

CRQBench: A Benchmark of Code Reasoning Questions

Anonymous EMNLP submission

Abstract

Large Language Models have demonstrated exceptional proficiency on coding tasks, but it is challenging to precisely evaluate their code reasoning ability. Existing benchmarks are insufficient as they are unrealistic and conflate semantic reasoning ability with performance on software engineering tasks. We introduce CRQBench, a benchmark of 100 C++ code reasoning questions and answers derived from contextualized code review comments. To curate CRQBench, we use an LLM assistant alongside human inspection, reducing manual effort. We conduct an evaluation of GPT-4 on CRQBench and find that it produces correct responses grounded in the given context for 65 of the 100 questions.

1 Introduction

Large Language Models (LLMs) have demonstrated effectiveness in coding tasks and appear to understand deep semantic properties of code (Chen et al., 2021; Chowdhery et al., 2022; Touvron et al., 2023). However, evaluations across various tasks (Jimenez et al., 2023; Zhong and Wang, 2023) show less promising results, suggesting that models may have a limited syntactic understanding of programs. To evaluate a model’s semantic reasoning ability in isolation, a benchmark specifically tailored for code reasoning question answering is needed.

The predominant benchmarks for evaluating LLMs trained on code are HumanEval (Chen et al., 2021) and MBPP (Austin et al., 2021). They measure a model’s ability to synthesize programs from docstrings. These text-to-code benchmarks are synthetic, handwritten, and involve generating a standalone function. Other, more realistic, benchmarks (Jimenez et al.,

2023; Zhong and Wang, 2023) are designed to evaluate code reasoning indirectly through a software engineering task, and as a result conflate the model’s ability to perform reasoning with the ability to perform the downstream task. In this work, we set out to curate a real-world, contextualized, benchmark for evaluating semantic reasoning ability in isolation.

Ideally, a benchmark for evaluating semantic reasoning ability should reflect real-world programming scenarios. Code review comments present an appealing target for this as they are non-synthetic and tied to a surrounding code context. Through a study of contextualized code review comments at a Corporation¹, we find that a subset embody semantically deep questions about code, but a majority are superficial (related to refactoring or style). Furthermore, we find that comments are rarely concise and unambiguous questions. Although code review comments provide a source of authentic semantic queries, it is non-trivial to extract clean questions.

We present CRQBench: a benchmark of real-world, contextualized, code reasoning questions. To reduce human curator effort, we propose a cooperative LLM and human-in-the-loop approach which leverages in-context learning (Brown et al., 2020) to filter and rephrase code reasoning questions from code review comments. We reproduce our Corporate results for open source release using Github pull request comments in the CodeReviewer dataset (Li et al., 2022).

In summary, our work presents a benchmark of 100 C++ (code reasoning questions, answer, code context) tuples derived from pull request comments in the CodeReviewer dataset. In addition, we present our curation technique as

¹Anonymized for double-blind review

a re-usable methodology and evaluate its effectiveness in reducing manual effort in benchmark curation. Lastly, we evaluate GPT-4 (OpenAI et al., 2024) on CRQBench and find that it produces correct responses grounded in the given context for 65 of the 100 questions.

2 Motivating Examples

In this section, we illustrate the presence of code reasoning questions in code review comments, while highlighting the challenges in extracting them. Reviewers’ identities are anonymized.

Observation 1: Most code review comments are not related to code reasoning. Through a manual analysis, we find that a majority of Github (65%) and Corporate (80%) code review comments are not related to code reasoning. We consider a comment to be related to code reasoning if in order to ask, answer, or address, it requires reasoning over reachability, data flow, control flow, or program variable and state. Instead, code review comments are often shallow edit suggestions related to style, structure, documentation, or syntactic reasoning. Consider Figures 1 (and 6 in appendix), in which the reviewers make shallow comments regarding style and *syntactic* reasoning respectively. During our analysis, we also found comments that are discussions of the intended behavior or specification (Figure 7 in appendix). We quantify the density of these comment categories in Table 1.

	Corp.	Github
CRQ	20%	35%
Shallow Edit Suggestion	60%	35%
Func Behavior Discussion	20%	30%

Table 1: Code Review Comments By Type.

Observation 2: Code review comments are often not phrased as questions. Through our manual analysis, we find that even when the comment is related to code reasoning, it is very rarely phrased as a concise and unambiguous code reasoning question. Consider Figure 2, in which the comment is phrased as an edit suggestion (removing the call to `std::move`) rather than the underlying code reasoning question: *Does calling `std::move` on the return value `s.releasePeerSet()` im-*

part the program’s behavior? Furthermore, the comment contains extraneous information, referencing another reviewer. In Figure 3, the comment is posed as a question, but it is overly verbose. It consists of two sentences, one of which is an extraneous edit suggestion related to functional behavior. The first sentence, although related to code reasoning, is ambiguous and not contextualized in the reviewed code. It does not explicitly state which program variables “*something else*” encompasses. A concise, unambiguous rephrasing could be: *Can `error_code` hold a value other than `ECONNREFUSED` or `ECONNRESET`?*

We also observe that code reasoning questions can be categorized into two types of queries that encompass all CRQs: VALUE and EQUIV queries. A VALUE query (Figure 3) asks about the value or possible value of a variable or expression at a program point. An EQUIV (equivalence) query (Figure 2) asks if two segments of code have differences in behavior. EQUIV queries typically underlay an edit suggestion. We find that in both Github and Corporate code review comments, 75% of code reasoning questions are EQUIV queries while 25% are VALUE queries.

