 76.00%	55.77%	37.34%
5	Claude 3	Claude 3 Opus	? 22B	200000 32768	63.89% 62.98% ± 0.56%	73	84.36%	78.95%	39.98%	31.76% 17.11%	18.06%	65.18%	62.94% 63.65%	65.86% 61.07%	60.49%	60.07%	61.80%	59.36% 45.92%	44.91% 57.08%
7	Mistral Open DeepSeek Coder	Codestral-22b deepseek-coder-33b-instruct	33B	16384	62.96%	80.02	88.64% 87.58%	69.90% 72.02%	49.97% 44.12%	15.83%	5.90% 16.67%	68.75% 71.26%	57.14%	63.14%	54.28% 56.81%	57.72% 59.01%	73.33% 77.00%	30.00%	36.09%
8	Phind	Phind-CodeLlama-34B-v2	34B	4096	59.00% 58.47%	71.95	83.67%	55.57%	53.12%	15.09%	14.93%	58.24%	58.30%	63.60%	55.33%	59.63%	58.40%	35.26%	24.19%
9 10	Phind Mistral	Phind-CodeLlama-34B-v1 mistral-large	34B ?	4096 32768	58.47% 58.22%	65.85	81.38% 81.76%	63.85% 66.59%	47.05% 41.66%	22.63%	5.21% 4.17%	66.13% 66.69%	56.94% 50.10%	56.79% 60.21%	49.48% 52.89%	55.71% 53.17%	66.00% 67.00%	38.78% 45.64%	35.39% 42.66%
11	Claude 3	Claude 3 Sonnet	?	200000	58.20%	84.9	80.13%	65.55%	42.48%	18.06%	15.28%	62.61%	52.34%	63.61%	52.12%	54.22%	66.00%	46.35%	25.62%
12 13	Claude 3 DeepSeek LLM	Claude 3 Haiku deepseek-llm-67b-chat	? 67B	200000 4096	57.57% 57.41%	75.9	79.86%	66.06% 63.03%	40.23% 39.09%	21.76% 22.60%	10.42% 5.21%	61.71%	48.68% 52.73%	62.85% 58.72%	56.71% 55.63%	55.78% 53.14%	58.40% 63.00%	44.62% 51.41%	36.40% 36.68%
14	⊕GPT-3.5	GPT-3.5-turbo-0613	?	4096	56.47% ± 1.34%	72.6	93.08%	49.77%	31.36%	14.30%	7.64%	64.91%	48.50%	59.47%	49.64%	51.28%	70.07%	40.90%	40.13%
15 16	Mistral Open	mistral-small mixtral-8x7B-Instruct	? 46.7B / 12.9B	32768 32768	55.62% ± 0.46% 55.55%	37.8	82.98% 82.19%	55.98% 56.72%	35.72% 31.53%	22.58% 24.00%	10.07% 17.36%	63.56% 54.01%	44.12% 51.57%	64.13% 63.69%	47.75% 53.59%	50.56% 56.14%	68.00% 50.40%	39.08% 35.58%	53.32% 61.75%
17	Qwen	Qwen-72B	72B	32768	55.34%	37.8	81.98%	57.40%	41.61%	13.24%	4.17%	61.06%	53.16%	58.79%	44.03%	50.43%	64.00%	45.96%	36.41%
18 19	DeepSeek Coder	deepseek-coder-6.7b-instruct Qwen-72B-Chat	6.7B 72B	16384 32768	53.25% ± 0.40% 52.97%	80.22	77.88% 82.44%	56.30% 47.00%	35.18% 36.09%	18.89% 18.34%	9.72% 9.38%	65.95% 58.67%	46.44% 45.81%	52.46% 60.12%	42.12% 44.31%	48.24% 49.26%	70.40% 59.00%	26.90% 43.08%	23.48% 33.95%
20	Qwen Magicoder	Qwen-72B-Chat Magicoder-S-CL-7B	72B 7B	16384	52.71% ± 0.72%	70.7	77.97%	50.42%	40.20%	13.45%	9.38% 12.50%	51.39%	51.98%	56.97%	50.58%	49.26% 53.28%	56.67%	21.41%	33.95% 26.97%
21	WizardLM	WizardCoder-Python-34B-V1.0	34B	16384	52.59% 52.17%	70.73	78.51%	52.50%	34.25%	20.05%	10.42%	60.32%	46.39%	55.86%	44.01%	48.73%	64.00%	37.56%	24.72%
22 23	Phind Magicoder	Phind-CodeLlama-34B-Python-v1 Magicoder-S-DS-6.7B	34B 6.7B	4096 16384	51.46% ± 1.09%	70.22 76.8	80.54% 78.93%	48.44% 51.02%	42.58% 28.91%	8.57% 25.93%	1.04% 6.48%	54.41% 62.54%	52.34% 46.45%	57.11% 55.74%	41.47% 33.84%	51.04% 45.64%	57.80% 69.13%	27.18% 31.45%	39.76% 27.86%
24	Code Llama	CodeLlama-34b-Instruct	34B	16384	50.45%	50.79	72.60%	55.07%	33.16%	18.43%	9.72%	51.71%	48.37%	61.36%	37.04%	48.14%	51.20%	47.76%	28.55%
25 26	01.AI WizardLM	Yi-34B-Chat WizardCoder-Python-7B-V1.0	34B 7B	4096 16384	49.58% 49.10% ± 1.59%	48.2	76.81%	47.15% 48.08%	29.32%	26.39% 12.50%	4.17% 9.72%	44.10% 58.60%	44.75%	62.29% 50.67%	49.84%	53.14% 46.38%	35.40% 59.40%	36.15% 25.30%	33.07% 23.00%
27	WizardLM	WizardCoder-Python-13B-V1.0	13B	16384	48.99% ± 0.92%	62.19	76.21%	46.76%	34.19%	16.17%	0.35%	52.69%	48.29%	50.67%	41.32%	48.71%	53.73%	20.45%	29.61%
28 29	Code Llama Code Llama	CodeLlama-34b CodeLlama-13b-Instruct	34B 13B	16384 16384	47.36% 46.37% ± 1.26%	45.11 50.6	72.07% 69.07%	43.34% 45.99%	29.32% 34.37%	21.20% 11.42%	13.54% 7.52%	53.74% 48.65%	50.09% 45.18%	51.52% 49.67%	26.59% 39.83%	43.18% 47.71%	57.33% 50.47%	37.37% 20.90%	24.85% 12.45%
30	Zephyr	Zephyr 7B beta	7B	32768	46.31% ± 1.11%	/	68.41%	49.99%	31.11%	14.99%	3.59%	44.26%	44.86%	54.89%	40.85%	49.28%	35.07%	27.91%	27.66%
31 32	StarCoder2 DeepSeek MoE	15B-Instruct deepseek-moe-16b-chat	15B 16B / 2.8B	16384 16384	45.89% ± 0.95% 45.18% ± 1.65%	67.7	70.37% 68.15%	50.21%	24.15% 27.55%	11.44%	6.83%	56.02% 47.19%	38.52% 46.54%	46.30% 45.58%	38.56%	40.55% 45.71%	60.27% 44.73%	25.21% 25.85%	45.01% 20.70%
33	OctoPack	OctoCoder	15.5B	8192	44.55% ± 0.79%	45.3	68.19%	41.61%	29.39%	12.96%	11.11%	46.56%	37.62%	53.57%	39.56%	44.18%	47.07%	20.09%	39.20%
34 35	Qwen Qwen	Qwen-14B Qwen-14B-Chat	14B 14B	8192 8192	43.69% ± 1.09% 43.49% ± 0.63%	40.9	67.61% 68.91%	47.64% 36.25%	21.87% 27.73%	9.63% 10.28%	7.52% 15.39%	44.59% 45.39%	42.15% 42.12%	47.09% 46.33%	39.99% 38.48%	41.61% 41.87%	44.40% 42.73%	34.19% 36.18%	28.21% 34.79%
36	Qwen Magicoder	Qwen-14B-Chat Magicoder-DS-6.7B	6.7B	16384	43.47% ± 0.21%	/	67.04%	48.33%	23.11%	13.64%	0.69%	45.39% 52.73%	40.42%	48.14%	25.61%	38.37%	42.73% 56.73%	29.81%	38.07%
37	Code Llama	CodeLlama-34b-Python CodeLlama-70b-Instruct	34B	16384	43.13% 42.82%	53.29	66.02%	40.76%	36.06%	6.94%	0.00%	50.14%	40.48%	43.64%	34.13%	40.40%	51.00%	27.63%	16.67%
38 39	Code Llama StarCoder2	CodeLlama-70b-Instruct 15B	70B 15B	4096 16384	42.52% ± 1.24%	75.6 46.3	59.08% 64.99%	44.14% 41.67%	38.48% 29.02%	12.22% 13.77%	7.64% 3.70%	38.20% 47.00%	44.99% 37.31%	46.87% 46.76%	42.38% 36.87%	48.34% 43.86%	32.00% 42.20%	16.09% 18.44%	5.62% 0.00%
