

# CODEMMLU: A MULTI-TASK BENCHMARK FOR ASSESSING CODE UNDERSTANDING CAPABILITIES OF CODELLMS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Recent advancements in Code Large Language Models (CodeLLMs) have predominantly focused on open-ended code generation tasks, often neglecting the critical aspect of code understanding and comprehension. To bridge this gap, we present CodeMMLU, a comprehensive multiple-choice question-answer benchmark designed to evaluate the depth of software and code understanding in LLMs. CodeMMLU includes over 10,000 questions sourced from diverse domains, encompassing tasks such as code analysis, defect detection, and software engineering principles across multiple programming languages. Unlike traditional benchmarks, CodeMMLU assesses models' ability to reason about code rather than merely generate it, providing deeper insights into their grasp of complex software concepts and systems. Our extensive evaluation reveals that even state-of-the-art models face significant challenges with CodeMMLU, highlighting deficiencies in comprehension beyond code generation. By underscoring the crucial relationship between code understanding and effective generation, CodeMMLU serves as a vital resource for advancing AI-assisted software development, ultimately aiming to create more reliable and capable coding assistants.

## 1 INTRODUCTION

In recent years, Code Large Language Models (CodeLLMs) have witnessed significant advancements (Wang et al., 2021; 2023b; Feng et al., 2020; Allal et al., 2023; Li et al., 2023; Lozhkov et al., 2024b; Guo et al., 2024; Pinnaparaju et al., 2024; Zheng et al., 2024c; Roziere et al., 2023; Nijkamp et al., 2022; Luo et al., 2023; Xu et al., 2022; Bui et al., 2023), demonstrating impressive performance across various benchmarks (Iyer et al., 2018; Lu et al., 2021; Chen et al., 2021; Austin et al., 2021; Lai et al., 2023; Hendrycks et al., 2021; da Silva Maldonado et al., 2017; Palomba et al., 2015). These sophisticated models, trained on extensive datasets of programming languages and coding patterns, have demonstrated impressive capabilities in generating, summarizing, debugging, and refactoring code. Despite the promising progress, existing benchmarks are often outdated and prone to data leakage (Matton et al., 2024), resulting in less rigorous evaluations. Moreover, in practice, CodeLLMs reveal various shortcomings, including bias and hallucination (Rahman & Kundu, 2024; Liu et al., 2024a), that have not been fully investigated due to the limitations of current benchmarks' ability.

Coding-related benchmarks primarily emphasize on code generation or completion (Iyer et al., 2018; Lu et al., 2021; Chen et al., 2021; Austin et al., 2021; Lai et al., 2023; Hendrycks et al., 2021; Ding et al., 2023; Zhuo et al., 2024), thereby restricting their capacity to assess the model's capabilities in handling a variety of tasks, including free-form and multiple-choice question answering (MCQA). Moreover, since code generation requires reliable evaluation based on test cases and executability, these benchmarks are typically limited in both quantity and diversity across different domains. Consequently, the measurement could be potentially biased and lack generalizations. Several efforts (Liu & Wan, 2021; Li et al., 2024) have been introduced to improve evaluation capabilities through free-form question answering and to expand domain diversity. However, these benchmarks frequently involve trade-offs that necessitate the use of less rigorous evaluation metrics, such as match-based metrics. As a result, the evaluations tend to lack robustness and reliability.

To address the aforementioned shortcomings, we present CodeMMLU - the first comprehensive MCQA benchmark specifically designed for evaluating code and software knowledge. We follow the scenario introduced in the well-known MMLU dataset (Hendrycks et al., 2020) from the natural language understanding domain offering a robust and simple-to-evaluate approach. CodeMMLU is specifically developed to overcome the limitations of current benchmarks in software field, providing a reliable evaluation for LLMs with the following key focuses:

**Comprehensiveness:** To address the issue of bias in a limited proportion of evaluation data, we carefully curated CodeMMLU from a diverse range of high-quality sources. In summary, we obtain over 10,000 questions for evaluation.

**Diversity in task, domain, and language:** CodeMMLU is designed to provide a comprehensive assessment of LLM's code understanding capabilities. To this end, we have diversified the tasks within the dataset, covering a spectrum from general software knowledge QA to real-world challenges like code generation, defect detection, and code repair. Additionally, the dataset encompasses a wide range of topic domains, including database management, APIs and frameworks, and data structures and algorithms, across more than 10 programming languages. This approach allows for a more fine-grained assessment of LLM capabilities, enabling us to observe the strengths and weaknesses of each model.

The dataset enables us to assess LLMs capabilities in coding and software tasks from a different commonsense perspective, extending beyond code generation and completion. As a result, we identify previously unexplored biases issue in CodeLLMs, align to the problem observed in natural language MCQA using LLMs as shown in prior work of Zheng et al. (2023). Furthermore, our experiments involving a large number of LLMs and various prompting techniques, such as few-shot and chain-of-thought, demonstrate noticeable patterns related to the impacts of model size, model family, and inconsistencies across different prompting strategies. Specifically, the closed-source model GPT-4 consistently achieves the highest average performance. Meanwhile, the Meta-Llama family demonstrates the greatest accuracy among the closed-source model families. Besides, Scaling law related to model size is partially observed within models of the same family but is not evident across different families. This suggests a substantial influence of pretraining datasets, methodologies, and model architectures. Additionally, the performance and advancements in prompting techniques illustrate a lack of improvement, implying concerns regarding hallucinations in the reasoning ability of each model. In summary, our contributions are outlined as follows:

1. We present CodeMMLU, the first MCQA benchmark for software and coding-related knowledge. Addressing the need for diverse evaluation scenarios in code as in the NLP field, CodeMMLU enables the ability to evaluate LLMs' alignment with human inference with MCQA in the software knowledge domain.
2. CodeMMLU provides a comprehensive evaluation of LLM capabilities, ensuring a substantial number of samples and diversity across tasks, domains, and languages. Hence, it enables a more nuanced understanding of an LLM's strengths and weaknesses, facilitating the development of models that are better aligned with the complexities and demands of the software domain.
3. To mitigate the bias issue in LLMs on MCQA to ensure a robust and reliable evaluation, we propose a rigorous evaluation scenario that requires the correct answer to be identified across all permutations of the choices. This increases the difficulty of our dataset, necessitating the use of more robust and sophisticated models.
4. Our experiments offer valuable insights into LLM performance, highlighting the impact of factors such as model size, model family, and prompting techniques used. This provides important information to the community on utilizing LLMs for specific tasks and domains.

## 2 RELATED WORK

**Code Generation and Understanding evaluation** To evaluate and advance the capabilities of Large Language Models (LLMs) in code-related tasks, numerous benchmark datasets have been developed, covering a wide range of programming challenges. Algorithm-style benchmarks include HumanEval Chen et al. (2021) and MBPP Austin et al. (2021), along with their extended versions

108 HumanEval+, Multi-PL and MBPP+ Liu et al. (2024b), which focus on standard algorithms. More  
 109 challenging algorithmic tasks are represented by CodeContests Li et al. (2022) and LiveCodeBench  
 110 Jain et al. (2024), which are based on competitive programming problems. DS-1000 Lai et al. (2023)  
 111 challenges LLMs with specific data manipulation and analysis tasks, MathQA-Python Austin et al.  
 112 (2021) focuses on mathematical problem-solving in Python. Repository-level benchmarks such as  
 113 RepoBench Liu et al. (2023), RepoEval Zhang et al. (2023), and SWE-Bench Jimenez et al. (2023)  
 114 present more complex, real-world scenarios in software development. Additionally, comprehensive  
 115 benchmarks like XCodeEval Khan et al. (2023), CRUXEval Gu et al. (2024a), and CodeXGLUE Lu  
 116 et al. (2021) evaluate LLMs across multiple dimensions of software development.

117 **Programming Comprehension Modeling** The community have actively research for a model of the  
 118 programmer behavior or cognitive process of programmer since the early day of software development  
 119 Storey (2005); Shneiderman & Mayer (1979); Xia et al. (2018). A cognitive model’s target to  
 120 describe the task of the cognitive structures and cognitive processes, which are involved programming  
 121 knowledge, concepts, techniques when comprehension; and processes the programmer involved when  
 122 building a problem solution. Shneiderman & Mayer (1979) introduce the multi-leveled of cognitive  
 123 structures which consist of measuring the semantic and syntactic knowledge. The semantic contains  
 124 programming concepts and technique (e.g. dynamic programming, recursion or sorting and merging  
 125 method); where syntactic related to programming language grammatical (e.g. format of iteration,  
 126 condition, library function). In other hands, to measuring the programmer comprehension in term of  
 127 cognitive processes, Shneiderman & Mayer (1979) propose 5 core foundation programming tasks  
 128 to any programmer behavior model which be accountable of, namely composition, comprehension,  
 129 debugging, modification, and learning. This model architecture are design to answer 2 important  
 130 question in programmer comprehension task, (1) measuring what kind of knowledge is available to  
 131 programmer and (2) what kind of process does the programmer go through during the designing  
 132 solution phase.