Observation 3: Answers to rephrased questions are not readily available. During our manual analysis, we inspected the developer’s responses to comments. Responses came in the form of a natural language reply and/or a code edit. Answers in the form of a developer reply suffer from the same ambiguities and verbosity as the reviewer comments. Answers in the form of an edit require careful manual inspection to connect the change to the underlying code reasoning question. Sometimes the comment is ignored and not addressed.

3 Technique

Figure 4 illustrates our overall technique, which leverages a Corporate code aware LLM in combination with human validation.

3.1 Classifying Comments

As discussed in **Observation 1**, a minority of code review comments are related to code reasoning. To reduce manual inspection, we create an LLM based Code Reasoning Classifier (Figure 12 in appendix) which takes the

```

libr/core/cmd_flag.c
591 + if (s3) {
592 +     *s3 = '\0';
593 +     if (!strncmp (s3+1, "base64:", 7)) {
594 +         comment = (char *) r_base64_decode_dyn (s3+8, -1);

```

Reviewer on Mar 31, 2019
spaces in s3+8

Figure 1: Shallow Edit Suggestion.²

```

src/ripple/app/misc/HashRouter.cpp
104 104 if (!s.shouldRelay(suppressionMap_.cLock().now(), holdTime_))
105 105     return boost::none;
106 106
107 - return std::move(s.peekPeers());
107 + return std::move(s.releasePeerSet());

```

Reviewer on Nov 16, 2016
I doubt the std::move is necessary. @Reviewer2 ?

Figure 2: Raw Code Review Comment.³

```

torch/Lib/c10d/Utils.cpp
144 + // ECONNREFUSED happens if the server is not yet listening.
145 + if (error_code == ECONNREFUSED) {
146 +     *anyRefused = true;
147 + }
148 + // ECONNRESET happens if the server's listen backlog is exhausted.
149 + if (error_code == ECONNRESET) {

```

Reviewer on Nov 25, 2019
can error_code ever be something else? Should we throw if we get an error_code that is not one of these two?

Figure 3: Raw Code Review Comment.⁴

raw reviewer comment and corresponding line of code and decides if it is related to code reasoning.

We evaluate the performance of our Code Reasoning Classifier prompt on 100 randomly selected, manually labeled comments as shown in Table 3. We also experiment with a keyword matching approach using a hand derived list of undesirable keywords (Figure 20 in appendix)⁵, but find it incurs significantly more false positives than our LLM classification. In summary, our classifier correctly identified 11 out of 20 Corporate and 22 out of 35 GitHub code review comments as related to code reasoning, while misidentifying only 6 and 9 comments respectively.

	Corp		Github	
	LLM	KW	LLM	KW
Precision	.64	.31	.71	.39
Recall	.52	.81	.63	.1
F1 Score	.57	.45	.67	.56

Table 2: Code Reasoning Classification performance of LLM and Keyword matching approaches.

3.2 Rephrasing Comments as CRQs

As discussed in **Observation 2**, comments are rarely phrased as concise questions

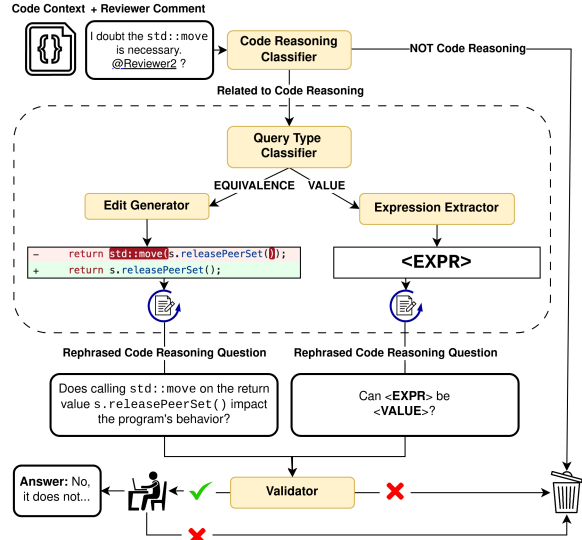


Figure 4: Benchmark Curation Methodology. Yellow boxes represent LLM prompts.

grounded over program elements. To avoid manual rephrasing, we again leverage the Corporate LLM. This portion of our technique is shown in the dotted box in Figure 4 and is invoked on samples that have been accepted by the CODE REASONING CLASSIFIER. Our technique invokes different rephrasing techniques for each query type. The Query Type Classifier (Figure 13 in appendix) classifies a comment as either an EQUIV query or VALUE query, triggering the appropriate rephrasing technique based on the classification.

When the Query Type Classifier decides the comment is an EQUIV query, we use chain of thought (Wei et al., 2022) reasoning to effectively rephrase. Since EQUIV queries are typically underlying edit suggestions, we employ an LLM based Edit Generator (Figure 14 in appendix) to perform the reviewer’s suggested edit as a link in a chain of thought. The edit is leveraged to rephrase the reviewer comment using a few shot prompt (Figure 16 in appendix).