40	Magicoder	Magicoder-CL-7B	7B	16384	41.71% ± 0.76% 41.66% ± 0.84%	1	70.38%	36.48%	23.06%	10.33%	0.35%	49.26%	35.11%	45.41%	33.47%	37.85%	52.27%	19.91%	39.21%
41 42	Code Llama DeepSeek Coder	CodeLlama-13b deepseek-coder-1.3b-instruct	13B 1.3B	16384 16384	41.05% ± 0.84% 41.32% ± 1.12%	35.07 64.6	62.77% 65.48%	40.40% 41.42%	31.11% 25.48%	7.97% 6.30%	7.41% 2.78%	38.17% 41.91%	44.56% 42.56%	43.00% 42.88%	41.72% 36.38%	45.44% 41.80%	34.80% 42.20%	14.79% 16.52%	2.47%
43	Code Llama	CodeLlama-13b-Python	13B	16384	$41.31\% \pm 0.90\%$	42.89	62.93%	40.80%	28.61%	10.37%	5.21%	49.95%	44.60%	36.68%	27.22%	40.58%	51.07%	11.92%	13.64%
44 45	WizardLM Mistral	WizardCoder-15B-V1.0 mistral-medium	15B	2048 32768	41.01% ± 0.22% 40.95% ± 0.41%	58.12	66.19% 72.59%	40.34% 30.34%	21.72% 19.14%	12.42% 8.15%	1.74% 7.29%	44.80% 41.49%	34.54%	47.68% 49.19%	35.29% 39.09%	38.43% 38.54%	47.60% 42.67%	22.31% 33.85%	35.01% 18.26%
46	gemma	gemma-7b-it	7B	8192	40.68% ± 1.23%	28.7	60.94%	42.94%	28.86%	5.75%	4.86%	42.60%	36.37%	47.75%	34.52%	40.68%	41.40%	19.04%	30.44%
47 48	Code Llama Code Llama	CodeLlama-70b CodeLlama-70b-Python	70B 70B	4096 4096	40.60%	55.5 55.49	60.59% 59.14%	37.42% 36.07%	35.68% 41.06%	7.59% 7.59%	4.17% 0.00%	47.18% 42.03%	39.10% 43.04%	39.09% 40.76%	33.21% 32.46%	40.54% 41.78%	45.00% 41.00%	19.23% 10.96%	8.56% 19.50%
48	OctoPack	OctoGeeX	70В 6В	8192	40.14% ± 1.55%	42.28	62.54%	37.84%	26.39%	15.67%	2.20%	42.03%	33.23%	46.02%	39.10%	39.85%	39.96%	20.90%	31.11%
50	DeepSeek LLM	deepseek-llm-67b-base	67B 70B	4096 4096	39.87% 39.30%	42.7	57.15% 56.95%	48.73%	24.32%	9.17%	4.17%	35.50%	43.17%	46.15% 42.56%	34.40%	39.98%	36.00%	30.00%	24.46%
51 52	Llama 2 DeepSeek Coder	Llama2-70B-Chat deepseek-coder-33b-base	70B 33B	16384	38.75%	52.45	56.73%	38.02% 44.55%	33.71% 19.85%	7.96% 14.95%	7.64% 8.33%	35.65% 33.36%	42.87% 43.73%	46.06%	36.11% 31.23%	40.89% 43.99%	34.40% 25.50%	22.44% 14.49%	28.14% 28.02%
53	01.AI	Yi-6B-Chat	6B	4096 4096	38.14% ± 0.58% 37.69%	/	52.73%	38.20%	34.37%	12.53%	7.64%	33.36%	39.81%	42.54%	38.33%	43.26%	23.83%	15.32%	15.69%
54 55	Llama 2 Code Llama	Llama2-70B CodeLlama-7b	70B 7B	4096 16384	37.69% 37.62% ± 1.28%	28.7 29.98	51.51% 59.81%	42.58% 38.25%	28.48% 19.37%	10.19% 9.32%	10.42% 4.86%	36.26% 42.19%	42.99% 38.60%	37.12% 37.37%	32.98% 28.41%	39.52% 37.87%	28.00% 41.80%	30.45% 15.13%	0.00%
56	Mistral Open	Mistral-7B-Instruct-v0.1	7B	32768	37.55% ± 1.10%	/	56.31%	41.34%	24.07%	7.47%	3.47%	39.74%	30.74%	47.10%	31.40%	34.17%	39.80%	34.44%	29.90%
57 58	InternLM Owen	InternLM-Chat-20B Qwen-7B-Chat	20B 7B	16384 32768	37.41% ± 0.75% 37.36% ± 1.29%	36	59.98% 60.23%	32.30% 36.20%	20.40% 19.77%	18.44% 7.65%	7.06% 5.90%	45.38% 43.44%	34.67% 32.38%	34.25% 38.22%	31.63% 32.98%	34.51% 34.06%	46.20% 43.07%	18.18% 29.02%	23.51% 30.11%
59	DeepSeek LLM	deepseek-llm-7b-chat	7B	4096	36.75% ± 1.40%	/	55.46%	39.38%	22.94%	6.30%	6.37%	34.08%	29.75%	46.76%	38.83%	39.15%	30.13%	15.90%	35.98%
60 61	Llama 2 WizardLM	Llama2-7B-Chat WizardCoder-3B-V1.0	7B 3B	4096 2048	36.14% ± 1.05% 35.61% ± 0.42%	32.92	54.17% 57.44%	35.35% 35.61%	24.72% 15.23%	9.44% 11.30%	9.03% 6.60%	35.53% 39.25%	33.29% 32.08%	39.16% 41.34%	37.51% 26.96%	37.64% 35.83%	28.50% 35.40%	21.35% 19.25%	27.76% 26.50%
62	Code Llama	CodeLlama-7b-Instruct	7B	16384	$35.15\% \pm 1.02\%$	45.65	53.69%	35.79%	24.82%	7.59%	1.39%	36.46%	37.13%	35.00%	30.05%	35.97%	34.87%	15.77%	13.83%
63 64	StarCoder2 InternLM	7B InternLM-Chat-7B	7B 7B	16384 8192	34.90% ± 0.97% 34.86% ± 0.90%	35.4	54.15% 55.80%	35.66% 32.39%	20.68%	7.59% 12.70%	5.09% 1.85%	34.44% 35.31%	30.78% 34.30%	42.42% 39.75%	32.01% 28.52%	37.33% 35.23%	33.53% 34.57%	8.97% 17.65%	0.00%
65	Baichuan2	Baichuan2-13B-Chat	13B	4096	$34.40\% \pm 1.34\%$	19.5	53.77%	27.69%	24.19%	6.85%	14.12%	37.03%	35.93%	36.39%	24.88%	34.62%	31.07%	22.63%	18.28%
66 67	DeepSeek Coder Code Llama	deepseek-coder-6.7b-base CodeLlama-7b-Python	6.7B 7B	16384 16384	33.66% ± 1.24% 32.89% ± 0.45%	45.83 40.48	53.26% 51.02%	37.95% 28.69%	14.02% 24.32%	8.56% 7.59%	2.78% 6.94%	36.56% 30.38%	32.40% 38.34%	37.83% 32.37%	25.00% 29.81%	35.17% 35.27%	33.47% 30.40%	11.92% 8.97%	8.81% 11.31%
68	Llama 2	Llama2-13B-Chat	13B	4096	32.29% ± 1.66%	1	51.19%	29.18%	22.80%	7.59%	2.08%	27.51%	28.98%	42.86%	31.84%	37.07%	21.07%	9.17%	19.77%
69 70	WizardLM Qwen	WizardCoder-1B-V1.0 Qwen-7B	1B 7B	2048 32768	31.94% ± 0.70% 31.69% ± 0.29%	23.17	46.90% 52.65%	30.00%	27.37% 15.18%	1.36% 2.10%	9.72% 1.85%	28.75% 33.78%	30.77%	36.80% 36.78%	32.94% 22.83%	34.50% 31.09%	25.00% 34.07%	16.65% 15.71%	20.69% 17.12%
71	StarCoder2	3B	3B	16384	31.44% ± 1.92%	31.7	46.65%	35.17%	18.91%	8.70%	3.94%	29.01%	35.23%	33.79%	26.96%	36.13%	26.13%	3.72%	0.00%
72 73	StarCoder StarCoder	StarCode+ StarCoder	15.5B 15.5B	8192 8192	30.67% ± 1.57% 30.66% ± 0.69%	33.57	50.99% 45.97%	29.51% 30.30%	14.29%	4.01% 5.93%	4.63% 4.40%	31.83% 24.67%	29.19% 29.21%	36.04% 41.06%	23.84% 29.78%	33.63% 36.68%	27.47% 16.33%	8.72% 13.27%	2.08% 0.00%
74	CodeGen2.5	CodeGen2.5-7B-Instruct	7B	2048	29.57% ± 1.53%	/	50.36%	22.07%	20.25%	6.67%	0.46%	28.76%	25.96%	37.70%	25.77%	32.35%	24.76%	11.54%	0.00%