133 Table 1: Comparison between common code understanding benchmarks for LLMs in term of coverage  
 134 five foundation tasks of programming comprehension model.  $\dagger$  and  $*$  denote the benchmark with the  
 135 free-flow generation and multiple-choice question answering format, respectively.

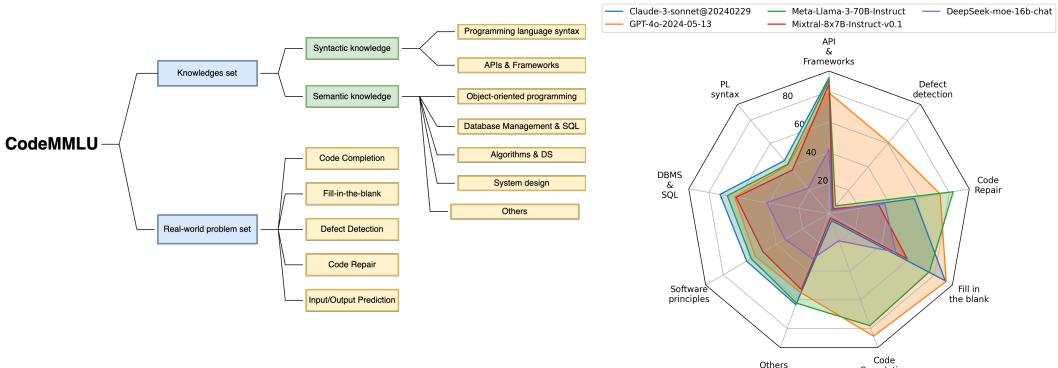
Benchmark	Programming Task					#Tasks	Data size
	Programming Knowledge	Code Composition	Code Comprehension	Code Modification	Code Debugging		
APPS $\dagger$ Hendrycks et al. (2021)		✓				1	5
MBPP $\dagger$ Austin et al. (2021)		✓				1	974
HumanEval $\dagger$ Chen et al. (2021)		✓				1	164
CRUXEval $\dagger$ Gu et al. (2024b)			✓			2	800
LiveCodeBench $\dagger$ Jain et al. (2024)		✓	✓		✓	4	-
CodeApex $\dagger*$ Fu et al. (2023)	✓	✓			✓	3	2.056
CodeMMLU $*$	✓	✓	✓	✓	✓	6	<b>19.912</b>

144 **Multi-task question answering benchmarks** Multiple Choice Questions (MCQs) have become  
 145 a significant tool for evaluating the capabilities of Large Language Models (LLMs) across various  
 146 domains. Recent trends in MCQ-based benchmarks are focused on testing advanced reasoning,  
 147 domain-specific knowledge, and robustness in LLMs, particularly in light of their expanding capabilities  
 148 Hendrycks et al. (2020); Lin et al. (2022); Zellers et al. (2019); Talmor et al. (2019). Since the  
 149 LLMs are evolving rapidly, MCQs method benefits are not only lie in enable scale evaluating, but  
 150 also provide a high-reliability result. Comparing with other open-ended judging method Chiang et al.  
 151 (2024) which heavily depend on LLM judge or human annotation effort, MCQs method enhance  
 152 the reliability by grounded the knowledge, problem context and defined possible available answers.  
 153 Despite the handy of MCQs, recent studies has shown LLMs display a significant sensitivity to  
 154 the order of answer options in MCQs Wang et al. (2023a); Robinson et al. (2023). Therefore raise  
 155 a request for applying appropriate debias methods on the MCQs benchmark to address the LLM  
 156 selection bias Zheng et al. (2024b); Pezeshkpour & Hruschka (2024)

### 3 MMCLU: DATA CREATION

160 The CodeMMLU benchmark is constructed to assess large language models’ (LLMs) comprehension  
 161 of programming tasks. We are inspired by programmer comprehension behavior models and integrate  
 multi-leveled cognitive structures and processes to measure the understandability of LLMs on software

problems (Shneiderman & Mayer, 1979). CodeMMLU is divided into two primary categories: (i) knowledge-based test sets containing syntactic and semantic tasks, and (ii) real-world programming problems. The overall CodeMMLU structure, as presented in figure 1, includes distinct approaches for data collection, filtering, and validation in both test sets.



**Figure 1: Overview of the CodeMMLU Structure.** The CodeMMLU benchmark is divided into two distinct subsets, each designed to evaluate different aspects of LLMs’ programming capability.

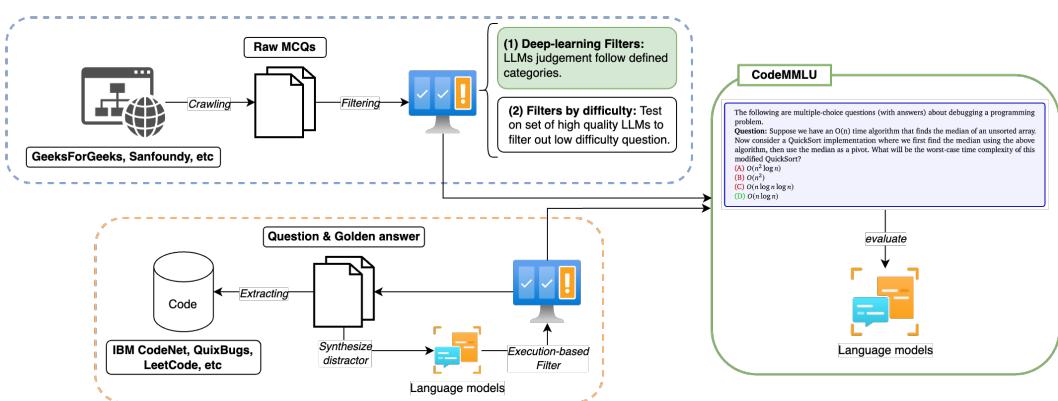
**Figure 2: Summary performance of LLMs on the CodeMMLU benchmark.** This radar chart presents the evaluation results (accuracy %) of different models across various CodeMMLU tasks.

### 3.1 KNOWLEDGE-BASED TASK CREATION

The knowledge-based subset covers wide range of topics, from high-level software principles concepts to low-level programming language grammars, targets to measuring the LLMs coding capability and comprehensibility of programming concepts. We collected programming-related MCQs from high-quality platforms, namely GeeksforGeeks, W3Schools, and Sanfoundry GeeksforGeeks (2024); W3Schools (2024); SanFoundry (2024).

- **Syntactic set.** Focused on programming language grammar and structural correctness, such as condition statement, format of iteration, common library usage, etc.
- **Semantic set.** Targeted more abstract programming concepts, such as algorithm, data structures, object-oriented principles, etc.

**Filtering Process.** We manually categorize questions into subject based on the topic of collected data. These processed questions are then applied a deep-learning-based filtering models to auto-



**Figure 3: Overview of CodeMMLU data creation pipeline.** The blue diagram describe the process of collecting raw multiple-choice questions (MCQs) from open source internet for a knowledge testset. Otherwise, the pipeline of real-world problem indicated in orange area.

216

217

Table 2: Summary of CodeMMLU Subject Categories and Task Distribution.

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

matically discard low-quality or irrelevant questions. For example, we removed duplicate or trivial question which did not challenge the code comprehension capabilities of LLMs (Appendix A.2.1). The final knowledge-based subset comprises of a manual and deep-learning-based filters to ensure each question met the desired quality standards, which must be clear, the unambiguous, and challenging to evaluate both semantic and syntactic understanding of LLMs. The filtering process returns a knowledge-based subset contains approximately 6,000 syntactic questions and over 3,000 question focus on semantic, cover 52 topics classified into 5 main subjects (Table 2).

### 3.2 REAL-WORLD TASK CREATION

Composition, comprehension, debugging, and modification are the core capabilities that any cognitive model of programmer behavior must address. We includes five distinct MCQ programming tasks, designed to assess the foundational abilities outlined in the cognitive process model of programmer comprehension, namely: Code completion; Code repair; Defect Detection; and Fill in the blank.