When the Query Type Classifier decides the comment is a VALUE query, we similarly use a two step inference process similar to a chain of thought. As a first step, our Expression Extractor uses a few-shot prompt (Figure 15 in appendix) leveraging the code context and reviewer comment to extract the relevant program expression: <EXPR>. The relevant <EXPR> is used as a link to rephrase

the reviewer comment as a code reasoning question over the given expression using a few-shot prompt (Figure 17 in appendix).

Lastly, the rephrased question is given to the Validator for a self-consistency (Wang et al., 2023) check to reduce the occurrence of poorly rephrased code reasoning questions that are not faithful to the original line comment. The Validator prompt (Figure 18 in appendix) asks the LLM to decide if, given the original code, the reviewer comment has the same meaning as the rephrased comment. If the LLM Validator confirms the consistency, the rephrased question is selected as a confident candidate, and given to a human inspector to validate.

We evaluate our technique’s effectiveness in rephrasing code review comments into concise and unambiguous code reasoning questions. Our rephrasing approach (entire dotted box component in Figure 4) is evaluated on both Corporate (150 samples) and Github (160 samples) code review comments that were flagged as related to code reasoning by our Code Reasoning Classifier. The samples were manually inspected and labelled as correct if they were concise, unambiguous, and faithful to the original reviewer comment. We achieved a precision of .66 on Corporate code review comments and .63 on Github pull request comments.

3.3 Evaluation

We evaluate our methodology in terms of manual human curation required. In a purely manual approach, a human curator would need to inspect and classify 500 Corporate (or 285 Github) code review comments and manually rephrase 100 questions. Using our proposed methodology, a human curator would need to inspect only 150 Corporate (or 160 Github) code review comments without the need for any manual rewriting. Figure 21 (in appendix) and Figure 5, respectively, illustrate this comparison. The pencil indicates a manual rephrasing while the magnifying glass indicates inspection using our proposed technique. In summary, our cooperative LLM + human validation approach reduces the number of samples required to inspect by 1.8x on Github pull request comments and 3.3x on

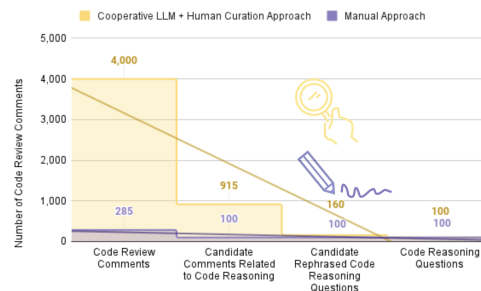


Figure 5: Manual vs Cooperative LLM Curation on Github Pull Request Comments.

Corporate code review comments.

4 GPT-4 Performance

We evaluate GPT-4 on our benchmark by prompting it with the surrounding function context (Figure 11 in Appendix). We evaluate the outputs manually considering a response to be accurate if it is both correct and contextually relevant. Correctness refers to the technical accuracy of the generated natural language response. Contextual relevance refers to the degree that the response is grounded in the given code context. We find the GPT-4 provides an accurate response on 65 of the 100 queries and are almost always (94%) grounded in the given code context.

	Acc	Total	%
	65	100	65%
VALUE	33	54	61%
EQUIV	32	46	70%

Table 3: Performance of GPT-4 on CRQBench.

Lastly, we conducted an error analysis to categorize the 35 incorrect responses. The majority of errors (25 instances) were due to the model lacking necessary context, such as usages of the given function, definitions of a used function or macro, or usages of a variable. Five errors were attributed to gaps in C++ knowledge, and the remaining five were due to incorrect evaluation of logic. Examples of each error scenario are shown in the Appendix (Figures 8 - 10).

We also experimented with evaluating the 7 billion parameter open source model Falcon (Almazrouei et al., 2023), but found it to have a much lower accuracy (25%) as it is a significantly smaller model.

5 Limitations

5.1 Extracting Answers to CRQs

To extract answers, we use an entirely manual based approach. A human curated an answer through a best effort approach by inspecting the cloned repository at the commit being reviewed. The answer is derived by reasoning over the code context, edit made (or not), and developer textual responses in the comment thread. In essence, our benchmark gathers the response which was implicitly provided by the developer, rather than an answer verified by a symbolic program analysis approach. We default to manual curation of answers due to the challenges presented in Section 2.

5.2 Size of Target Environment

Although the number of samples to inspect or rephrase is greatly diminished with our approach, the total number of comments needed to arrive at 100 code reasoning questions is much larger. Our cooperative approach requires greater than 10x more code review comments to derive 100 CRQs. This is due to false negatives in Code Reasoning Classifier and the Validator.

References

Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Coljocar, Mérouane Debbah, Étienne Goffinet, Daniel Hesslow, Julien Launay, Quentin Malaric, Daniele Mazzotta, Badreddine Noune, Baptiste Pannier, and Guilherme Penedo. 2023. The Falcon Series of Open Language Models. [arXiv:2311.16867](https://arxiv.org/abs/2311.16867)

Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. 2021. Program Synthesis with Large Language Models. [arXiv preprint arXiv:2108.07732](https://arxiv.org/abs/2108.07732) (2021).

