

Empirical Computation

Eric Tang
Carnegie Mellon University

Marcel Böhme
MPI for Security and Privacy

Abstract

In this vision paper, we explore the challenges and opportunities of a form of computation that employs an empirical (rather than a formal) approach, where the solution of a computational problem is returned as empirically most likely (rather than necessarily correct). We call this approach as *empirical computation* and observe that its capabilities and limits *cannot* be understood within the classic, rationalist framework of computation. Hence, we appeal to the software engineering community to develop the foundations and techniques required to analyze the properties of empirical computation as it generates solutions to computational problems.

While we take a very broad view of "computational problem", a classic, well-studied example is *sorting*: Given a set of n numbers, return these numbers sorted in ascending order.

- To run a classical, "formal computation", we might first think about a *specific* algorithm (e.g., merge sort) to solve the problem before developing a *specific* program that implements it. The program will expect the input to be given in a *specific* format, type, or data structure (e.g., unsigned 32-bit integers). In software engineering, we have many approaches to analyze the correctness of such programs. From complexity theory, we know that there exists no correct program that can solve the average instance of the sorting problem faster than $O(n \log n)$.
- In contrast, to run an *empirical computation*, we might directly ask a large language model (LLM) to solve *any* computational problem (which can be stated informally in natural language) and provide the input in *any* format (e.g., negative numbers written as Chinese characters). There is no (problem-specific) program that could be analyzed for correctness. Also, the time it takes an LLM to return an answer is entirely *independent* of the computational complexity of the problem that is solved.

What are the properties of empirical computation? How to inquire about its fundamental limits? How can we analyze, estimate, or predict the correctness of empirical computation in the general, in the problem-specific, or in the instance-specific? The purpose of this paper is establish empirical computation as a field in software engineering that is timely and rich with interesting problems.

ACM Reference Format:

Eric Tang and Marcel Böhme. 2025. Empirical Computation. In . ACM, New York, NY, USA, 5 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>



This work is licensed under a Creative Commons Attribution 4.0 International License. Conference'17, July 2017, Washington, DC, USA
© 2025 Copyright held by the owner/author(s).
ACM ISBN 978-x-xxxx-xxxx-x/YY/MM
<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 Introduction

Large Language Models (LLMs) are posed to change the field of software engineering. For instance, it was only two years ago that Microsoft announced the first LLM-based code synthesis tool CoPilot [13]. Yet, recently Google reported that *25% of all of its new code is LLM-generated* [4]—despite concerns about "hallucination" as potential source of bugs and security flaws [14, 15].

Today, we solve computational problems¹ by *programming*, i.e., using a programming language, we tell the machine *how* instances of that problem are to be solved: We think about a *specific* algorithm before developing a *specific* program (or function) that implements it. The program will expect the input to be given in a *specific* format, type, or data structure (e.g., `uint32_t`). LLM-based code synthesis [13] allows developers to automate parts of that *how*: The developer tells the LLM informally in natural language what they need, and the LLM generates the corresponding code. At this point, it is not difficult to imagine that, for some problems (e.g., parsing from natural language [10, 12, 13]), the source code might not be needed at all, and is *substituted by a direct call to an LLM*.

In the future, we expect some computational problems to be solved by *prompting* (e.g., using an informal encoding as a natural language prompt to an LLM). The result of this empirical approach to computation is returned as empirically most likely rather than necessarily correct. In this vision paper, we explore the challenges and opportunities of this *empirical computation*.

We argue that the capabilities and limits of empirical computation cannot be understood within the classic, rationalist framework of computation and call on the software engineering community to develop new instruments for the analysis of empirical computation.

What are the properties of empirical computation? How to analyze the fundamental limits? Are there any guarantees? How can we analyze, estimate, or predict the correctness of empirical computation in the general, in the problem-specific, or in the instance-specific? The purpose of this paper is establish empirical computation as an emerging and interesting area in software engineering.

Programming vs prompting to solve problems. Today, we write programs to solve problems. As mentioned, a program represents *how* that problem is solved. Just like the problem can often be decomposed into smaller sub-problems, a program can be composed from many smaller computational units (e.g., functions). A program is a *specific* solution for a *specific* problem. The software engineering community has developed many approaches and techniques to analyze the properties (incl. correctness) of a program [2]. However, in the future, we might simply prompt an LLM to solve a problem. There is no program to analyze w.r.t. the problem it solves.

¹In this paper, we take a very broad view of the term *computational problem*. For instance, the computational problem that is solved by a given program is the relationship between the inputs and outputs of that program (or the "purpose" of that program).

What can we say about the result? How do we even test the correctness of LLMs for various computational problems?