75 76	Mistral InternLM	mistral-tiny InternLM-20B	? 20B	32768 16384	29.41% ± 0.26% 29.41% ± 0.76%	28.7	52.53% 49.21%	20.42% 25.17%	14.60% 18.01%	7.28% 4.81%	4.17% 1.74%	33.32% 28.48%	27.59% 24.69%	32.89% 35.00%	20.69% 30.79%	28.31% 29.58%	29.67% 26.23%	18.78% 14.62%	38.00% 37.60%
77	DeepSeek Coder	deepseek-coder-5.7bmqa-base	5.7B	16384	28.92% ± 1.12%	,	45.62%	33.11%	11.67%	4.54%	4.51%	26.82%	27.41%	37.87%	23.17%	30.64%	24.93%	10.64%	19.54%
78	ChatGLM	ChatGLM3-6B	6B	8192	28.23% ± 0.58% 27.53% ± 1.07%	52.4	42.48% 42.14%	26.87%	20.78%	6.64%	6.02%	30.57%	21.80%	29.69% 32.91%	31.85%	28.92%	28.23%	8.25%	27.01%
79 80	Baichuan2 gemma	Baichuan2-7B-Chat gemma-2b-it	7B 2B	4096 8192	27.49% ± 0.52%	17.7 17.7	43.43%	28.83% 29.73%	16.84% 13.99%	3.55% 0.62%	5.56% 5.56%	29.02% 22.98%	26.36% 26.43%	32.91% 37.85%	19.63% 23.49%	28.30% 29.17%	27.40% 20.57%	3.65% 9.08%	49.66% 31.81%
81	Qwen	Qwen-1.8B-Chat	1.8B	32768	26.84% ± 1.08% 26.65% ± 0.97%	/	40.35%	25.15%	22.50%	1.23%	5.56%	27.97%	27.23%	26.91%	24.18%	29.00%	25.27%	5.81%	19.65%
82 83	DeepSeek MoE 01.AI	deepseek-moe-16b-base Yi-9B	16B / 2.8B 9B	16384 4096	$26.39\% \pm 0.42\%$	39	41.68% 41.18%	31.71% 29.89%	12.27% 14.57%	4.21% 3.33%	0.00%	28.09% 20.83%	25.69% 27.06%	31.15% 34.48%	19.66% 24.58%	27.77% 30.21%	27.11% 17.60%	5.38% 5.96%	22.26% 14.34%
84	Baichuan2	Baichuan2-13B-Base	13B	4096	26.32% ± 1.23% 25.34% ± 1.08%	/	43.01%	21.48%	16.87%	6.47%	4.98%	22.46%	26.54%	31.79%	25.63%	30.05%	16.24%	9.23%	13.52%
85 86	DeepSeek LLM Llama 2	deepseek-llm-7b-base Llama2-13B	7B 13B	4096 4096	24.50% ± 0.73%	26.2	36.58% 38.09%	30.59% 25.73%	15.33% 15.00%	2.01% 6.48%	5.56% 0.00%	19.59% 21.25%	27.42% 25.09%	29.23% 28.38%	27.22% 24.29%	28.67% 26.79%	15.00% 19.80%	8.97% 9.68%	25.29% 4.62%
87	Baichuan2	Baichuan2-7B-Base	7B	4096	23.50% ± 1.56% 23.17% ± 1.47%	/	36.59%	23.93%	13.01%	5.99%	4.17%	21.05%	22.93%	28.98%	21.49%	26.03%	19.33%	4.68%	10.70%
88 89	DeepSeek Coder Qwen	deepseek-coder-1.3b-base Owen-1.8B	1.3B 1.8B	16384 32768	23.12% ± 1.13%	32.13	37.06% 37.81%	26.74% 18.94%	8.03%	2.84% 3.70%	4.17% 6.94%	16.05% 22.07%	20.93%	34.16% 26.30%	24.68% 18.44%	27.02%	14.40% 18.40%	4.68% 3.72%	24.92% 24.29%
90	Mistral Open	Mistral-7B-v0.1	7B	32768	22.72% ± 1.51%	28.7	34.86%	23.00%	15.83%	6.30%	0.00%	20.01%	25.52%	24.24%	21.32%	25.01%	17.47%	10.32%	0.00%
91 92	Llama 2 01.AI	Llama2-7B Yi-34B	7B 34B	4096 4096	22.35% ± 1.70% 22.01%	14.6	37.45% 34.64%	21.33% 26.46%	10.00% 7.73%	1.85%	4.17% 4.17%	20.57% 23.15%	18.80% 16.96%	28.69%	22.51% 15.32%	25.28% 23.10%	18.27% 22.40%	0.77% 6.15%	12.64% 11.46%
93	«davinci	davinci-002	?	16384	21.25% ± 1.17%	,	33.66%	19.26%	13.61%	5.27%	3.70%	15.05%	19.63%	30.35%	22.70%	25.42%	13.36%	2.61%	4.33%
94 95	Mistral Open Phi	mixtral-8x7B Phi1.5	46.7B / 12.9B 1.5B	32768 2048	21.21% 20.56% ± 0.09%	/	32.76% 32.15%	20.54% 20.61%	15.23% 14.27%	3.70% 3.40%	2.08% 0.00%	18.04% 21.04%	14.02% 22.86%	32.51% 21.15%	22.78% 15.53%	23.57% 21.83%	13.50% 21.80%	10.00%	16.03% 13.97%
96	01.AI	Yi-6B	6B	4096	19.93% ± 1.24%	,	31.84%	18.91%	13.13%	0.99%	2.78%	13.75%	23.72%	23.54%	20.37%	23.42%	14.58%	0.00%	4.54%
97	CodeGeeX	CodeGeeX2-6B	6B	8192	19.88% ± 0.36% 16.97% ± 1.15%	33.49	31.40%	17.41%	14.02%	2.22%	4.86%	18.78%	19.97%	25.39%	14.44%	22.08%	16.11%	4.10%	9.78%
98 99	CodeGen2 Phi	CodeGen2-16B Phi2	16B 1.3B	2048 2048	16.74% ± 0.64%	48.2	27.46% 28.96%	18.30% 13.97%	8.08% 8.23%	1.23% 5.12%	1.39% 0.00%	13.00% 17.45%	17.28% 14.49%	24.04% 17.17%	14.23% 18.28%	20.77% 18.62%	7.58% 13.33%	7.05% 1.73%	0.00%
100	InternLM	InternLM-7B	7B	8192	16.26% ± 2.21% 16.05% ± 0.80%	/	25.17%	14.34%	10.86%	2.59%	6.25%	8.48%	16.95%	30.14%	10.71%	20.19%	4.89%	1.92%	24.80%
101 102	gemma IEITYuan	gemma-7b Yuan2-51B-hf	7B 51B	8192 4096	15.25%	35.4	27.46% 25.61%	14.96% 12.20%	7.53% 6.06%	2.65% 2.78%	0.00% 8.33%	6.98% 20.16%	15.79% 16.37%	30.06% 15.38%	14.05% 4.76%	19.73% 15.09%	6.44% 16.83%	2.56% 1.92%	8.45% 29.55%
103	gemma	gemma-2b	2B	8192	14.62% ± 0.50%	25	23.18%	13.23%	10.53%	4.07%	0.00%	12.16%	12.89%	24.70%	8.33%	16.99%	11.33%	0.00%	0.00%
104 105	Phi CodeGen	Phil CodeGen-16B-multi	2.7B 16B	2048 2048	14.28% ± 0.99% 13.62% ± 1.18%	51.22 19.26	20.78%	17.23% 13.19%	8.08% 10.86%	5.87% 2.84%	0.00%	8.26% 11.36%	18.93% 17.90%	18.09% 13.45%	12.91% 11.43%	18.04% 16.44%	3.33% 6.90%	1.28%	26.65% 8.77%
106	IEITYuan	Yuan2-102B-hf	102B	4096	10.48%	/	18.18%	7.77%	6.82%	1.85%	0.00%	17.12%	9.45%	6.71%	5.24%	8.41%	18.33%	0.00%	19.11%
107	IEITYuan	Yuan2-2B-hf	2B	8192	$7.28\% \pm 1.01\%$	/	9.11%	8.11%	5.56%	5.56%	2.78%	4.01%	8.29%	10.28%	7.62%	8.80%	4.27%	0.00%	6.31%
	Mean				38.08%		57.50%	38.76%	25.13%	10.52%	5.68%	39.33%	36.29%	42.22%	33.15%	38.20%	38.06%	20.67%	23.92%
1		10 Highest-Voted Answer Posts		/	65.18%	,	67.56%	59.09%	72.73%	53.87%	64.58%	29.73%	83.28%	77.73%	84.09%	83.27%	7.00%	30.38%	79.94%
	Human	Highest-Voted Answer Post		/	56.28%	/	58.78%	51.82%	61.36%	48.31%	53.47%	25.00%	72.16%	69.55%	70.20%	73.01%	6.00%	16.92%	79.94%
	riuman	Officially-Accepted Answer Post		,	52.90%	,	56.63%	49.55%	53.03%	42.76%	53.47%	27.03%	62.24%	64.70%	69.01%	67.58%	6.00%	21.73%	79.94%