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen

Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language Models are Few-Shot Learners. [CoRR abs/2005.14165](https://arxiv.org/abs/2005.14165) (2020). [arXiv:2005.14165](https://arxiv.org/abs/2005.14165) <https://arxiv.org/abs/2005.14165>

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harrison Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Heben Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Joshua Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. 2021. Evaluating Large Language Models Trained on Code. [CoRR abs/2107.03374](https://arxiv.org/abs/2107.03374) (2021). [arXiv:2107.03374](https://arxiv.org/abs/2107.03374) <https://arxiv.org/abs/2107.03374>

Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2022. PaLM: Scaling Language Modeling with Pathways. [arXiv:cs.CL/2204.02311](https://arxiv.org/abs/2204.02311)

²https://github.com/radareorg/radare2/pull/13555#discussion_r270676569

³https://github.com/XRPLF/rippled/pull/1904#discussion_r88226072

⁴https://github.com/pytorch/pytorch/pull/30354#discussion_r350485189

⁵We also experimented with a hand derived desirable keyword list (Figure 19 in appendix). We report results for the undesirable keyword approach as it achieved a higher F1 score.

414	Carlos E. Jimenez, John Yang, Alexander Wet-	Bob McGrew, Scott Mayer McKinney, Christine	476
415	tig, Shunyu Yao, Kexin Pei, Ofir Press, and	McLeavey, Paul McMillan, Jake McNeil, David	477
416	Karthik Narasimhan. 2023. SWE-bench: Can	Medina, Aalok Mehta, Jacob Menick, Luke	478
417	Language Models Resolve Real-World GitHub	Metz, Andrey Mishchenko, Pamela Mishkin,	479
418	Issues? arXiv:cs.CL/2310.06770	Vinnie Monaco, Evan Morikawa, Daniel Moss-	480
		ing, Tong Mu, Mira Murati, Oleg Murk, David	481
419	Zhiyu Li, Shuai Lu, Daya Guo, Nan Duan, Shailesh	Mély, Ashvin Nair, Reiichiro Nakano, Ra-	482
420	Jannu, Grant Jenks, Deep Majumder, Jared	jeev Nayak, Arvind Neelakantan, Richard Ngo,	483
421	Green, Alexey Svyatkovskiy, Shengyu Fu, and	Hyeonwoo Noh, Long Ouyang, Cullen O’Keefe,	484
422	Neel Sundaresan. 2022. Automating Code	Jakub Pachocki, Alex Paino, Joe Palermo, Ash-	485
423	Review Activities by Large-Scale Pre-training.	ley Pantuliano, Giambattista Parascandolo, Joel	486
424	arXiv:cs.SE/2203.09095	Parish, Emy Parparita, Alex Passos, Mikhail	487
		Pavlov, Andrew Peng, Adam Perelman, Filipe	488
425	OpenAI, Josh Achiam, Steven Adler, Sandhini	de Avila Belbute Peres, Michael Petrov, Hen-	489
426	Agarwal, Lama Ahmad, Ilge Akkaya, Flo-	rique Ponde de Oliveira Pinto, Michael, Poko-	490
427	rencia Leoni Aleman, Diogo Almeida, Janko	rny, Michelle Pokrass, Vitchyr H. Pong, Tolly	491
428	Altenschmidt, Sam Altman, Shyamal Anad-	Powell, Alethea Power, Boris Power, Elizabeth	492
429	kat, Red Avila, Igor Babuschkin, Suchir Bal-	Proehl, Raul Puri, Alec Radford, Jack Rae,	493
430	aji, Valerie Balcom, Paul Baltescu, Haiming	Aditya Ramesh, Cameron Raymond, Francis	494
431	Bao, Mohammad Bavarian, Jeff Belgum, Irwan	Real, Kendra Rimbach, Carl Ross, Bob Rotsted,	495
432	Bello, Jake Berdine, Gabriel Bernadett-Shapiro,	Henri Roussez, Nick Ryder, Mario Saltarelli,	496
433	Christopher Berner, Lenny Bogdonoff, Oleg	Ted Sanders, Shibani Santurkar, Girish Sas-	497
434	Boiko, Madelaine Boyd, Anna-Luisa Brakman,	try, Heather Schmidt, David Schnurr, John	498
435	Greg Brockman, Tim Brooks, Miles Brundage,	Schulman, Daniel Selsam, Kyla Sheppard, Toki	499
436	Kevin Button, Trevor Cai, Rosie Campbell, An-	Sherbakov, Jessica Shieh, Sarah Shoker, Pranav	500
437	drew Cann, Brittany Carey, Chelsea Carlson,	Shyam, Szymon Sidor, Eric Sigler, Maddie	501
438	Rory Carmichael, Brooke Chan, Che Chang, Fo-	Simens, Jordan Sitkin, Katarina Slama, Ian	502
439	tis Chantzis, Derek Chen, Sully Chen, Ruby	Sohl, Benjamin Sokolowsky, Yang Song, Natalie	503
440	Chen, Jason Chen, Mark Chen, Ben Chess,	Staudacher, Felipe Petroski Such, Natalie Sum-	504
441	Chester Cho, Casey Chu, Hyung Won Chung,	mers, Ilya Sutskever, Jie Tang, Nikolas Tezak,	505
442	Dave Cummings, Jeremiah Currier, Yunxing	Madeleine B. Thompson, Phil Tillet, Amin	506