**Code Completion** covers the composition aspect in programming behavior model and evaluates a model’s ability by request to complete partially written code based on the provided requirement. We adapted the HumanEval Chen et al. (2021), originally designed for code generation, to and MCQ. From their 164 unique programming problems, we employed large language models (LLMs) to generate alternative, plausible but incorrect solutions, known as distractors. These alternatives were crafted to closely resemble plausible but incorrect solutions. All generated solutions, including both correct (migrate from HumanEval solutions) and incorrect options, were tested to ensure they are executable. Some of the incorrect solutions were designed to pass certain test cases but fail others. We target to add a layer of complexity, challenging the model to distinguish between correct and nearly-correct solutions based on semantic and syntactic understanding.

**Code Repair** is designed to evaluate model’s debugging capability by requiring to identify and fix errors in provided code snippets. We create Code Repair upon the QuixBugs Lin et al. (2017), which originally designed for debugging algorithmic programs. We use a "diff" operation on buggy and corrected version in QuixBugs in Python and Java to identify the specific changes needed to fix the bugs, which later served as the correct solution. To create plausible distractors for the MCQ format, we targeted components are known to be frequently involved in bugs, namely return statement, loop condition, if/else/switch expressions, and guided a LLMs to generate alternative fixes based on these components. These alternatives were intentionally designed to mislead by suggesting changes that seem plausible but do not (fully) resolve the bug. Each distractor was verified for incorrectness, and all options were executable to ensure that the model needed a deep understanding of the code to correctly identify and apply the necessary fix.

	Subject	Topic	Source	Testsize
Syntactic knowledge	API & Frameworks usage	Jquery, Django, Pandas, Numpy, Scipy, Azure, Git, AWS, svg, xml, Bootstrap, NodeJS, AngularJS, React, Vue, C, C#, C++, Java, Javascript, PHP, Python, R, Ruby, MatLab, HTML, CSS, TypeScript.	W3Schools, Geeks4Geeks, Sanfoundry	740
	Programming language syntax			6,220
Semantic knowledge	DBMS & SQL	DBMS, MySQL, PostgreSQL, SQL, Data structure & Algorithm, Object-oriented programming, Compiler design, Computer organization and Architecture, Software Development & Engineering, System Design.	W3Schools, Geeks4Geeks, Sanfoundry	393
	Software principles			3,246
	Others	Program accessibility, Computer networks, Computer science, Cybersecurity, Linux, Web technologies, AWS.		1,308
Real-world task		Code completion Fill in the blank Code repair Defect detection	HumanEval LeetCode QuixBugs IBM CodeNet	163 2,129 76 6,006

270 **Defect Detection** aims to evaluate a model’s ability to identify and understand defects within a  
 271 code snippet, focusing on both logical and syntactical errors. We measure the comprehension and  
 272 debugging capability of the LLMs by requiring it to predict the execution outcome of a given piece of  
 273 code. The task include two sub-tasks: detecting any defection/flaw in provided code, and comprehend  
 274 the output of certain test sample. The Defect Detection set derived from IBM CodeNet Puri et al.  
 275 (2021), a large-scale, high-quality benchmark designed for algorithmic coding tasks. IBM Project  
 276 CodeNet include over 14 million code samples submitted by users to solve a variety of algorithmic  
 277 challenges in multiple programming languages. From the vast pool of user-submitted solutions, we  
 278 focused on the Python and Java subsets and collected both accepted and buggy versions of code. To  
 279 ensure a high-quality and representative dataset, we filtered out duplicate submissions and retained a  
 280 diverse set of code samples. To create the Defect Detection task, we mixed the buggy versions with  
 281 their correct counterparts, ensuring a balanced representation of various types of errors. For each  
 282 code snippet, there will be a corresponding execution result (golden answer) and a random of three  
 283 others distracting option, which could be one of several possible outcomes, including: (i) Compile  
 284 Error, (ii) Time Limit Exceeded, (iii) Memory Limit Exceeded, (iv) Internal Error, (v) Runtime Error  
 285 and (vi) The code not contain any issue.

286 **Fill in the blank** asks the model to fill in missing parts of a code snippet, given documentation and  
 287 incomplete code snippet. We designed this task to assess a model’s ability to comprehend code, not  
 288 only because filling in the empty task, but models also be tested the understanding of both high-level  
 289 programming concepts and low-level grammatical structures, since problem requires models to  
 290 correctly identify and complete the incomplete parts of a program. We collected approximately 2,000  
 291 coding problem from LeetCode, covering solution from three wisely-use programming language  
 292 (Python, Java, C++). From each problem’s solution, we randomly selected key components to be  
 293 blanked out. These components were chosen for their importance in the logic and flow of the program,  
 294 for example loop conditions, expression statements, conditional statements.

295 To create plausible but incorrect options for the multiple-choice question (MCQ) format, we used  
 296 large language models (LLMs) to generate alternative solutions for the blanked-out components.  
 297 These distractors were designed to be contextually relevant but incorrect, adding complexity to the  
 298 task. After generating these distractors, we executed them to verify their incorrectness, ensuring that  
 299 they do not solve the problem as intended. To further enhance the difficulty of the task, we normalized  
 300 all variable names and function names within the code snippets. This normalization process involved  
 301 replacing original names with generic placeholders, such as ‘var1’, ‘var2’, ‘func1’, etc. This step  
 302 reduces the reliance on specific naming conventions and forces the model to focus purely on the logic  
 303 and structure of the code, rather than relying on context clues from variable or function names.

## 304 4 EXPERIMENTAL RESULTS

### 305 4.1 SETUP

306 **Model selection.** We evaluated 35 top-tier open-source models, covering a wide range of parameter  
 307 sizes and architectures. The models were selected from 7 different families, with parameters  
 308 ranging from 1 billion to over 70 billion. Each family included base, instructed, and chat versions:  
 309 MetaLlama3.1/8B/70B Dubey et al. (2024), MetaLlama3/8B/70B AI@Meta (2024), CodeLla-  
 310 MA/7B/13B/34B Rozière et al. (2024), DeepSeek-ai/6.7B/7B/33B Guo et al. (2024), MistralAI/8x7B  
 311 Jiang et al. (2024), Qwen2/7B qwe (2024), CodeQwen1.5/7B Bai et al. (2023), Yi/6B/9B/34B AI et al.  
 312 (2024), StarCoder2/7B/15B Lozhkov et al. (2024a). In addition to open-source models, we included 3  
 313 proprietary models from OpenAI and Claude to ensure comprehensive coverage of the state-of-the-art  
 314 in language modeling: GPT-3.5/GPT-4 OpenAI et al. (2024), Claude-3-opus/Claude-3.5-sonnet The.  
 315

316 **Answer extraction.** CodeMMLU leverages the MCQ format for scalability and ease of evaluation.  
 317 In order to maintain this advantage, we only apply simple regex methods to extract the selection  
 318 answer (i.e. extract by directly answering or contains the pattern “answer is A|B|C|D”). The model  
 319 response is required to be parsable; otherwise, it will be marked as unanswered.

320 **Prompting.** We employed various prompt strategies to test model performance across different  
 321 scenarios, namely: Zeroshot, Fewshot, CoT, and CoT Fewshot.

324 All experiments were conducted on a cluster of 8 A100 GPUs.  
 325

## 326 4.2 KEY INSIGHTS 327

328 **Table 3: Summary performance of LLM family on CodeMMLU.** The evaluation results (accuracy  
 329 %) of different language models across CodeMMLU task.