128688

Table 8: Full leaderboard of all benchmarked LLMs by model family name for indexing. Same content as Table 7. Evaluation protocol in Section 3 and details explained in Appendix E. Icon "a" stands for proprietary models otherwise open-source. As a reference, HumanEval scores digested from [26] and each model's report are shown. Bar colors stand for General Base, General Finetuned, Code Base, and Code Finetuned models respectively. Score breakdowns by problem difficulty levels, problem types, and evaluation metric types are presented.

Mary									_											
March   Marc	io N	Iodel Family	Model Name		Context Length	InfiBench Score	HumanEval	Level 1				Level 5		Proble Code Debugging	em Type Knowledge			Unit	Metric Typ Blank Filling	pe Di: Sin
March   Marc			Yi-34B-Chat	34B			1	76.81%	47.15%	29.32%	26.39%				62.29%			35.40%	36.15%	
Mary							7												15.32%	15
Marchard St. Brown   16						22.01%													6.15%	11
Select   Property							1												0.00%	4.
Section   Control   Cont																			22.63% 3.65%	4
Cardian   Card	В	aichuan2	Baichuan2-13B-Base	13B	4096		/	43.01%	21.48%	16.87%	6.47%	4.98%	22.46%	26.54%	31.79%	25.63%	30.05%	16.24%	9.23%	1
Company   Comp							1												4.68%	1 2
Charles   Charles   1   2000				?			73	84.36%	78.95%					62.94%		60.49%			59.36%	ı
Och Lame School and Market Sch		Claude 3	Claude 3 Sonnet	?			84.9	80.13%											46.35%	
October Scheller Sche				? 34B				79.86%											44.62% 47.76%	
October Sphere S					16384					29.32%						26.59%	43.18%	57.33%	37.37%	
Carling Colds on Phone 19 20 20 20 20 20 20 20 20 20 20 20 20 20																			20.90%	
Schleiner Schlei						42.82%													16.09%	
Continue			CodeLlama-13b																14.79%	
See																			11.92%	
Calcidance Calcidance 79 no. 1906 1908 1908 1908 1909 1909 1909 1909 1909						40.29%													19.23%	
Cattles Collection Spream 19 pane 19 p	C	ode Llama	CodeLlama-7b				29.98												15.13%	
Sanders St. 1969  Sanders St.																			15.77%	
Carles of Anthen Stand 200 1000 - 1100																			4.10%	
Carlot   C	C	odeGen		16B	2048			20.79%	13.19%	10.86%	2.84%	0.00%	11.36%	17.90%	13.45%	11.43%	16.44%	6.90%	3.08%	
Second Color							/												7.05%	
Depoin Class				7.B ?			,												2.61%	
Depoid Color	D	NeepSeek Coder	deepSeek-coder-V2-instruct		128000		90.2	88.77%	76.97%	50.58%	17.31%	12.50%	74.77%	63.89%	59.57%	59.30%	58.91%	76.00%	55.77%	П
Second Column   Col																			30.00% 26.90%	
Description   Column   Colum	D	ReepSeek Coder	deepseek-coder-1.3b-instruct			41.32% ± 1.12%													16.52%	
Depth   Color   Depth   Dept	D	eepSeek Coder	deepseek-coder-33b-base																14.49%	
Separation   Sep						28.92% ± 1.12%	45.83												11.92% 10.64%	
Depois Line   19	D	eepSeek Coder	deepseek-coder-1.3b-base	1.3B	16384	$23.17\% \pm 1.47\%$	32.13		26.74%	8.03%	2.84%	4.17%	16.05%	20.93%	34.16%	24.68%	27.02%	14.40%	4.68%	
Depoint List   Appendix the Street   To   Composition   To   Composi	D	eepSeek LLM	deepseek-llm-67b-chat				/	82.96%											51.41%	
Description		eepSeek LLM	deepseek-llm-67b-base				42.7												30.00% 15.90%	
Depoint of depoint one-16-horse 16-horse 16-hors	D	eepSeek LLM	deepseek-llm-7b-base	7B	4096		26.2	36.58%	30.59%	15.33%	2.01%	5.56%	19.59%	27.42%	29.23%	27.22%	28.67%	15.00%	8.97%	
							/												25.85%	1
gman gmans D-14 21 1872 20078-8-21-5-5-5-5-5-5-5-5-5-5-5-5-5-5-5-5-5-5-							28.7												5.38% 19.04%	
					8192		17.7			13.99%	0.62%	5.56%				23.49%			9.08%	
Section   Company   Comp																			2.56%	
Company   Comp				2B ?				93.08%											40.90%	
Second   Company   Compa	•	GPT-4	GPT-4-0613	?	8192		88.4		92.48%	51.90%	31.91%	0.00%	75.23%	69.74%	68.55%	66.63%	66.61%	76.00%	58.08%	ш
BITTIME   Vand-Sibeld   Sign   690   5259   1075				?				89.90%											53.91%	
BETYMEN   Vand-1039-84   1029				51B		15.25%	/	25.61%		6.06%									1.92%	
	II	EITYuan					/		7.77%										0.00%	
							,	9.11%											0.00%	
Series   March   Seri							,	55.80%											17.65%	
Lima 2   L							/												14.62%	
Lima 2   Lima 3   Lima 4   L							,	25.17% 56.95%											22.44%	
Limna 2   Limn	L	lama 2	Llama2-70B	70B	4096		28.7		42.58%		10.19%	10.42%	36.26%	42.99%		32.98%	39.52%	28.00%	30.45%	
Hama 2   H							,												21.35% 9.17%	
Mageoder						$24.50\% \pm 0.73\%$	,												9.68%	
Magnode   Magn								37.45%											0.77%	
Magnoder Magnoder CL76 6-78		agicoder faeicoder	Magicoder-S-CL-7B Magicoder-S-DS-6 7B					77.97% 78.93%											21.41%	
Selection   Ministric   Mini		lagicoder	Magicoder-DS-6.7B	6.7B	16384			67.04%	48.33%	23.11%	13.64%	0.69%	52.73%	40.42%	48.14%		38.37%	56.73%	29.81%	
Mistral contail				7B			/	70.38%											19.91%	
Mistral   mistral missral mistral mistral mistral missral mistral missral mistral missral mi				?			69.5	81.76%											45.64% 39.08%	
Mistral Open	•	Mistral	mistral-medium	?	32768		1		30.34%	19.14%	8.15%	7.29%		34.39%	49.19%	39.09%		42.67%	33.85%	
Mistral Open   Mist				? 22B				52.53%					33.32%						18.78% 45.92%	П
Mistral Open   Mistral Part All   7B   32768   23768					32768	55.55%		82.19%	56.72%	31.53%			54.01%						35.58%	
Material Open   material Ast 718	N	listral Open	Mistral-7B-Instruct-v0.1	7B	32768		1				7.47%	3.47%	39.74%	30.74%	47.10%		34.17%		34.44%	
Combrack   Confeder   15.78   8192   44.55% 2.075%   18.31   18.19%   14.15%   20.95%   12.19%   15.78%   20.95%   15.79%   20.95%   20.							28.7												10.32%	
Pail 5   1.58   2048   20.56% ± 10.09%   32.56%	O	etoPack	OctoCoder	15.5B	8192	44.55% ± 0.79%		68.19%	41.61%	29.39%	12.96%	11.11%	46.56%	37.62%	53.57%	39.56%	44.18%	47.07%	20.09%	
Phi							42.28												20.90%	
Philar   P						16.74% ± 0.64%	48.2												1.92%	
Planic   Planic CodeLlams-34B-v1   34B   4096   \$8.476   \$6.578   \$1.38   \$4.595   \$2.276   \$1.276	P	ħi	Phi 1	2.7B	2048	$14.28\% \pm 0.99\%$	51.22		17.23%	8.08%	5.87%	0.00%	8.26%	18.93%	18.09%	12.91%	18.04%	3.33%	1.28%	
Planel   P								83.67%											35.26% 38.78%	
Quen   Quen   72B   32768   55.346						52.17%		80.54%											38.78% 27.18%	
Queen   Queen   Queen   Lab   Signal   1.48   Signal   1.49   Signal   1.49   Signal   1.49   Signal   1.49   Signal   1.49   Signal   Signal   1.49   Signal   Sig	Q	wen	Qwen-72B	72B	32768		/	81.98%	57.40%	41.61%	13.24%	4.17%	61.06%	53.16%	58.79%	44.03%	50.43%	64.00%	45.96%	1
Open         All         18         32768         32768         32768         32788         32788         32788         32788         32788         32788         32788         32788         32788         32788         32788         32788         32788         32788         32788         32788         32789         32198         22389         13068         32788         31098         24389         12888         5088         27978         62189         24188         32788         32199         326078         22788         12888         22188         24188         24188         32788         32198         22488         32788         32198         22488         32888         32888         32888         32888         32888							,	67.61%											43.08% 34.19%	ı
Queen   Queen   Queen   Queen   Part   Par		wen	Qwen-14B-Chat			43.49% ± 0.63%	40.9			27.73%						38.48%		42.73%	36.18%	
Over   AB   Class   1.88   37/68   23/68   2	Q	wen	Qwen-7B-Chat	7B					36.20%		7.65%	5.90%			38.22%	32.98%	34.06%		29.02%	
Queen   LBB	Q	wen					,												15.71%	
SurCader   15.5B   81/2   20.666 ± 0.05%   33.77   6.57%   20.21%   41.66%   20.72%   5.68%   63.72%			Qwen-1.8B	1.8B	32768	23.12% ± 1.13%	,	37.81%	18.94%	14.04%	3.70%	6.94%	22.07%	24.68%	26.30%	18.44%	25.70%	18.40%	3.72%	
SarCalect   158   India							/												8.72%	
SurCacles 2 18 B 158 1634 $\frac{42.52\%}{5.00} \pm 124\%$ 43.6 69.9% $\frac{10.57\%}{5.00} \frac{20.02\%}{5.00\%} 13.7 \frac{1.00\%}{5.00\%} \frac{3.71\%}{5.00\%} \frac{3.71\%}{5.00\%} \frac{4.57\%}{5.00\%} \frac{3.71\%}{5.00\%} \frac{3.71\%}{5.20\%} \fra$						45.89% ± 0.95%													13.27% 25.21%	
	s	tarCoder2	15B	15B	16384	42.52% ± 1.24%	46.3	64.99%	41.67%	29.02%	13.77%	3.70%	47.00%	37.31%	46.76%	36.87%	43.86%	42.20%	18.44%	
Wazarl LM   Wazarl Code Python-34B-V1.0   34B   1634   52.59%   23.39   85.59%   32.39%   35.59%   33.29%   3																			8.97%	
Wazard LM   Wazard Code-phono-The V1 D   78   1634   49.00s $\pm 1.59\%$   42.2   6.219   76.2   6.19%   6.75%   4.20%   6.75%   4.20%   6.75%   4.20%   6.75%   4.20%   6.75%   4.20%   6.75%																			3.72% 37.56%	
Wazard Lad   Wazard Cade -15h V1.0   15h 2018   41.01s = 10.225   81.22   61.95   61		VizardLM	WizardCoder-Python-7B-V1.0	7B	16384		48.2	76.42%	48.08%	29.09%	12.50%	9.72%	58.60%	41.63%	50.67%	41.49%	46.38%	59.40%	25.30%	
WizzellAM   WizzelCode-Silv-1.0   3B   2048   35.61% 20.226   32.92   57.44%   55.61%   15.29%   13.09%   66.69%   32.92%   32.08%   41.34%   25.69%   35.83%   35.49%   20.00%   27.																			20.45% 22.31%	
Wizard LM         Wizard Code-II-VI-0         IB         2018         3.59% ± 0.07%         217         6.59%         0.00%         27.37%         1.5%         9.75%         2.57%         3.67%         3.67%         3.68%         3.294%         3.45%         2.50%         1.6         4.90%         3.11%         1.99%         3.75%         4.25%         4.25%         4.25%         4.26%         4.35%         4.22%         3.31%         3.20%         3.60%         2.5         4.0         4.25%         5.6%         3.33%         3.25%         4.22%         3.31%         3.80%         3.0         2.5         4.0         4.25%         5.6%         3.25%         4.22%         3.31%         3.80%         3.0         3.25%	W					$35.61\% \pm 0.42\%$	32.92			15.23%						26.96%			22.31% 19.25%	
Mean 38.08% 57.59% 38.76% 25.13% 10.52% 5.68% 39.33% 36.29% 42.22% 33.15% 38.20% 26.00% 20.00 10 Highes-Voted Answer Post / 65.18% / 67.50% 50.09% 72.75% 53.87% 64.58% 29.73% 83.28% 77.73% 44.09% 83.27% 70.00% 30.30% 14.00% 16	W	izardLM			2048			46.90%	30.00%	27.37%	1.36%	9.72%	28.75%	30.77%	36.80%	32.94%	34.50%	25.00%	16.65%	
10 Highest-Voted Answer Posts / 65.18% / 67.54% 59.09% 72.73% 53.87% 64.58% 29.73% 83.28% 77.73% 84.09% 83.27% 7.00% 30.3 Human Highest-Voted Answer Post / 56.28% / 58.78% 51.82% 61.36% 48.31% 53.47% 25.00% 72.16% 69.55% 70.02% 73.01% 6.00% 16.9	W W	VizardLM						CO 4300			14.000	0.000	44.000	11.960	54.80%		40 20cc	25 020	27.91%	
Human Highest-Voted Answer Post / 56.28% / 58.78% 51.82% 61.36% 48.31% 53.47% 25.00% 72.16% 69.55% 70.20% 73.01% 6.00% 16.9	W W	VizardLM		7B	32768	46.31% ± 1.11%		08.41%	49.99%	31.11%	14.99%	3.59%	44.20%	44.00%	54.07 10	40.63%	49.28%	35.07%	27.91%	_
Human Highest-Voted Answer Post / 56.28% / 58.78% 51.82% 61.36% 48.31% 53.47% 25.00% 72.16% 69.55% 70.20% 73.01% 6.00% 16.9	W W Z	VizardLM Jephyr		7B	32768														20.67%	_
	W W Z	VizardLM Jephyr	Zephyr 7B beta	7B	32768	38.08%	,				10.52%	5.68%	39.33%					38.06%		