443	Dai, Cory Decareaux, Thomas Degry, Noah	Tootoonchian, Elizabeth Tseng, Preston Tuggle,	507
444	Deutsch, Damien Deville, Arka Dhar, David Do-	Nick Turley, Jerry Tworek, Juan Felipe Cerón	508
445	han, Steve Dowling, Sheila Dunning, Adrien	Uribe, Andrea Vallone, Arun Vijayvergiya,	509
446	Ecoffet, Atty Eleti, Tyna Eloundou, David	Chelsea Voss, Carroll Wainwright, Justin Jay	510
447	Farhi, Liam Fedus, Niko Felix, Simón Posada	Wang, Alvin Wang, Ben Wang, Jonathan	511
448	Fishman, Juston Forte, Isabella Fulford, Leo	Ward, Jason Wei, CJ Weinmann, Akila Wel-	512
449	Gao, Elie Georges, Christian Gibson, Vik	hinda, Peter Welinder, Jiayi Weng, Lilian Weng,	513
450	Goel, Tarun Gogineni, Gabriel Goh, Rapha	Matt Wiethoff, Dave Willner, Clemens Winter,	514
451	Gontijo-Lopes, Jonathan Gordon, Morgan Graf-	Samuel Wolrich, Hannah Wong, Lauren Work-	515
452	stein, Scott Gray, Ryan Greene, Joshua Gross,	man, Sherwin Wu, Jeff Wu, Michael Wu, Kai	516
453	Shixiang Shane Gu, Yufei Guo, Chris Hal-	Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming	517
454	lacy, Jesse Han, Jeff Harris, Yuchen He, Mike	Yuan, Wojciech Zaremba, Rowan Zellers, Chong	518
455	Heaton, Johannes Heidecke, Chris Hesse, Alan	Zhang, Marvin Zhang, Shengjia Zhao, Tian-	519
456	Hickey, Wade Hickey, Peter Hoeschele, Bran-	hao Zheng, Juntang Zhuang, William Zhuk, and	520
457	dson Houghton, Kenny Hsu, Shengli Hu, Xin	Barret Zoph. 2024. GPT-4 Technical Report.	521
458	Hu, Joost Huizinga, Shantanu Jain, Shawn	arXiv:2303.08774	522
459	Jain, Joanne Jang, Angela Jiang, Roger Jiang,		
460	Haozhun Jin, Denny Jin, Shino Jomoto, Bil-	Hugo Touvron, Thibaut Lavril, Gautier Izac-	523
461	lie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz	ard, Xavier Martinet, Marie-Anne Lachaux,	524
462	Kaiser, Ali Kamali, Ingmar Kanitscheider, Ni-	Timothée Lacroix, Baptiste Rozière, Naman	525
463	tish Shirish Keskar, Tabarak Khan, Logan	Goyal, Eric Hambro, Faisal Azhar, Aurelien	526
464	Kilpatrick, Jong Wook Kim, Christina Kim,	Rodriguez, Armand Joulin, Edouard Grave,	527
465	Yongjik Kim, Jan Hendrik Kirchner, Jamie	and Guillaume Lample. 2023. LLaMA: Open	528
466	Kiros, Matt Knight, Daniel Kokotajlo, Łukasz	and Efficient Foundation Language Models.	529
467	Kondraciuk, Andrew Kondrich, Aris Konstan-	arXiv:cs.CL/2302.13971	530
468	tinidis, Kyle Kosic, Gretchen Krueger, Vishal		
469	Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc	531
470	Leike, Jade Leung, Daniel Levy, Chak Ming	Le, Ed Chi, Sharan Narang, Aakanksha Chowd-	532
471	Li, Rachel Lim, Molly Lin, Stephanie Lin,	hery, and Denny Zhou. 2023. Self-Consistency	533
472	Mateusz Litwin, Theresa Lopez, Ryan Lowe,	Improves Chain of Thought Reasoning in Lan-	534
473	Patricia Lue, Anna Makanju, Kim Malfacini,	guage Models. arXiv:cs.CL/2203.11171	535
474	Sam Manning, Todor Markov, Yaniv Markovski,		
475	Bianca Martin, Katie Mayer, Andrew Mayne,	Jason Wei, Xuezhi Wang, Dale Schuurmans,	536

537 Maarten Bosma, brian ichter, Fei Xia, Ed Chi,
538 Quoc V Le, and Denny Zhou. 2022. Chain-
539 of-Thought Prompting Elicits Reasoning in
540 Large Language Models. In Advances in Neural
541 Information Processing Systems, S. Koyejo,
542 S. Mohamed, A. Agarwal, D. Belgrave, K. Cho,
543 and A. Oh (Eds.), Vol. 35. Curran Associates,
544 Inc., 24824–24837. [https://proceedings.
545 neurips.cc/paper_files/paper/2022/file/
546 9d5609613524ecf4f15af0f7b31abca4-Paper-Conference.
547 pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/9d5609613524ecf4f15af0f7b31abca4-Paper-Conference.pdf)


548 Li Zhong and Zilong Wang. 2023. Can ChatGPT
549 replace StackOverflow? A Study on Robustness
550 and Reliability of Large Language Model Code
551 Generation. arXiv:cs.CL/2308.10335

A Appendix

```

engines/wintermute/ad/ad_scene.cpp
1117 1117  bool AdScene::updateFreeObjects() {
1118 1118      AdGame *adGame = (AdGame *)_gameRef;
1119      -      // 3D-code removed
1120      -      // bool is3DSet;
1119  +      bool is3DSet = false;

```

 **Reviewer** on Jun 24, 2020


This would produce a warning about unused variable, no?