330     Family	331     Model name	332     Size (B)	333     Knowledge test	334     Syntactic	335     Semantic	336     Real-world tasks	337     CodeMMLU
<i>Closed-source models</i>							
333     Anthropic	Claude-3-sonnet@20240229	-	67.22	<b>66.08</b>		38.26	53.97
	GPT-4o-2024-05-13	-	60.41	57.82		77.18	<b>67</b>
334     OpenAI	GPT-3.5-turbo-0613	-	61.68	84.88		58.52	51.7
	<i>Open-source models</i>						
335     Meta Llama	CodeLlama-34b-Instruct-hf	34	56.81	46.93		23.55	38.73
	Meta-Llama-3-70B	70	63.38	86.02		63.1	48.98
	Meta-Llama-3-70B-Instruct	70	64.9	<b>87.59</b>		67.68	<b>62.45</b>
	Meta-Llama-3.1-70B	70	64.09	87.02		65.65	37.56
	Meta-Llama-3.1-70B-Instruct	70	64.42	87.45		<b>69.21</b>	60
341     Mistral	Mistral-7B-Instruct-v0.3	7	54.42	51.25		31.85	43.33
	Mistral-8x7B-Instruct-v0.1	46.7	61.17	85.02		61.83	42.96
	Codestral-22B-v0.1	22	60.34	82.17		58.52	47.6
343     Phi	Phi-3-medium-128k-instruct	14	58.54	54.56		37.89	48.03
	Phi-3-mini-128k-instruct	3.8	53.01	74.18		54.2	37.93
345     Qwen	Qwen2-57B-A14B-Instruct	57	61.34	57.48		30.48	46.34
	CodeQwen1.5-7B-Chat	7	49.66	67.62		46.06	49.82
347     Yi	Yi-1.5-34B-Chat	34	58.32	55.59		40.27	49.39
	Yi-1.5-9B-Chat	9	55.64	75.46		60.05	47.23
349     Deep Seek	DeepSeek-coder-7b-instruct-v1.5	7	56.67	47.9		28.46	41.21
	DeepSeek-coder-33b-instruct	33	53.65	74.89		52.16	36.6
	DeepSeek-moe-16b-chat	16.4	31.74	41.94		41.48	31.01
	DeepSeek-Coder-V2-Lite-Instruct	16	59.91	81.74		64.38	46.51
352     InternLM	InternLM2-5-20b-chat	20	57.85	55.51		30.44	44.89
353     StarCoder2	StarCoder2-15b-instruct-v0.1	15	56.58	49.07		42.79	47.94

355 The CodeMMLU benchmark revealed significant performance differences across models, as follow  
 356 in Tab. 3. The GPT-4o from OpenAI outperformed all models on CodeMMLU, reflecting its  
 357 quality across diverse tasks (Fig. 2). On the other hand, despite not being the latest model, the  
 358 Meta-Llama-3-70B instructed version achieves the highest score among open-source models from 8  
 359 families.

360 The Fig. 4 indicates the capability of measuring LLMs coding knowledge and skill in a wide range of  
 361 subjects of CodeMMLU. Our benchmark provides clear, distinct rankings that bring a higher hierarchy  
 362 of models than other code generation benchmarks. Interestingly, the result are not following the  
 363 scaling laws Kaplan et al. (2020), where models with larger parameter sizes outperform the smaller.  
 364 This shows the effect of data quality in the pretraining LLMs process to LLMs ability, where recently  
 365 released models are likely to achieve comparable performance with the large size model from previous  
 366 version. The importance of instruction tuning in improving the model’s performance in complex tasks  
 367 also indicated in CodeMMLU, while the prompt technique (i.e. Chain-of-thought Wei et al. (2023))  
 368 shows ineffective in boosting the comprehension of the models (Fig. 8. Models with instruction  
 369 tuning achieve substantial higher than its non-instruct counterpart, like DeepSeek-Coder-33b is  
 370 superior then its base model by approximately 29%.

371 **Correlation Between Software Knowledge and Real-World Performance** Our experiments  
 372 found a strong correlation between performance on knowledge-based tasks and real-world coding  
 373 tasks. The correlation between model ranks on the knowledge test set and real-world problems, by  
 374 **TODO: x** in Pearson’s correlation score, indicated a strong alignment between these two. Models  
 375 that demonstrated a deeper understanding of software principles consistently outperformed others on  
 376 real-world tasks.

377 **Selection bias in MCQs format** We experimented evaluating with multiple answer order permutations  
 378 (follow Zheng et al. (2024a)), the result displayed significant inconsistent behavior exhibited by

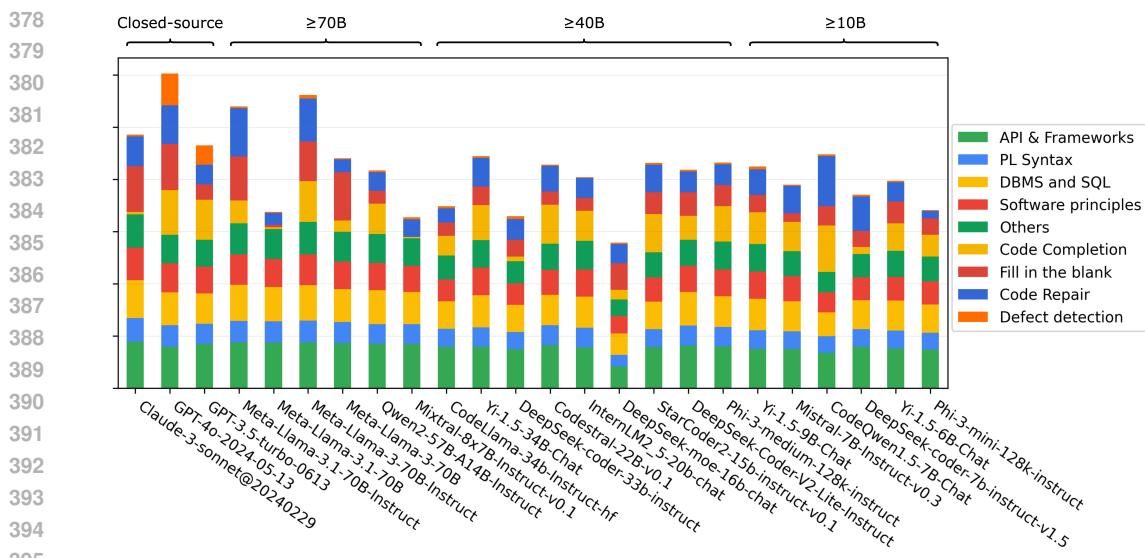


Figure 4: **CodeMMLU accuracy by task on LLMs.** While knowledge tasks are following the scaling law, real-world tasks offer more challenges to LLMs which indicate the performance of instruction tuning and data quality when evaluating on CodeMMLU.

LLMs when swapping golden answer positions. As presented in Tab. 7, the model’s performance changes dramatically in each answer order configuration, which is based on the correct answer’s position. The LLMs accuracy fluctuates between different permutations (i.e. DeepSeek-Coder-34B  $\Delta\sigma = 36.66$ ), demonstrating how sensitive it can be to the structure and ordering of answers (Fig. 6).

Figure 5: **Performance Comparison between HumanEval and MCQ Code Completion Tasks.** The performance fluctuation highlights the selection biases observed when the correct (golden) answer is moved to positions A, B, C, or D.

Models	HumanEval	Code Completion MCQ			
		A	B	C	D
CodeLlama-7B-Python	40.48	0.00 (-40.48)	90.24 (+49.76)	14.02 (-26.46)	0.61 (-39.87)
CodeLlama-7B-Instruct	45.65	3.66 (-41.99)	1.22 (-44.43)	93.90 (+48.25)	15.85 (-29.80)
CodeLlama-13B-Python	42.89	0.61 (-42.28)	54.88 (+11.99)	70.12 (+27.23)	12.20 (-30.69)
CodeLlama-13B-Instruct	50.6	2.44 (-48.16)	68.29 (+17.69)	72.56 (+21.96)	29.88 (-20.72)
CodeLlama-34B-Python	45.11	0.61 (-44.50)	77.44 (+32.33)	70.73 (+25.62)	49.39 (4.28)
CodeLlama-34B-Instruct	50.79	9.15 (-41.64)	84.76 (+33.97)	65.24 (+14.45)	46.34 (-4.45)
Deepseek-Coder-7B-base-v1.5	43.2	40.85 (-2.35)	74.39 (+31.19)	64.02 (+20.82)	39.02 (-4.18)
DeepSeek-Coder-33B-base	56.1	1.22 (-54.88)	82.32 (+26.22)	75.00 (+18.90)	56.10 (0.00)
Phind-CodeLLama-34B-v2	71.95	6.10 (-65.85)	90.85 (+18.90)	75.00 (+3.05)	46.34 (-25.61)
Mixtral-8x7B-Instruct-v0.1	40.2	22.56 (-17.64)	74.39 (+34.19)	71.95 (+31.75)	63.41 (+23.21)

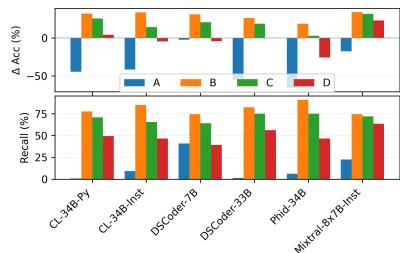


Figure 6: **Task-Specific Accuracy and Performance Fluctuations Across Answer Options** Models exhibit marked fluctuations in accuracy depending on the position of the correct answer in Code Completion in CodeMMLU. Revealing the bias and inconsistencies in related coding multiple-choice question (MCQ) task and how sensitive LLMs are to answer ordering.