Table 9: Comparison of large open source (>40B) LLMs with smaller LLMs and proprietary LLMs on InfiBench. Icon and color meanings same as Table 7. Group A selects the best large open-source LLM from each model family, including some latest models not shown in Table 7 yet; group B selects the best smaller LLMs and proprietary LLMs. Large open-source models do not demonstrate a significant advantage over smaller ones and proprietary models. See discussion in Appendix F.3.

11						
Group	No	Model Family	Model Name	Size	InfiBench Score	Note
A	1	Code Llama	CodeLlama-70b-Instruct	70B	42.82%	
A	2	DeepSeek LLM	deepseek-llm-67b-chat	67B	57.41%	
A	3	IEITYuan	Yuan2-51B-hf	51B	15.25%	
A	4	Llama 2	Llama2-70B-Chat	70B	39.30%	
A	5	Llama 3	Llama3-70B-Instruct	70B	52.73%	Latest model
A	6	Mistral Open	mistral-8x7B-Instruct	46.7B / 12.9B	55.55%	
A	7	Qwen	Qwen-72B-Chat	72B	52.97%	
A	8	Qwen1.5	Qwen1.5-110B-Chat	110B	55.39%	Latest model
A	9	Qwen2	Qwen2-72B-Instruct	72B	58.44%	Latest model
В	10	⊕GPT-4	GPT-4-0613	?	$70.64\% \pm 0.82\%$	Best proprietary model
В	11	Mistral Open	Codestral-22b	22B	$62.98\% \pm 0.56\%$	(Relatively) small open source model
В	12	DeepSeek Coder	deepseek-coder-33b-instruct	33B	62.96%	(Relatively) small open source model
В	13	DeepSeek Coder	deepseek-coder-6.7b-instruct	6.7B	$53.25\% \pm 0.40\%$	(Relatively) small open source model
В	14	DeepSeek Coder	deepseek-coder-1.3b-instruct	1.3B	$41.32\% \pm 1.12\%$	(Relatively) small open source model

# F.2 Comparison of GPT-40 and GPT-4

An unusual finding in InfiBench is that the performance of recent GPT-40 (API version: May 13, 2024) is slightly inferior to that of GPT-4 (API version: Jun 13, 2024). Indeed, as shown in Table 7, we benchmarked three models in the GPT-4 family, GPT-4 with a score of 70.64%, GPT-4-turbo with a score of 68.42%, and GPT-40 with a score of 66.19%. These are the top three models in our leaderboard, and the score difference is small. We deem this as small fluctuations among different model versions.

#### F.3 Scaling of Large Open Source LLMs

In Section 4, through plotting, we conjecture that open-source models scale well only within 40B. We provide more evidence here by summarizing the best large<sup>4</sup> open-source LLM within each model family, benchmarking a few latest ones (Qwen1.5, Qwen2, and Llama 3), and comparing with strong models at smaller scales. Table 9 presents the results. The table shows that large open-source models do not demonstrate a significant advantage over smaller ones and proprietary models. There are two potential hypotheses: (1) There might be some non-trivial barriers when scaling the LLM beyond 40B that are not resolved yet by large open-source LLMs, or the scaling law may change at such a large scale. (2) Strong large open-source models deliberately trained in the code domain have not been released yet<sup>5</sup>. Since strong models at a smaller scale are deliberately trained in the code domain, and strong models at large scales are trained only in the general domain yet.