Figure 6: Shallow Edit Suggestion.⁵

```

src/pickup.cpp
157 +      g->u.invoke_item( &pseudo );
158 +      pseudo.ammo_consume( pseudo.ammo_required(), g->u.pos() );
159 +      return true;
160 +  };

```

 **Reviewer** on Oct 9, 2016

This function doesn't work for variable charge items, such as water purifier:

- `ammo_required` can't be relied on to provide the number of charges required, since it only depends on the tool itself
- Fails to return non-consumed charges.

Figure 7: Functional Behavior Discussion.⁶

⁵https://github.com/scummvm/scummvm/pull/2347#discussion_r445103384

⁶https://github.com/CleverRaven/Cataclysm-DDA/pull/18691#discussion_r82533193

⁷The prompt used in our technique uses Corporate internal code review comments. To share the prompt, comments are substituted with open source code reviews from Github. We attempt to find substitutions which are similar to our Corporate code based prompt.

Prompt: In the following code, Is there a difference between using `!WI_IsFlagSet` and `WI_IsFlagClear`?

```
...
case WM_WINDOWPOSCHANGING:
{
    // GH#10274 - if the quake window gets moved to another monitor via aero
    // snap (win+shift+arrows), then re-adjust the size for the new monitor.
    if (IsQuakeWindow())
    {
        // Retrieve the suggested dimensions and make a rect and size.
        LPWINDOWPOS lpwpos = (LPWINDOWPOS)lparam;

        // We only need to apply restrictions if the position is changing.
        // The SWP_flags are confusing to read. This is
        // "if we're not NOT moving the window"
        if (!WI_IsFlagSet(lpwpos->flags, SWP_NOMOVE))
        {
            // Figure out the suggested dimensions and position.
            RECT rcSuggested;
        }
    }
}
...
```

Response: Yes. `WI_IsFlagSet` and `WI_IsFlagClear` are (presumably, given they are not defined in the given code) utility methods for bit flag checking in the code.

`WI_IsFlagSet` checks whether a particular bit (or bits) is set in a value. If the bit is set, it means that the corresponding feature or property is enabled or true.

`WI_IsFlagClear`, on the other hand, checks whether a particular bit (or bits) is not set in a value. If the bit is not set, it means that the corresponding feature or property is disabled or false.

So the difference between them is that they check for opposite things:

- `WI_IsFlagSet` -> Checks if the bit(s) is/are set (1).
- `WI_IsFlagClear` -> Checks if the bit(s) is/are not set (0).

Which to use depends on what you are trying to check in your code.

Figure 8: Incorrect GPT Response: Logic Error. The model recognized that the flags check for opposite conditions but failed to deduce that the inverse of one flag is equivalent to the other.

Prompt: Can `uchar_from_utf8(&ch, &*iter, normalized_src.end() - iter)` be less than or equal to zero?

```
void imgtool::simple_charconverter::from_utf8(std::ostream &dest, const std::string &src) const
{
    // normalize the incoming unicode
    std::string normalized_src = normalize_unicode(src, m_norm);

    auto iter = normalized_src.begin();
    while(iter != normalized_src.end())
    {
        // get the next character
        char32_t ch;
        iter += uchar_from_utf8(&ch, &*iter, normalized_src.end() - iter);
    }
    ...
}
```

Response: No. The function `uchar_from_utf8()` is not defined or included in the provided code context, so it's a bit speculative to determine its exact behavior. However, normally, such a function would return the number of bytes processed while converting a UTF-8 sequence into a unicode code point (character).

By convention, functions like these generally return a non-negative integer indicating the number of bytes processed. However, it's not impossible for such a function to return a negative value or zero in case of error conditions like an invalid UTF-8 sequence or encountering the end of the input prematurely.

Therefore, whether `uchar_from_utf8(&ch, &*iter, normalized_src.end() - iter)` can be less than or equal to zero mostly depends on the implementation details of `uchar_from_utf8()`.

In general, good coding practice would suggest handling such possible error conditions. However, assuming well-formed UTF-8 data and no I/O errors, it would typically always return a positive number.

Without the code or documentation of `uchar_from_utf8()`, it's challenging to give a definitive answer.

Figure 9: Incorrect GPT Response: Missing Context. The model was unable to determine if the expression could be less than or equal to zero as it did not have access to the function definition.

Prompt: In the following code, is there a difference between using `std::make_pair` with `emplace_back` as opposed to not using `std::make_pair`?

```
...
if (timeout.count() != 0) {
    // Compute the expiration time for this message based on the timeout
    auto expirationTime = computeRpcMessageExpiryTime(timeout);

    // Add the Future to the right vector in the timeoutMap_
    {
        std::unique_lock<std::mutex> lock(timeoutMapMutex_);
        auto& timeoutFuturesVector = timeoutMap_[expirationTime];
        timeoutFuturesVector.emplace_back(
            std::make_pair(futureResponseMessage, timeout));
    }
    timeoutThreadCV_.notify_one();
}
...
```

Response: Yes, there is a difference between using `std::make_pair` with `emplace_back` versus not using `std::make_pair`.