**Disagreement between Traditional Benchmarks and CodeMMLU** A notable finding from our experiments is the discrepancy between performance on traditional benchmarks when comparing side-by-side the question HumanEval, and our code completion test set, which is HumanEval MCQ-based format. Many LLMs that perform well on HumanEval do not consistently translate their success to CodeMMLU. For instance, when comparing sample questions from HumanEval with their MCQ equivalents, the number of cases where models answered both correctly or both incorrectly was surprisingly low (Fig. 7).

This lack of alignment suggests that traditional benchmarks might overestimate a model’s understanding by focusing too narrowly on code generation. The MCQ format in CodeMMLU forces models to engage with more complex reasoning and contextual understanding, exposing weaknesses that remain hidden in generative tasks.

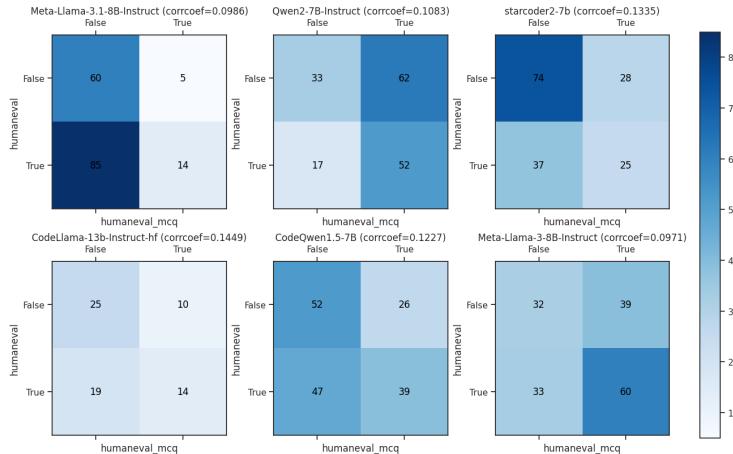


Figure 7: **CodeMMLU Task accuracy.**

## 5 CONCLUSIONS, LIMITATIONS, AND FUTURE WORK

In this work, we introduced CodeMMLU, a comprehensive and scalable benchmark designed to evaluate large language models’ (LLMs) capabilities across a wide range of software knowledge and real-world programming tasks. Our experiments highlighted the benchmark’s key advantages, including its cost-effectiveness, scalability, and extensive task coverage. The insights gained revealed a strong correlation between software knowledge and real-world task performance, demonstrating that models with deeper comprehension outperform those relying purely on probabilistic generation.

Additionally, CodeMMLU provides more accurate and detailed rankings of LLMs, particularly in open-source models, where significant reordering of performance was observed. The benchmark also revealed inconsistencies in model comprehension when compared to traditional evaluations like HumanEval, emphasizing the need for more robust benchmarks that go beyond simple code generation.

**Limitations.** While CodeMMLU offers a broad and diverse evaluation, there are some limitations. First, the MCQ format, though effective at testing comprehension, might not fully capture creative aspects of code generation or models’ ability to optimize code. Second, the current scope of languages and tasks could be expanded to include more specialized domains or additional programming languages to better assess models’ versatility.

**Future Work.** Looking forward, we plan to release CodeMMLU as an open-source benchmark for the research community. This release will include the full dataset, along with tools for automated evaluation, allowing for widespread adoption and further improvements. Future updates will focus on adding more complex tasks, refining the balance between real-world scenarios and theoretical knowledge, and incorporating user feedback to make the benchmark even more robust for next-generation LLMs.

486 REFERENCES  
487488 The claude 3 model family: Opus, sonnet, haiku. URL <https://api.semanticscholar.org/CorpusID:268232499>.  
489490 Qwen2 technical report. 2024.  
491492 01. AI, :, Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li,  
493 Jiangcheng Zhu, Jianqun Chen, Jing Chang, Kaidong Yu, Peng Liu, Qiang Liu, Shawn Yue, Senbin  
494 Yang, Shiming Yang, Tao Yu, Wen Xie, Wenhao Huang, Xiaohui Hu, Xiaoyi Ren, Xinyao Niu,  
495 Pengcheng Nie, Yuchi Xu, Yudong Liu, Yue Wang, Yuxuan Cai, Zhenyu Gu, Zhiyuan Liu, and  
496 Zonghong Dai. Yi: Open foundation models by 01.ai, 2024. URL <https://arxiv.org/abs/2403.04652>.  
497498 AI@Meta. Llama 3 model card. 2024. URL [https://github.com/meta-llama/llama3/blob/main/MODEL\\_CARD.md](https://github.com/meta-llama/llama3/blob/main/MODEL_CARD.md).  
499500 Loubna Ben Allal, Raymond Li, Denis Kocetkov, Chenghao Mou, Christopher Akiki, Carlos Munoz  
501 Ferrandis, Niklas Muennighoff, Mayank Mishra, Alex Gu, Manan Dey, et al. Santacoder: don't  
502 reach for the stars! *arXiv preprint arXiv:2301.03988*, 2023.  
503504 Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan,  
505 Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language  
506 models. *arXiv preprint arXiv:2108.07732*, 2021.  
507508 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge,  
509 Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu,  
510 Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan,  
511 Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin  
512 Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng  
513 Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou,  
514 Xiaohuan Zhou, and Tianhang Zhu. Qwen technical report. *arXiv preprint arXiv:2309.16609*,  
2023.515 Nghi DQ Bui, Hung Le, Yue Wang, Junnan Li, Akhilesh Deepak Gotmare, and Steven CH Hoi. Codetf:  
516 One-stop transformer library for state-of-the-art code llm. *arXiv preprint arXiv:2306.00029*, 2023.  
517518 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared  
519 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large  
520 language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.  
521522 Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li, Dacheng  
523 Li, Hao Zhang, Banghua Zhu, Michael Jordan, Joseph E. Gonzalez, and Ion Stoica. Chatbot arena:  
524 An open platform for evaluating llms by human preference, 2024. URL <https://arxiv.org/abs/2403.04132>.  
525526 Everton da Silva Maldonado, Emad Shihab, and Nikolaos Tsantalis. Using natural language pro-  
527 cessing to automatically detect self-admitted technical debt. *IEEE Transactions on Software  
Engineering*, 43(11):1044–1062, 2017.  
528529 Yangruibo Ding, Zijian Wang, Wasi Uddin Ahmad, Hantian Ding, Ming Tan, Nihal Jain, Murali Kr-  
530 ishna Ramanathan, Ramesh Nallapati, Parminder Bhatia, Dan Roth, and Bing Xiang. Crosscodee-  
531 val: A diverse and multilingual benchmark for cross-file code completion. In Alice Oh, Tristan  
532 Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), *Advances in  
533 Neural Information Processing Systems 36: Annual Conference on Neural Information Processing  
Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023.  
534535 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha  
536 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn,  
537 Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston  
538 Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron,  
539 Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris  
McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton

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

- 594 Meghan Keneally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov,  
 595 Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat,  
 596 Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White,  
 597 Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning  
 598 Dong, Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli,  
 599 Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux,  
 600 Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao,  
 601 Rachel Rodriguez, Rafi Ayub, Raghatham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li,  
 602 Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott,  
 603 Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru  
 604 Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng  
 605 Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang,  
 606 Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen,  
 607 Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin Cho, Sunny  
 608 Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best,  
 609 Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook  
 610 Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar,  
 611 Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li,  
 612 Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaofang Wang,  
 613 Xiaoqian Wu, Xiaolan Wang, Xide Xia, Xilun Wu, Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi,  
 614 Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian,  
 615 Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei  
 Zhao. The llama 3 herd of models, 2024. URL <https://arxiv.org/abs/2407.21783>.
- 616 Zhangyin Feng, Daya Guo, Duyu Tang, Nan Duan, Xiaocheng Feng, Ming Gong, Linjun Shou, Bing  
 617 Qin, Ting Liu, Dixin Jiang, et al. Codebert: A pre-trained model for programming and natural  
 618 languages. *arXiv preprint arXiv:2002.08155*, 2020.
- 619 Lingyue Fu, Huacan Chai, Shuang Luo, Kounianhua Du, Weiming Zhang, Longteng Fan, Jiayi Lei,  
 620 Renting Rui, Jianghao Lin, Yuchen Fang, Yifan Liu, Jingkuan Wang, Siyuan Qi, Kangning Zhang,  
 621 Weinan Zhang, and Yong Yu. Codeapex: A bilingual programming evaluation benchmark for large  
 622 language models. *ArXiv*, abs/2309.01940, 2023. URL <https://api.semanticscholar.org/CorpusID:261530384>.
- 625 GeeksforGeeks. GeeksforGeeks | Quiz Hub: Test Your Knowledge, 2024. URL <https://www.geeksforgeeks.org/quizzes/>.
- 626 Alex Gu, Baptiste Rozière, Hugh Leather, Armando Solar-Lezama, Gabriel Synnaeve, and Sida I  
 627 Wang. Cruxeval: A benchmark for code reasoning, understanding and execution. *arXiv preprint  
 628 arXiv:2401.03065*, 2024a.
- 631 Alex Gu, Baptiste Rozière, Hugh Leather, Armando Solar-Lezama, Gabriel Synnaeve, and Sida I.  
 632 Wang. Cruxeval: A benchmark for code reasoning, understanding and execution, 2024b. URL  
 633 <https://arxiv.org/abs/2401.03065>.
- 635 Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao  
 636 Bi, Y Wu, YK Li, et al. Deepseek-coder: When the large language model meets programming—the  
 637 rise of code intelligence. *arXiv preprint arXiv:2401.14196*, 2024.
- 638 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and  
 639 Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint  
 640 arXiv:2009.03300*, 2020.
- 642 Dan Hendrycks, Steven Basart, Saurav Kadavath, Mantas Mazeika, Akul Arora, Ethan Guo, Collin  
 643 Burns, Samir Puranik, Horace He, Dawn Song, and Jacob Steinhardt. Measuring coding challenge  
 644 competence with apps. *NeurIPS*, 2021.
- 646 Srinivasan Iyer, Ioannis Konstas, Alvin Cheung, and Luke Zettlemoyer. Mapping language to code in  
 647 programmatic context. In *Proceedings of the 2018 Conference on Empirical Methods in Natural  
 Language Processing*, pp. 1643–1652, 2018.

- 648 Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando  
 649 Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free  
 650 evaluation of large language models for code. *arXiv preprint arXiv:2403.07974*, 2024.
- 651
- 652 Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris  
 653 Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand,  
 654 Gianna Lengyel, Guillaume Bour, Guillaume Lample, Lélio Renard Lavaud, Lucile Saulnier, Marie-  
 655 Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le  
 656 Scao, Théophile Gervet, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed.  
 657 Mixtral of experts, 2024. URL <https://arxiv.org/abs/2401.04088>.
- 658
- 659 Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik  
 660 Narasimhan. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint*  
 661 *arXiv:2310.06770*, 2023.
- 662
- 663 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child,  
 664 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models,  
 665 2020. URL <https://arxiv.org/abs/2001.08361>.
- 666
- 667 Mohammad Abdullah Matin Khan, M Saiful Bari, Xuan Long Do, Weishi Wang, Md Rizwan Parvez,  
 668 and Shafiq Joty. xcodeeval: A large scale multilingual multitask benchmark for code understanding,  
 669 generation, translation and retrieval. *arXiv preprint arXiv:2303.03004*, 2023.
- 670
- 671 Yuhang Lai, Chengxi Li, Yiming Wang, Tianyi Zhang, Ruiqi Zhong, Luke Zettlemoyer, Wen-tau Yih,  
 672 Daniel Fried, Sida Wang, and Tao Yu. Ds-1000: A natural and reliable benchmark for data science  
 673 code generation. In *International Conference on Machine Learning*, pp. 18319–18345. PMLR,  
 674 2023.
- 675
- 676 Linyi Li, Shijie Geng, Zhenwen Li, Yibo He, Hao Yu, Ziyue Hua, Guanghan Ning, Siwei Wang,  
 677 Tao Xie, and Hongxia Yang. Inficode-eval: Systematically evaluating the question-answering  
 678 capabilities of code large language models. *arXiv preprint arXiv:2404.07940*, 2024.
- 679
- 680 Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou,  
 681 Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, et al. Starcoder: may the source be with  
 682 you! *arXiv preprint arXiv:2305.06161*, 2023.
- 683
- 684 Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittweiser, Rémi Leblond, Tom  
 685 Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, et al. Competition-level code generation  
 686 with alphacode. *Science*, 378(6624):1092–1097, 2022.
- 687
- 688 Derrick Lin, James Koppel, Angela Chen, and Armando Solar-Lezama. Quixbugs: a multi-lingual  
 689 program repair benchmark set based on the quixey challenge. *Proceedings Companion of the 2017*  
 690 *ACM SIGPLAN International Conference on Systems, Programming, Languages, and Applications:*  
 691 *Software for Humanity*, 2017. URL <https://api.semanticscholar.org/CorpusID:7158771>.
- 692
- 693 Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic human  
 694 falsehoods, 2022. URL <https://arxiv.org/abs/2109.07958>.
- 695
- 696 Chenxiao Liu and Xiaojun Wan. Codeqa: A question answering dataset for source code comprehen-  
 697 sion. *arXiv preprint arXiv:2109.08365*, 2021.
- 698
- 699 Fang Liu, Yang Liu, Lin Shi, Houkun Huang, Ruifeng Wang, Zhen Yang, and Li Zhang. Exploring  
 700 and evaluating hallucinations in llm-powered code generation. *arXiv preprint arXiv:2404.00971*,  
 701 2024a.
- 702
- 703 Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. Is your code generated by  
 704 chatgpt really correct? rigorous evaluation of large language models for code generation. *Advances*  
 705 *in Neural Information Processing Systems*, 36, 2024b.
- 706
- 707 Tiansyang Liu, Canwen Xu, and Julian McAuley. Repobench: Benchmarking repository-level code  
 708 auto-completion systems. *arXiv preprint arXiv:2306.03091*, 2023.

- 702 Anton Lozhkov, Raymond Li, Loubna Ben Allal, Federico Cassano, Joel Lamy-Poirier, Nouamane  
 703 Tazi, Ao Tang, Dmytro Pykhtar, Jiawei Liu, Yuxiang Wei, Tianyang Liu, Max Tian, Denis Kocetkov,  
 704 Arthur Zucker, Younes Belkada, Zijian Wang, Qian Liu, Dmitry Abulkhanov, Indraneil Paul,  
 705 Zhuang Li, Wen-Ding Li, Megan Risdal, Jia Li, Jian Zhu, Terry Yue Zhuo, Evgenii Zheltonozhskii,  
 706 Nii Osae Osae Dade, Wenhao Yu, Lucas Krauß, Naman Jain, Yixuan Su, Xuanli He, Manan  
 707 Dey, Edoardo Abati, Yekun Chai, Niklas Muennighoff, Xiangru Tang, Muhtasham Oblokulov,  
 708 Christopher Akiki, Marc Marone, Chenghao Mou, Mayank Mishra, Alex Gu, Binyuan Hui, Tri  
 709 Dao, Armel Zebaze, Olivier Dehaene, Nicolas Patry, Canwen Xu, Julian McAuley, Han Hu, Torsten  
 710 Scholak, Sébastien Paquet, Jennifer Robinson, Carolyn Jane Anderson, Nicolas Chapados, Mostafa  
 711 Patwary, Nima Tajbakhsh, Yacine Jernite, Carlos Muñoz Ferrandis, Lingming Zhang, Sean Hughes,  
 712 Thomas Wolf, Arjun Guha, Leandro von Werra, and Harm de Vries. Starcoder 2 and the stack v2:  
 713 The next generation, 2024a. URL <https://arxiv.org/abs/2402.19173>.
- 714 Anton Lozhkov, Raymond Li, Loubna Ben Allal, Federico Cassano, Joel Lamy-Poirier, Nouamane  
 715 Tazi, Ao Tang, Dmytro Pykhtar, Jiawei Liu, Yuxiang Wei, et al. Starcoder 2 and the stack v2: The  
 716 next generation. *arXiv preprint arXiv:2402.19173*, 2024b.
- 717 Shuai Lu, Daya Guo, Shuo Ren, Junjie Huang, Alexey Svyatkovskiy, Ambrosio Blanco, Colin  
 718 Clement, Dawn Drain, Dixin Jiang, Duyu Tang, et al. Codexglue: A machine learning benchmark  
 719 dataset for code understanding and generation. *arXiv preprint arXiv:2102.04664*, 2021.
- 720 Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xiubo Geng, Wenxiang Hu, Chongyang Tao, Jing  
 721 Ma, Qingwei Lin, and Dixin Jiang. Wizardcoder: Empowering code large language models with  
 722 evol-instruct. *arXiv preprint arXiv:2306.08568*, 2023.
- 723 Alexandre Matton, Tom Sherborne, Dennis Aumiller, Elena Tommasone, Milad Alizadeh, Jingyi He,  
 724 Raymond Ma, Maxime Voisin, Ellen Gilsean-McMahon, and Matthias Gallé. On leakage of code  
 725 generation evaluation datasets. *arXiv preprint arXiv:2407.07565*, 2024.
- 726 Erik Nijkamp, Bo Pang, Hiroaki Hayashi, Lifu Tu, Huan Wang, Yingbo Zhou, Silvio Savarese,  
 727 and Caiming Xiong. Codegen: An open large language model for code with multi-turn program  
 728 synthesis. *arXiv preprint arXiv:2203.13474*, 2022.
- 729 OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni  
 730 Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor  
 731 Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian,  
 732 Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny  
 733 Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks,  
 734 Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea  
 735 Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen,  
 736 Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung,  
 737 Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch,  
 738 Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty  
 739 Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte,  
 740 Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel  
 741 Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua  
 742 Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike  
 743 Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon  
 744 Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne  
 745 Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo  
 746 Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar,  
 747 Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik  
 748 Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich,  
 749 Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy  
 750 Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie  
 751 Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini,  
 752 Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne,  
 753 Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David  
 754 Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie  
 755 Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély,  
 Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo

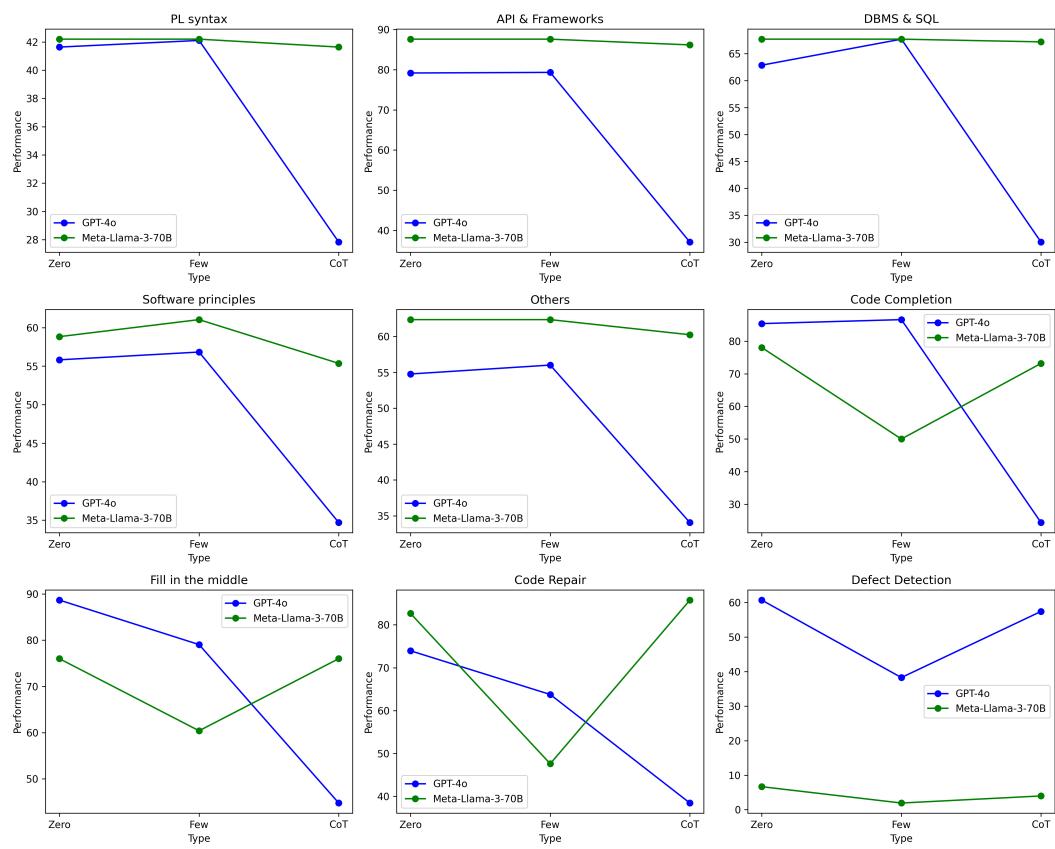
- 756 Noh, Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano,  
 757 Giambattista Parascandolo, Joel Parish, Emry Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng,  
 758 Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto,  
 759 Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power,  
 760 Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis  
 761 Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted  
 762 Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel  
 763 Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon  
 764 Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky,  
 765 Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie  
 766 Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng,  
 767 Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun  
 768 Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang,  
 769 Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian  
 770 Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren  
 771 Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming  
 772 Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao  
 773 Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. Gpt-4 technical report, 2024. URL  
<https://arxiv.org/abs/2303.08774>.
- 774 Fabio Palomba, Dario Di Nucci, Michele Tufano, Gabriele Bavota, Rocco Oliveto, Denys Poshyvanyk,  
 775 and Andrea De Lucia. Landfill: An open dataset of code smells with public evaluation. In *2015*  
 776 *IEEE/ACM 12th Working Conference on Mining Software Repositories*, pp. 482–485. IEEE, 2015.
- 777 Pouya Pezeshkpour and Estevam Hruschka. Large language models sensitivity to the order of options  
 778 in multiple-choice questions. pp. 2006–2017, 01 2024. doi: 10.18653/v1/2024.findings-naacl.130.
- 779 Nikhil Pinnaparaju, Reshith Adithyan, Duy Phung, Jonathan Tow, James Baicoianu, Ashish Datta,  
 780 Maksym Zhuravinskyi, Dakota Mahan, Marco Bellagente, Carlos Riquelme, et al. Stable code  
 781 technical report. *arXiv preprint arXiv:2404.01226*, 2024.
- 782 Ruchir Puri, David Kung, Geert Janssen, Wei Zhang, Giacomo Domeniconi, Vladimir Zolotov,  
 783 Julian T Dolby, Jie Chen, Mihir Choudhury, Lindsey Decker, Veronika Thost, Veronika  
 784 Thost, Luca Buratti, Saurabh Pujar, Shyam Ramji, Ulrich Finkler, Susan Malaika, and  
 785 Frederick Reiss. Codenet: A large-scale ai for code dataset for learning a diversity of  
 786 coding tasks. In J. Vanschoren and S. Yeung (eds.), *Proceedings of the Neural Information  
 787 Processing Systems Track on Datasets and Benchmarks*, volume 1, 2021. URL [https://datasets-benchmarks-proceedings.neurips.cc/paper\\_files/paper/2021/file/a5bfc9e07964f8dddeb95fc584cd965d-Paper-round2.pdf](https://datasets-benchmarks-proceedings.neurips.cc/paper_files/paper/2021/file/a5bfc9e07964f8dddeb95fc584cd965d-Paper-round2.pdf).
- 788 Mirza Masfiquur Rahman and Ashish Kundu. Code hallucination. *arXiv preprint arXiv:2407.04831*,  
 789 2024.
- 790 Joshua Robinson, Christopher Michael Rytting, and David Wingate. Leveraging large language  
 791 models for multiple choice question answering, 2023. URL <https://arxiv.org/abs/2210.12353>.
- 792 Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi  
 793 Adi, Jingyu Liu, Tal Remez, Jérémie Rapin, et al. Code llama: Open foundation models for code.  
 794 *arXiv preprint arXiv:2308.12950*, 2023.
- 795 Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi  
 796 Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, Jérémie Rapin, Artyom Kozhevnikov, Ivan Evtimov,  
 797 Joanna Bitton, Manish Bhatt, Cristian Canton Ferrer, Aaron Grattafiori, Wenhan Xiong, Alexandre  
 798 Défossez, Jade Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier, Thomas  
 799 Scialom, and Gabriel Synnaeve. Code llama: Open foundation models for code, 2024. URL  
<https://arxiv.org/abs/2308.12950>.
- 800 SanFoundry. SanFoundry, 2024. URL <https://www.sanfoundry.com/>.
- 801 Ben Shneiderman and Richard Mayer. Syntactic/semantic interactions in programmer behavior: A  
 802 model and experimental results. *International Journal of Parallel Programming*, 8:219–238, 06  
 803 1979. doi: 10.1007/BF00977789.