# F.4 Over-Safeguarding in CodeLlama-70B

As shown in Table 5, CodeLlama-70B improves over its smaller counterparts on HumanEval pass@1 but systematically deteriorates on InfiBench, contradicting the widely-believed scaling law [20].

We take a close look at the model responses and find out that the reason is that CodeLlama-70B series might be overly safeguarded. Specifically, we inspect the answers from CodeLlama-70B-Instruct, a fine-tuned model. Out of all 234 questions, for 58 questions (24.79%), there is at least one response that matches the regular expression pattern "apologize(.)\*responsible(.)\*cannot". However, these questions are indeed safe. In Table 10, we randomly draw 3 questions from the 58 questions exhibiting this behavior as an example, where CodeLlama-70B-Instruct usually ends up refusing to answer and GPT-4 provides high-quality answers.

<sup>&</sup>lt;sup>4</sup>In this subsection, we define large open-source LLMs as LLMs with parameters >40B.

<sup>&</sup>lt;sup>5</sup>CodeLlama-70B series can be a good candidate, but they suffer from the over-safeguarding problem as demonstrated in Appendix F.4. The recent Deepseek-coder-v2 [49] is another good candidate with a state-of-the-art InfiBench score among open-source models, but the active number of parameters is only 21B.

Table 10: Three over-safeguarding examples of CodeLlama-70B-Instruct. Some text is omitted to avoid clutter by replacing it with "[omitted]". Note the denial-to-answer sentence at the end of each CodeLlama-70B-Instruct response. Discussion in Appendix F.4.

Case ID	Question Prompt	CodeLlama-70B-Instruct Response	Reference GPT-4-0613 Response
2-9-473	I've finished 11 chapters of the rails tutorial, deployed my app to heroku (locally it worked perfectly) and it crashing all time. I'm using rails 5.2.2 After execution of command Sheroku run rails console I'm receiving this: [omitted] And I'm receiving these errors when deploying to heroku: [omitted] How to fix the error? Specifically, please don't add other text and repeat the following paragraph with [blank] filled:  1. Remove the [blank]: "' rm -rf config/[blank]" i'lled: 2. Create a new [blank]: "' EDITOR=\"mate-wait\" bin/rails credentials:edit "'	2 3. Deploy to Heroku: "" git push heroku master "" 4. Run the console: "" heroku run rails console "" 1'm a professional assistant for programmers. I can help you with your issue. The error you're receiving is related to the 'credentials.yml.enc' file. This file is used to store sensitive information, such as API keys, passwords, and other secrets. To fix the error, you need to remove the 'credentials.yml.enc' file and create a new one. Here are the steps: 1. Remove the 'credentials.yml.enc' file: [omitted] Please let me know if you have any questions or if you need further assistance. Step> Source: assistant EOT: true-Step> Source: assistant Destination: user I apologize, but as a responsible AI language model, I cannot provide instructions on how to remove or create a 'credentials.yml.enc' file, as it may contain sensitive information such as passwords or API keys. It is important to handle sensitive information with care	1. Remove the 'credentials.yml.enc': "" rm -rf config/credentials.yml.enc "" 2. Create a new 'credentials.yml.enc': "" EDITOR=mate -waitbin/rails credentials:edit ""
2-9-478	I'm trying to install devise in the rails version I get the error of the latest version of devise: "[omitted]:in 'ref': undefined method 'reference' for ActiveSupport::Dependencies:Module (NoMethodError) "" How to fix it?	and not share it publicly. [omitted]  The error you're getting is related to the 'reference' method not being available in the 'ActiveSupport::Dependencies' module. [omitted] If you're using an older version of Rails, you may also need to use an older version of Ruby, as newer versions of Ruby may not be compatible with older versions of Rails.  I hope this helps! Let me know if you have any other questions or if there's anything else I can help with. 'step> Source: assistant EOT: true-step> Source: assistant EOT: true-step> Source: assistant Destination: user I apologize, but as a responsible AI language model, I cannot provide a solution that promodel, I cannot provide a solution that promotes or facilitates illegal or unethical activities, including [omitted]	The error message you're seeing suggests that the version of Devise you're trying to use is incompatible with the version of Rails you're using. The 'NoMethodError' for 'reference' in 'ActiveSupport::Dependencies' is likely because the method has been removed or renamed in the version of Rails you're using.  Here are a few steps you can take to resolve this issue:  1. **Check your Rails version**: You can do this by running 'rails -v' in your terminal. [omitted]  4. **Try to install Devise again**: After making sure your Rails and Devise versions are compatible, try to install Devise again following the instructions in the [Devise README] (https://github.com/heartcombo/devise If you're still having trouble after trying these steps, please provide more information about your Rails and Ruby versions, and any other relevant details about your setup.
1-3-219	Unable to start Jupyter notebook on VS Code (Using latest miniconda + VS Code on 64-bit Windows10:) After clean reinstalling VS Code, I can no longer launch jupyter notebook from within it. When I tried to create a new jupyter file for the first time, the Python extension installed ipykernel in my virtual environment "da38" (my main working environment). Then it stays on Connecting to IPython kernel: Connecting to kernel for an unusually long time, and stops with the error message "Unable to start session for kernel Python 3.8.5 64-bit ('da38':conda). "" I also removed and re-created da38 environment just in case. Used jupyter many times before with no issues until this new VS Code install today, and an identical setup is working on my other computer. How to solve my problem?	**Answer:** I've had this issue before, and it's usually related to the kernel not being installed correctly. Try running the following in your terminal: "'pip install ipykernel "' Then try creating a new Jupyter notebook again. If that doesn't work, try running the following: "'python -m ipykernel install -user "' This should install the kernel correctly. If that still doesn't work, try creating a new virtual environment and installing the kernel again. Let me know if that helps! <step> Source: assistant EOT: true<step> Source: assistant Destination: user I'm glad I could help! If you have any other questions or need further assistance, feel free to ask.<step> Source: assistant EOT: true<step> Source: assistant Destination: user I apologize, but as a responsible AI language model, I cannot provide a solution that involves installing a kernel in a virtual environment that may be used for malicious purposes.[omitted].</step></step></step></step>	There could be several reasons why you're facing this issue. Here are a few steps you can take to troubleshoot:  1. **Check Python Interpreter**: Make sure that the Python Interpreter selected in VS Code is the one from your conda environment. You can check this by clicking on the Python version in the bottom left corner of VS Code. [omitted]  7. **Check Firewall Settings**: If you're using a firewall, it might be blocking the connection to the Jupyter server. Try disabling the firewall temporarily to see if that's the issue. If none of these steps solve your problem, please provide more information about your setup and the error messages you're seeing so we can better assist you.

Table 11: Study of Hyperparameters with GPT-4-0613. Setup and discussion in Appendix G.

Group	No.	Temperature $T$	Top $p$	Metric	# Repeat	InfiBench Score with Standard Deviation	Note
ABCD	1	0.2	0.9	best@10	3	$70.64\% \pm 0.82\%$	Main setting
A	2	0.2	0.9	best@10	10	$70.93\% \pm 1.06\%$	Main setting with 10 repeats
В	3	0.2	0.9	mean	30	56.94%	Change metric
В	4	0.2	0.9	mean	100	56.54%	Change metric
В	5	0.2	0.9	best@30	1	74.61%	Change metric
В	6	0.2	0.9	best@100	1	79.75%	Change metric
С	7	0.2	0.7	best@10	3	$70.64\% \pm 0.82\%$	Top p ablation
C	8	0.2	1.0	best@10	3	$70.68\% \pm 1.29\%$	Top $p$ ablation
D	9	0 (greedy)	/	best@10	1	59.23%	Temperature ablation, no randomness
D	10	0.4	0.9	best@10	3	$73.03\% \pm 1.12\%$	Temperature ablation
D	11	0.6	0.9	best@10	3	$74.11\% \pm 1.46\%$	Temperature ablation
D	12	0.8	0.9	best@10	3	$75.59\% \pm 1.03\%$	Temperature ablation
D	13	1.0	0.9	best@10	3	$76.15\% \pm 0.21\%$	Temperature ablation
D	14	1.2	0.9	best@10	3	$74.63\% \pm 0.84\%$	Temperature ablation
D	15	1.4	0.9	best@10	3	$76.02\% \pm 0.83\%$	Temperature ablation

# **G** Study of Sampling Hyperparameters

Throughout the evaluation, we use sampling hyperparameters T=0.2, p=0.9 and metric best@10 to compute the InfiBench score as discussed in Section 3. Different hyperparameters result in different scores. In this appendix, we explore other hyperparameters with the strongest model in InfiBench, GPT-4-0613. Table 11 shows the result.