In the context of this code,

```
timeoutFuturesVector.emplace_back(std::make_pair(futureResponseMessage, timeout));
```

is adding a `std::pair` containing the `futureResponseMessage` and `timeout` to the `timeoutFuturesVector`.

`std::make_pair` is a utility function that constructs a `pair` object with its first element set to `x` and its second element set to `y`.

`emplace_back` is used to construct and insert elements directly into a vector without the need for an intermediate pair object.

Here are the two cases for clarity:

1. **With `std::make_pair`:**

In this case, `std::make_pair` will construct a pair object from the passed-in arguments and then `emplace_back` just inserts this pair object into the vector.

2. **Without `std::make_pair`:**

If you don't use `std::make_pair` and directly pass arguments to `emplace_back`, it will try to construct an object of the vector's element type from the arguments.

That being said, in this case, without `std::make_pair`, you may get a compile error if the vector's element type can't be constructed directly from the arguments you pass to `emplace_back`. So if your vector is of `std::pair` type (like in your case it seems to be), then using `std::make_pair` is the right way.

Figure 10: Incorrect GPT Response: Incorrect C++ knowledge. The model incorrectly responded Yes. In particular, the model was incorrect that without using `std::make_pair`, a compile error may occur.

Given the following code, {QUESTION}
{SURROUNDING_FUNCTION}

Figure 11: Prompt to GPT-4 for CRQBench evaluation.

```

You are a senior, expert C/C++ programmer with a lot of experience analyzing C/C
++ code.
Given a C/C++ line comment:

{COMMENT}

associated with this line of code:

'''
{CODE_SEGMENT}
'''

please determine which of the following categories best classifies that
line comment:

1. unknown: unknown
2. reasoning: a code reasoning question such as requiring control-flow or value
   propagation
3. explanation: a question asking what the code does
4. structure: request to reorganize or refactor code, such as extracting classes
   or methods
5. planning: prioritization, planning, or choosing what to work on next
6. style: a code style, or readability question including things like renaming
   symbols
7. format: a comment or question about code formatting

and provide an explanation of why that comment should be in that specific
category.
If you are uncertain about the category, respond with "0".

Generate your output in YAML format like this:

Explanation: <why you chose a specific category>
Line comment category: <category>

Response:
'''yaml

```

Figure 12: Code Reasoning Classifier Prompt.

You are an expert software engineer, asked to determine whether or not a reviewer comment is about substituting code.

Here are some examples.

Example 1:

```
LINE_COMMENT:
'''
also check if the tf example is empty?
'''
```

1. No.

Example 2:

```
LINE_COMMENT:
'''
This stores a pointer to a temporary object. This pointer becomes invalid right
after this statement. I suggest you initialise it with a 'nullptr' and check
for that.
'''
```

1. No.

Example 3:

```
LINE_COMMENT:
'''
Please use stolen(fduri) instead of hardcoding 5
'''
```

1. Yes.

Example 4:

```
LINE_COMMENT:
'''
The copy of a single int (previous type of fg, bg) is fine here, however copying
the vector is not good. Use a const reference instead.
'''
```

1. Yes.

Example 5:

```
LINE_COMMENT:
'''
{COMMENT}
'''
```

1.

Figure 13: Query Type Classifier Prompt.⁷

You are an expert C++ coder and know how to review C++ code at a Corporation,
and how to
respond to reviewer comments.
In the code below, a reviewer has left a comment marked by LINE_COMMENT.

Code:
```\n{CODE\_SEGMENT\_WITH\_COMMENT}\n```

The comment is asking to consider the following original code and modified code:

Figure 14: Edit Generator Prompt.

You are an expert C++ coder and know how to review C++ code at a Corporation, and how to respond to reviewer comments. In the code below, a reviewer has left a comment marked by `LINE_COMMENT`.

Please answer the following question.

1. Can you please list the relevant program expression?

Here are some examples.

Example 1:

Code:

```
'''
int comfort = 0;
comfort_response_t comfort_response;
comfort_response.aid = &item("null");
[*] [LINE_COMMENT] This stores a pointer to a temporary object. This pointer
 becomes invalid right after this statement. I suggest you initialise it with
 a nullptr and check for that.
bool plantsleep = has_trait(trait_CHLOROMORPH);
bool fungaloid_cosplay = has_trait(trait_M_SKIN3);
bool websleep = has_trait(trait_WEB_WALKER);
'''
```

`LINE_COMMENT`: This stores a pointer to a temporary object. This pointer becomes invalid right after this statement. I suggest you initialise it with a `nullptr` and check for that.

1. `comfort_response.aid`

Example 2:

Code:

```
'''
Value *Arg = Call->getArgOperand(0);
Value *LoweredArg = getLoweredByValOperand(Arg, Builder);
for (Value *A : Call->arg_operands()) {
 DXASSERT(A == Arg, "oops");
[*] [LINE_COMMENT] how could A be different unless there's a fundamental problem
 with operand iteration or something?
}
HLModule::MarkPreciseAttributeOnValWithFunctionCall(LoweredArg, Builder, *
 m_pModule);
addToDeadInsts(Call);
'''
```

`LINE_COMMENT`: how could A be different unless there's a fundamental problem with operand iteration or something?

1. A

Example 3:

Code:

```
'''
{CODE_SEGMENT_WITH_COMMENT}
'''
```

`LINE_COMMENT`: {COMMENT}

Figure 15: Expression Extractor Prompt.<sup>7</sup>

Can you please rephrase a reviewer LINE\_COMMENT that suggests the following code change? The rephrased comment should include program elements. Here are a few examples.

Example 1:

Before:

```
void gyroDataAnalyse(const gyroDev_t *gyroDev, biquadFilter_t *notchFilterDyn) {
 static FAST_RAM float fftAcc[XYZ_AXIS_COUNT] = {0, 0, 0};
 static FAST_RAM uint32_t fftAccCount = 0;
}
```

After:

```
void gyroDataAnalyse(const gyroDev_t *gyroDev, biquadFilter_t *notchFilterDyn) {
 static FAST_RAM float fftAcc[XYZ_AXIS_COUNT];
 static FAST_RAM uint32_t fftAccCount = 0;
}
```

[\*] [LINE\_COMMENT] Zero-initialized by default as well.

Rephrased: Is there a difference between initializing fftAcc to {0, 0, 0} and not initializing fftAcc?

Example 2:

Before:

```
if (isIntegralType(iter.dtype(), false)) {
 AT_DISPATCH_INTEGRAL_TYPES(iter.dtype(), "fmod_cpu", [&]() {
 cpu_kernel(iter, [](scalar_t a, scalar_t b) -> scalar_t {
 return std::fmod(a, b);
 });
 });
}
```

After:

```
if (isIntegralType(iter.dtype(), false)) {
 AT_DISPATCH_INTEGRAL_TYPES(iter.dtype(), "fmod_cpu", [&]() {
 cpu_kernel(iter, [](scalar_t a, scalar_t b) -> scalar_t {
 return a % b;
 });
 });
}
```

[\*] [LINE\_COMMENT] Should we use fmod or % here?

Rephrased: Is there a difference between computing the modulus of a and b using std::fmod vs using the binary operator %?

Example 3:

Before:

```
{BEFORE}

```

After:

```
{AFTER}

```

[\*] [LINE\_COMMENT] {COMMENT}

Figure 16: Equiv Rewriter Prompt.<sup>7</sup>



You are an expert C++ coder and know how to review C++ code at a Corporation,  
and how to respond to reviewer comments.  
A reviewer has left a comment marked by LINE\_COMMENT.

Can you please rephrase the LINE\_COMMENT as a question over program expressions?  
The Question should start with the word "Can".

Here are some examples.

Example 1:

Comment:

'''

This stores a pointer to a temporary object. This pointer becomes invalid right  
after this statement.

I suggest you initialise it with a nullptr and check for that.

'''

Program Expression: comfort\_response.aid

Question: Is there a difference between initializing comfort\_response.aid with &  
item( "null" ) as opposed to nullptr?

Example 2:

Comment:

'''

how could A be different unless there's a fundamental problem with operand  
iteration or something?

'''

Program Expression: A

Question: Can A ever be equal to Arg?

Example 3:

Comment:

'''

{COMMENT}

'''

Program Expression: {EXPR}

Figure 17: Value Rewriter Prompt.<sup>7</sup>

You are an expert C++ coder and know how to review C++ code at a Corporation,  
and how to respond to reviewer comments.  
In the code below, a reviewer has left a comment marked by LINE\_COMMENT.

Code:

'''

{CODE}

'''

Given the LINE\_COMMENT, answer the following question with YES or NO:

Does '''{Q\_C}''' have the same meaning as the LINE\_COMMENT?

Figure 18: Validator Prompt.

Interesting Keywords:

```
"except"
"segfault"
" fault"
"precondition"
"assumption"
"undefined behavior"
" ub "
"null",
"reach"
"ever be true"
"ever be false"
"branch taken"
"branch not taken"
"deref"
"reference"
```

Figure 19: Desirable Keywords for Positive Matching.

Uninteresting Keywords:

```
"test"
\nit"
\follow up"
\log a higher level"
"log a lower level"
"logging"
\naming"
\readability"
\TODO"
\description"
\comment"
"typo"
"clang"
"style guide"
\period"
"restructure"
"restructuring"
"refactor"
"move"
"offline"
"space"
"spacing"
```

Figure 20: Undesirable Keywords for Negative Matching.

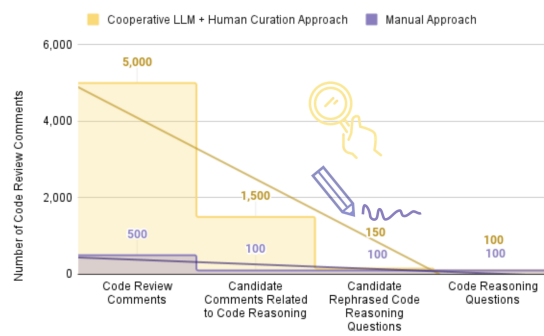


Figure 21: Manual vs Cooperative LLM Curation approaches on Corporate Code Review Comments.