- 810 Zayne Sprague, Fangcong Yin, Juan Diego Rodriguez, Dongwei Jiang, Manya Wadhwa, Prasann  
 811 Singhal, Xinyu Zhao, Xi Ye, Kyle Mahowald, and Greg Durrett. To cot or not to cot? chain-of-  
 812 thought helps mainly on math and symbolic reasoning, 2024. URL <https://arxiv.org/abs/2409.12183>.
- 813
- 814 M.-A. Storey. Theories, methods and tools in program comprehension: past, present and future. In  
 815 *13th International Workshop on Program Comprehension (IWPC'05)*, pp. 181–191, 2005. doi:  
 816 10.1109/WPC.2005.38.
- 817
- 818 Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. Commonsenseqa: A question  
 819 answering challenge targeting commonsense knowledge, 2019. URL <https://arxiv.org/abs/1811.00937>.
- 820
- 821 W3Schools. W3Schools.com, 2024. URL <https://www.w3schools.com/quiztest/>.
- 822
- 823 Peiyi Wang, Lei Li, Liang Chen, Zefan Cai, Dawei Zhu, Binghuai Lin, Yunbo Cao, Qi Liu, Tianyu  
 824 Liu, and Zhifang Sui. Large language models are not fair evaluators, 2023a. URL <https://arxiv.org/abs/2305.17926>.
- 825
- 826 Yue Wang, Weishi Wang, Shafiq Joty, and Steven CH Hoi. Codet5: Identifier-aware unified pre-trained  
 827 encoder-decoder models for code understanding and generation. *arXiv preprint arXiv:2109.00859*,  
 828 2021.
- 829
- 830 Yue Wang, Hung Le, Akhilesh Deepak Gotmare, Nghi DQ Bui, Junnan Li, and Steven CH Hoi.  
 831 Codet5+: Open code large language models for code understanding and generation. *arXiv preprint  
 832 arXiv:2305.07922*, 2023b.
- 833
- 834 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le,  
 835 and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models, 2023.  
 URL <https://arxiv.org/abs/2201.11903>.
- 836
- 837 Xin Xia, Lingfeng Bao, David Lo, Zhenchang Xing, Ahmed E. Hassan, and Shaping Li. Measuring  
 838 program comprehension: A large-scale field study with professionals. In *2018 IEEE/ACM 40th  
 839 International Conference on Software Engineering (ICSE)*, pp. 584–584, 2018. doi: 10.1145/  
 3180155.3182538.
- 840
- 841 Frank F Xu, Uri Alon, Graham Neubig, and Vincent Josua Hellendoorn. A systematic evaluation of  
 842 large language models of code. In *Proceedings of the 6th ACM SIGPLAN International Symposium  
 843 on Machine Programming*, pp. 1–10, 2022.
- 844
- 845 Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine  
 846 really finish your sentence?, 2019. URL <https://arxiv.org/abs/1905.07830>.
- 847
- 848 Fengji Zhang, Bei Chen, Yue Zhang, Jacky Keung, Jin Liu, Daoguang Zan, Yi Mao, Jian-Guang Lou,  
 849 and Weizhu Chen. Repocoder: Repository-level code completion through iterative retrieval and  
 850 generation. *arXiv preprint arXiv:2303.12570*, 2023.
- 851
- 852 Chujie Zheng, Hao Zhou, Fandong Meng, Jie Zhou, and Minlie Huang. On large language models'  
 853 selection bias in multi-choice questions. *arXiv preprint arXiv:2309.03882*, 2023.
- 854
- 855 Chujie Zheng, Hao Zhou, Fandong Meng, Jie Zhou, and Minlie Huang. Large language models are  
 856 not robust multiple choice selectors. In *International Conference on Learning Representations*,  
 857 2024a. URL <https://openreview.net/forum?id=shr9PXz7T0>.
- 858
- 859 Chujie Zheng, Hao Zhou, Fandong Meng, Jie Zhou, and Minlie Huang. Large language models are not  
 860 robust multiple choice selectors, 2024b. URL <https://arxiv.org/abs/2309.03882>.
- 861
- 862 Tianyu Zheng, Ge Zhang, Tianhao Shen, Xueling Liu, Bill Yuchen Lin, Jie Fu, Wenhui Chen, and  
 863 Xiang Yue. Opencodeinterpreter: Integrating code generation with execution and refinement. *arXiv  
 864 preprint arXiv:2402.14658*, 2024c.
- 865
- 866 Terry Yue Zhuo, Minh Chien Vu, Jenny Chim, Han Hu, Wenhao Yu, Ratnadira Widayarsi, Imam  
 867 Nur Bani Yusuf, Haolan Zhan, Junda He, Indraneil Paul, et al. Bigcodebench: Benchmarking code  
 868 generation with diverse function calls and complex instructions. *arXiv preprint arXiv:2406.15877*,  
 869 2024.
- 870

## 864 A APPENDIX

### 865 A.1 PROMPT COMPARISON

866 We experimented on 48 LLM models from over 12 family with 3 different settings, namely: Zero-shot,  
 867 Few-shot and Chain-of-Thought (CoT). The overall trend can be demonstrated through GPT-4o and  
 868 MetaLlama 3 70B (Fig. 8), the SOTA models in closed-source and open-source families. The LLMs  
 869 show a significant decrease in performance in the Chain-of-Thought prompt setting compare to  
 870 few-shot and zero-shot settings. As concluded in Sprague et al. (2024) final result, CoT bring extra  
 871 "thinking" step into the original task which only effective on task involved math or logic while the  
 872 result on benchmark like MMLU are identical with and without CoT prompting. The decreasing that  
 873 happened mostly in CodeLLM knowledge test set align with Sprague et al. (2024) conclusion of the  
 874 inefficient of CoT for non-reasoning task.  
 875



905 **Figure 8: Comparison between GPT4o and Meta Llama-3 70B on various prompt settings.** We  
 906 experiment with zero-shot, 1-shot, and CoT prompt configuration, where the result indicates the  
 907 ineffectiveness of CoT in boosting the models' performance. Comparing to zeroshot config, 1-shot  
 908 prompt slightly increase the performance in knowledge tasks but falls shorter in real tasks.

### 910 A.2 MULTIPLE-CHOICE QUESTION ANSWERING FILTERING

#### 911 A.2.1 KNOWLEDGE SET: RAW DATA FILTER BASED LLM

912 **TODO:** This section describe how knowledge subset is filtering out

913 This includes:

- 914 • **Incomplete Data:** Instances that lack a complete problem statement, code snippet, or  
 915 solution are discarded.

- 918     • **Media Content:** Data containing images, videos, or other non-textual elements are excluded  
 919       to maintain consistency and ensure that the benchmark focuses solely on text-based coding  
 920       tasks.  
 921     • **Irrelevant or Redundant Content:** Any data that does not contribute to the assessment of  
 922       programming knowledge or coding ability is filtered out.  
 923

924     A.2.2 KNOWLEDGE SET: CONTROLLING DATA CHALLENGING  
 925

926     To ensuring the knowledge set meets the quality and challenging criteria, in order to provide a reliable  
 927       measure of LLM performance in the coding and programming domain. we pass the filtered subset  
 928       through a meticulous process to enhance the robust and comprehensive.

- 929     • **Common Errors Identification:** We scan the dataset for common errors or issues that  
 930       could lead to misleading assessments, such as ambiguous problem statements or incorrect  
 931       solutions.  
 932     • **Difficulty Level Assessment:** We evaluate the difficulty level of each instance, ensuring  
 933       that the tasks are neither too trivial nor exceedingly difficult. Instances deemed too easy  
 934       for LLMs are removed, as they do not contribute significantly to evaluating the model's  
 935       advanced problem-solving abilities.  
 936       This process is done after we benchmark through variety of LLMs to cherry pick most  
 937       common correct instances, and then spend manual effort to validate them.  
 938     • **Final Filtering:** Any data that fails to meet our stringent quality criteria is filtered out,  
 939       resulting in a final dataset that is both high-quality and well-suited for benchmarking.  
 940

941     By focusing on a wide range of software design principles subject, we target the benchmark to set a  
 942       new standard for evaluating LLMs in this critical area of AI development.  
 943

944  
 945  
 946  
 947  
 948  
 949  
 950  
 951  
 952  
 953  
 954  
 955  
 956  
 957  
 958  
 959  
 960  
 961  
 962  
 963  
 964  
 965  
 966  
 967  
 968  
 969  
 970  
 971

972 A.3 CODEMMLU ANALYSIS  
973

974 **TODO:** - Visualize CodeMMLU - Analysis CodeMMLU properties:  
975

- 976 • Question Length
- 977 • Difficulty
- 978 • Filter/Execution statistic

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

999

1000

1001

1002

1003

1004

1005

1006

1007

1008

1009

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1024

1025