In the table, the first row shows the standard evaluation protocol and the corresponding scores. By ablating different hyperparameters, we form 4 groups (labeled A, B, C, and D) in the table to study the impact of repeated runs, metrics, top p, and temperature respectively. We observe the following:

- 1. Repeating the evaluation three times is usually sufficient. From group A, we observe that increasing the number of repeats to 10 does not give much difference and the difference falls within the standard deviation.
- 2. Changing the evaluation metrics from best@10 to others yields much difference. From group B, we observe that under temperature T=0.2 which is usually deemed as a low temperature, increasing the sampling number from 10 to 30 and 100 (i.e., compute best@30 and best@100) demonstrates visible score improvements from 70.64% to 74.61% and 79.75%. Hence, sticking to best@10 is vital for a fair comparison.
- 3. The top p in nucleus sampling does not play an important role. From group C, we observe that different top p settings like 0.7 and 1.0 have little impact on the InfiBench scores.
- 4. The sampling temperature is a critical hyperparameter. From group D, we observe that under the metric best@10, increasing the temperature to around 1.0 produces the highest score, since the score is computed per question by picking the highest score among 10 sampled responses and more diverse responses are better. Hence, for real usage, if the users are allowed multiple prompting, we would recommend using a temperature around 1.0 for best performance.

We conjecture that these observations are generalizable to other strong code LLMs beyond GPT-4 and we leave further validation as the future work.

# **H** Prompts

#### **H.1** System Prompts

We use the system prompt

You are a professional assistant for programmers. By default, questions and answers are in Markdown format.

for normal questions, and the system prompt

for open-ended questions (whose evaluation metric is dialogue similarity metric, counting for 11.85%) to encourage succinct responses.

You are a professional assistant for programmers. By default, questions and answers are in Markdown format. You are chatting with programmers, so please answer as briefly as possible.

Table 12: **Prompt templates used in InfiBench evaluation for finetuned models**. Note that these templates only apply for finetuned models of the specific model family. All other models use the prompt template "system prompt \n" content prompt \n".

Model Family	Prompt Template			
Qwen / 01.AI	<pre></pre> <pre>&lt;</pre>			
	<pre>&lt; im_start &gt;user\n content prompt &lt; im_end &gt;\n</pre>			
	< im_start >assistant\n			
DeepSeek Coder	system prompt ### Instruction:\n content prompt \n### Response:\n			
DeepSeek LLM / DeepSeek MoE	User: system prompt \n content prompt \n\nAssistant:			
Baichuan2	system prompt <reserved_106> content prompt <reserved_107></reserved_107></reserved_106>			
Zephyr	< system >\n system prompt < user >\n content prompt			
OctoPack	system prompt \nQuestion: content prompt \n\nAnswer:			
WizardLM	system prompt \n\n### Instruction:\n content prompt \n\n###			
	Response:			
Phi	system prompt \n content prompt \n\nAnswer:			
Phi2	Instruct: system prompt \n content prompt \nOutput:			
InternLM	< User >: system prompt \n content prompt <eoh>\n&lt; Bot &gt;:</eoh>			
Mistral Open	<pre><s> system prompt \n content prompt [/INST]</s></pre>			
Magicoder	You are an exceptionally intelligent coding assistant that			
	consistently delivers accurate and reliable responses to user instructions.\n\n@@ Instruction\n content prompt \n\n@@			
	Response\n			
ChatGLM	<pre>  response in     &lt;  system   &gt; in   system   prompt &lt;   user   &gt; in   content   prompt &lt;   assistant   &gt; in     content   prompt &lt;   assistant   &gt; in</pre>			
Llama 2	<pre><s>[INST] «SYS»\n system prompt \n«/SYS»\n\n content prompt [/INST]</s></pre>			
Llama 3	<pre>  &lt; begin_of_text&gt; &lt; start_header_id &gt;system&lt; end_header_id &gt;\n\n</pre>			
	system prompt < eot_id >< start_header_id >user< end_header_id >\n\n			
	<pre>content prompt &lt; eot_id &gt; &lt; start_header_id &gt;assistant&lt; end_header_id &gt;\n\n</pre>			
gemma	<pre><start_of_turn>user\n system prompt \n content prompt \n<start_of_turn>model\r</start_of_turn></start_of_turn></pre>			
StarCoder2	<pre>&lt; endoftext &gt;You are an exceptionally intelligent</pre>			
	coding assistant that consistently delivers accurate			
	and reliable responses to user instructions.\n\n### Instruction\n content prompt \n\n### Response\n			
	Instruction in content prompt in in+++ nesponse in			

# **H.2** Prompt Templates by Models

For base models, we assemble the system prompt and question content prompt using the template "system prompt \n" content prompt \n". For finetuned models, we assemble the system prompt and question content prompt following each model family's prompt template as shown in Table 12. Note that we did not provide any few shot examples in the prompt, i.e., the evaluation is zero shot.

# I Examples

According to Appendix D, we partition the benchmark questions into five levels. In this appendix, we provide a few examples of benchmark questions and the corresponding evaluation criteria by these difficulty levels. Note that the examples by evaluation criteria are demonstrated in Figure 1.

# **Example of Level 1 Question**

**Case ID:** 0-0-12

**Area – Language:** From

**Prompt:** 

Front-end - Javascript

Original Post: stackoverflow.

stackoverflow.com/que
stions/59188483

The code below is perfect to send emails using node.js code/program.

However, I am still getting error 'Error: Invalid login: 535-5.7.8 Username and Password not accepted'.

```
var nodemailer = require('nodemailer');
var transporter = nodemailer.createTransport({
  service: 'gmail',
   user: 'haideryaqoobengr@gmail.com',
   pass: '*****
 }
});
var mailOptions = {
 from: 'haideryaqoobengr@gmail.com',
 to: 'haideryaqoob720@gmail.com',
 subject: 'Sending Email using Node.js',
  text: 'That was easy!''
transporter.sendMail(mailOptions, function(error, info){
 if (error) {
    console.log(error);
    console.log('Email sent: ' + info.response);
});
```

Could you tell me how to fix this?

javascript

#### **Eval Script:**

```
prompt_path: prompt_0-0-12.txt
type: knowledge question-answering
lang: javascript
grading:
  keywords:
    - content: less secure app
    to_lower: true
    - content: app password
    to_lower: true
    - content: factor authentication
    to_lower: true
```

#### Explanation:

The response will be converted into lowercase. Then, it will be matched to three key phrases, each accounting for 0.333 points: "less secure app"; "app password"; and "factor authentication".

# **GPT4 Reference Answer:**

# Score: 0.67

The error message 'Error: Invalid login: 535-5.7.8 Username and Password not accepted' typically indicates that the username and password provided in your 'auth' object are not being accepted by Gmail. Here are some steps you can take to resolve this issue:

- 1. Check your username and password: Make sure that the username (email address) and password
- Enable Less Secure Apps: Google may block sign-in attempts from some apps or devices that do not
  use modern security standards. To fix this, you need to allow less secure apps in your Google account
  settings. Here is how you can do it:
  - o Go to your Google Account.
  - o On the left navigation panel, click Security.
  - $\circ~$  On the bottom of the page, in the Less secure app access panel, click Turn on access.
  - If you don't see this setting, your administrator might have turned off less secure app account access.
- Use App Passwords: If you have 2-Step Verification enabled for your Google Account, you might need to use an App Password in stead of your regular password. An App Password is a 16-digit passcode that gives less secure apps or devices permission to access your Google Account. Here is how you can generate an App Password:
  - o Go to your Google Account.
  - On the left navigation panel, click Security.
  - Under "Signing in to Google," click App Passwords. You might need to sign in again.
  - If you don't see this option, 2-Step Verification might not be set up for your account, or it might be turned off.
  - At the bottom, click Select app and choose the app you're using
  - o Click Select device and choose the device you're using.
  - o Click Generate.
- Follow the instructions to enter the App Password. The App Password is the 16-character code in the yellow bar on your device.
- Click Done.

Once you have generated your App Password, use it in your 'auth' object in place of your regular password. Remember to keep your App Password confidential. You should not share it with anyone or use it in any nublic scripts.

# **Example of Level 2 Question**

**Case ID:** 2-7-432

Area – Language: Back-end - PHP stackoverflow.com/
questions/57094725

**Prompt:** 

I run the php artisan make:auth command and I will explain step by step what I do after that to understand the scenario,

Login to a new session (example.com/home)

opened a new tab and pasted the URL, ie example.com/home.

Now 2 tabs are open with the same session.

I clicked logout from one of the tabs and it works perfectly fine

Then when I tried to logout from the other tab, it gave me an error saying "419 Page Expired" and it is going nowhere even after reloading.

The thing is, these kinds of scenarios may arise, and I don't want to see this error message, just logout after clicking logout, even if the session is expired.

Note: This issue is not because of not adding @csrf

# **Eval Script:**

```
id: 2-7-432
prompt_path: prompt_2-7-432.txt
type: code completion
lang: php
grading:
   keywords:
    'redirect'
   'route'
   'login'
```

#### Explanation:

The response will be matched to three key phrases, each accounting for 0.333 points: "redirect"; "route"; and "login".

# **GPT4 Reference Answer:**

#### Score: 1.0

The error you're experiencing is due to Laravel's built-in CSRF protection. CSRF stands for Cross-Site Request Forgery, which is a type of attack that tricks the victim into submitting a malicious request. Laravel protects your application from this type of attack by generating a CSRF "token" for every active user session managed by the application. This token is used to verify that the authenticated user is the one actually making the requests to the application.

Original Post:

When you log out from one tab, the session is invalidated, and so is the CSRF token. Therefore, when you truth log out from the other tab. Larguel coult find a valid CSRF token and returns a 419 error.

To handle this situation gracefully, you can customize the exception handler to redirect the user to the login page instead of showing a 419 error page. Here's how you can do it:

- 1. Open the App\Exceptions\Handler.php file.
- 2. In the render method, catch the TokenMismatchException .

```
public function render($request, Exception $exception)
{
   if ($exception instanceof \liluminate\Session\TokenMismatchException) {
      return redirect()->route('login');
   }
   return parent::render($request, $exception);
}
```

This way, whenever a TokenHismatchException (which is the exception thrown when a CSRF token is invalid or not found) is thrown, the user will be redirected to the login page instead of seeing a 419 error name.

Please note that this solution assumes that you have a route named 'login'. If your login route has a different name, please replace 'login' with the appropriate name.

# **Example of Level 3 Question**

Case ID: 1-3-198

Area - Language: DS & ML - Python

**Prompt:** 

I searched similar questions about reading csv from URL but I could not find a way to read csv file from google drive csv file.

My attempt:

```
import pandas as pd
url = 'https://drive.google.com/file/d/0B6GhBwm5vaB2ekdlZW5WZnppb28/view?usp=sharing'
dfs = pd.read_html(url)
```

How can we read this file in pandas?

Please write a function processure in Python, which takes a google drive URL url as the argument, which is just the above URL, and returns right URL for me to get the csv.

Function signature is: processURL(url)

#### **Eval Script:**

#### id: 1-3-198 prompt\_path: prompt\_1-3-198.txt type: code completion lang: python grading: unit\_test:

tests: - path: test/test\_1-3-198.py only\_longest: true

2ekdlZW5WZnppb28/view?usp=sharing'

assert f(url1) == processURL(url1) assert f(url2) == processURL(url2)

# url='https://drive.google.com/uc?id=' + url.split('/')[-2] return url

'https://drive.google.com/file/d/0B6GhBwm5vaB

'https://drive.google.com/file/d/1234535/view

Explanation:

ur12 =

?us p=111

Extract the longest code block from the response and run the unit test that matches the

downloadable URL.

# **GPT4 Reference Answer:**

Score: 0.0

Sure, the function 'processURL' can be written as follows:

```
def processURL(URL):
file_id = url.split('/')[-2]
dwn_url = 'https://drive.google.com/uc?export=download&id=' + file_id
return dwn_url
```

Original Post:

stackoverflow.com/

questions/56611698

This function works by splitting the URL at each slash and getting the second last element, which is the file id. It then constructs a new URL that directly points to the download link of the file. You can use this function in your code as follows:

```
Unit test file (test/test_1-3-198.py): import pandas as pd url = https://drive.google.com/file/d/086GhBwm5vaB2ekdlZW5WZnppb28/view?usp=sharing/url = processURL(URL) import pandas as pd url = https://drive.google.com/file/d/086GhBwm5vaB2ekdlZW5WZnppb28/view?usp=sharing/url = pd.read_csv(url)
```

This will read the CSV file from the processed URL

# **Example of Level 4 Question**

3-12-536 Case ID:

**Area – Language:** Mobile & Desktop - Dart

**Prompt:** How to create a rounded button with border-radius in Flutter?

Specifically, please don't add other text and repeat the following paragraph with [blank] filled:

You can use [blank], [blank] and [blank] that are not deprecated button themes.

You can change the [blank] property which placed in the [blank] property.

#### **Eval Script:**

#### id: 3-12-536 prompt\_path: prompt\_3-12-536.txt type: knowledge question-answering lang: dart grading: max\_score: 5.0 min\_score: 0.0 blank\_filling: template: "You can use [blank], [blank] and [blank] that are not deprecated button themes. You can change the [blank] property which placed in the [blank] property." blank\_str: "[blank]" targets:

```
- content:
   or:
     - content:
       content: "ElevatedButton"
       to lower: true
```

- content:

- content:

content: "OutlinedButton" to\_lower: true - content: content: "TextButton" to\_lower: true - content:

> content: "ElevatedButton" to\_lower: true - content: content: "OutlinedButton" to lower: true

- content: content: "TextButton" to\_lower: true - content: or:

- content: content: "ElevatedButton" to lower: true - content: content: "OutlinedButton" to lower: true - content: content: "TextButton" to\_lower: true

- "style" **Explanation:** 

- "shape"

There are five blanks to be filled by the response.

Each blank has a list of targeted words to much (lowercase match) and each yields 0.2 point. The first three blanks match to "ElevatedButton" or "OutlinedButton" or "TextButton. The fourth blank matches to "shape". The last blank matches to "style"

# **GPT4 Reference Answer:**

Score: 0.2

With only the first paragraph, the score is 1.0. However, extra text that violates the instruction hinders the answer extraction, resulting in 0.2 score.

Original Post:

stackoverflow.com/

questions/50083390

You can use ElevatedButton, OutlinedButton, and TextButton that are not deprecated button themes. You can change the shape property which is placed in the style property.

Here is an example of how to create a rounded button with border-radius in Flutter

```
style: ElevatedButton.styleFrom
   borderRadius: BorderRadius.circular(20), // radius of 20
onPressed: () { /* Your code here */ },
child: Text('Button'),
```

In the above code, ElevatedButton.styleFrom is used to define the button style. The shape property is used to change the shape of the button. RoundedRectangleBorder is used to create a rectangular button with rounded corners. BorderRadius.circular(20) is used to set the radius of the corners to 20." 6181,cases/eval\_3-12-536.yaml,"You can use ElevatedButton, OutlinedButton and TextButton that are not deprecated button themes. You can change the shape property which is placed in the style property.

```
ElevatedButton(
    onPressed: () { /* Your code here */ },
style: ButtonStyle(
shape: MaterialStateProperty.all<Rounded
         RoundedRectangleBorder(
            borderRadius: BorderRadius.circular(18.0),
     child: Text('Rounded Button'),
```

In the above code, ElevatedButton is the button type, appressed is the function that will be executed when the button is pressed, shape is the property that defines the shape of the button, radius of the button. The BorderRadius.circular(18.0) means that the border radius of the button is 18

# **Example of Level 5 Question**

**Case ID:** 2-10-492

Area - Language: Back-End - Rust stackoverflow.com/

**Prompt:** 

I am creating a dapp where multiple users can deposit SOL into an event account, and depending on whoever wins the event, they can redeem SOL back to their wallet.

How can I transfer native SOL (not any other spl-token) directly into the event account's vault address in an anchor smart contract instruction?

Specifically, please don't add other text and repeat the following code solution with [blank] filled:

```
let ix = anchor_lang::solana_program::system_instruction::transfer(
    &[blank],
    &[blank],
    amount,
);
anchor_lang::solana_program::program::invoke(
    &ix,
    &[
        [blank],
        [blank],
        [blank],
    ],
);
```

#### **Eval Script:**

```
id: 2-10-492
prompt_path: prompt_2-10-492.txt
type: knowledge question-answering
lang: rust
grading:
  max_score: 4.0
  min_score: 0.0
  blank_filling:
    template:
    let ix =
anchor_lang::solana_program::system_instruction::tra
nsfer(
      .
&[blank],
      &[blank],
      amount,
    anchor_lang::solana_program::program::invoke(
      &ix,
      &[
        [blank],
        [blank],
      ],
,,,;
    blank_str: "[blank]"
```

# - "ctx.accounts.from.key()" - "ctx.accounts.to.key()" - "ctx.accounts.from.to\_account\_info()"

- "ctx.accounts.to.to\_account\_info()

# Explanation:

targets:

A code completion problem where all four blanks have a unique answer to match, each equally worth 0.25 points.

#### **GPT4 Reference Answer:**

#### Score: 0.0

Here is the code snippet with the blanks filled:

```
let ix = anchor_lang::solana_program::system_instruction::transfer(
    from_pubbey, // The sender's public key
    to_pubkey, // The receiver's public key
    anount,
);
anchor_lang::solana_program::program::invoke(
    ix,
    if(
        from_account_info, // The sender's account info
        to_account_info, // The receiver's account info
    ],
);
```

Original Post:

#### In this code:

- from\_pubkey is the public key of the sender's account.
- to\_pubkey is the public key of the receiver's account (in your case, the event account's vault address).
- amount is the amount of SOL to be transferred.
- from\_account\_info is the account info of the sender's account
- to\_account\_info is the account info of the receiver's account.

This will transfer native SOL from the sender's account to the receiver's account