Programming gives us absolute, direct, and precise control about the correctness and efficiency of a computation. The procedure of the computation is explicit and predefined. The fault of an incorrect computation can be localized and repaired. In contrast, prompting offers more flexibility particularly when requirements are uncertain. Without constraints on language, input structure, or algorithm, an “empirical computer” returns the result after interpreting the user-provided, informal description of the problem statement using the available context. If we found the result to be incorrect, **how can we programmatically improve correctness for future instances** (without introducing regressions for other problems)?

Formal vs informal representations of problem instances. In programming, inputs are provided using a *specific* format, type, or data structure. Everyone who, as part of a software testing class, has been asked to implement and to test a program for the triangle classification problem knows that we have to be very precise about that format. When parsing from a file, do we give lengths or angles, separated by commas or dashes? Are spaces and tabs allowed? When receiving as function input, should I use a primitive type or implement my own class? Are negative numbers allowed? What about floating point?

In *prompting*, the input can be provided informally in any format. The correct interpretation arises from context. Unlike in programming, there is no predefined “contract” between caller and callee that specifies the precise format or data structure for the input. This offers flexibility at the interface, but at the cost of ambiguity.

We experimentally explore the correctness of empirical sorting when numbers are provided not using digits but expressed entirely differently, e.g., using words in German (zweiundvierzig [speak, two-and-forty; 42]) or as characters in Korean. We find that the result of empirical sorting is more likely correct if numbers are provided in a language that is well-represented on the internet.

Limits in terms of efficiency vs correctness. The computational limits of programming to solve computational problems are well-studied in theoretical computer science. However, an LLM finds answers to any computational problem in a time that is *independent* of the computational complexity of the problem. For instance, an LLM might take as much time to predict the name of the last US president as would to predict the plaintext from a given SHA-512 hash. Of course, correctness is a different matter. Are there problems that are particularly amenable to empirical computation? **How can we quantify the limits of empirical computation in the general or in the problem-specific?**²

We experimentally explore the efficiency and correctness of empirical computation to solve various simple, well-studied computational problems, like sorting or searching. Our results confirm that average execution time is independent of the computational complexity of the problem.³ Yet, correctness decreases as input size increases. For instance, *empirical sorting* or *empirical searching* (in sorted or unsorted lists) returns the correct result for a list of 50 random numbers with 50% probability (which is quite impressive

²We note that PAC learning [18] and learnability theory offers insight only on the requirements for (or dependence on) the training data (i.e., sample complexity).

³In fact, time seems dependent on the size of the *output* more than the input.

without an algorithm). We also find that correctness increases if the problem instance is “more familiar” (e.g., given by another LLM).

Formal vs informal notions of correctness. In programming, it is upon the programmer to elicit and to understand the computational problem to be solved (i.e., the requirements) before implementing a program that is meant to solve *all* instances of that problem. If the computational problem is precisely understood (such as sorting or searching), it can be encoded as formal specification, and given the specification, the programmer could formally verify, i.e., *guarantee* the correctness of her program w.r.t. that specification.

However, sometimes the problem is not precisely understood and can only be described vaguely (cf. requirements elicitation). In such cases, prompting offers more flexibility and allows a contextual interpretation of a vaguely described computational problem (which may even be further elicited interactively). Yet, unlike in formal language, a natural language description can be ambiguous, and the way the prompt is designed can substantially impact the correctness of the response (cf. prompt engineering). **What are effective ways to describe a computational problem to maximize correctness?**

2 Preliminary Experiments

To explore the challenges and opportunities of the empirical approach to computation, we conducted a large number of experiments. While we take a very broad view of the term “computational problem”, for our experiments, we focus on simple, classic, well-studied problems, like sorting or searching, whose correctness is well-defined and computational complexity is well-studied.

2.1 Experimental Setup

Implementation. For our evaluation, we used LLM-GPT-Sort,⁴ an existing, independently developed Python tool that prompts an LLM to sort a sequence of numbers. The concrete *prompt* is

“Sort the elements in the given collection in ascending order, and return only the sorted collection in the list format: <numbers>.”

The input was provided as a Python list which LLM-GPT-Sort translated into a string (numbers). The output was again provided as a (hopefully sorted) Python list which LLM-GPT-Sort generated by parses the LLM response. We extended the tool to also support

- searching sorted lists; $O(\log n)$,
- searching unsorted lists; $O(n)$,
- computing the longest palindromic substring; $O(n)$ [11], and
- finding a subset with a given subset sum; $O(2^{\frac{n}{2}})$ [6].

The LLM used for all our experiments is GPT-3.5 Turbo. In order to mitigate the impact of randomness, we repeated each experiment at least 30 times and report averages.

Methodology. For sorting, we generated Python lists of random numbers and of random length using existing functions from the standard Python library `random` (e.g., `random.sample`). As discussed in more detail in the individual sections, we varied properties like precision and magnitude of the numbers, as well as the length of

⁴<https://github.com/njmarko/llm-gpt-sort>

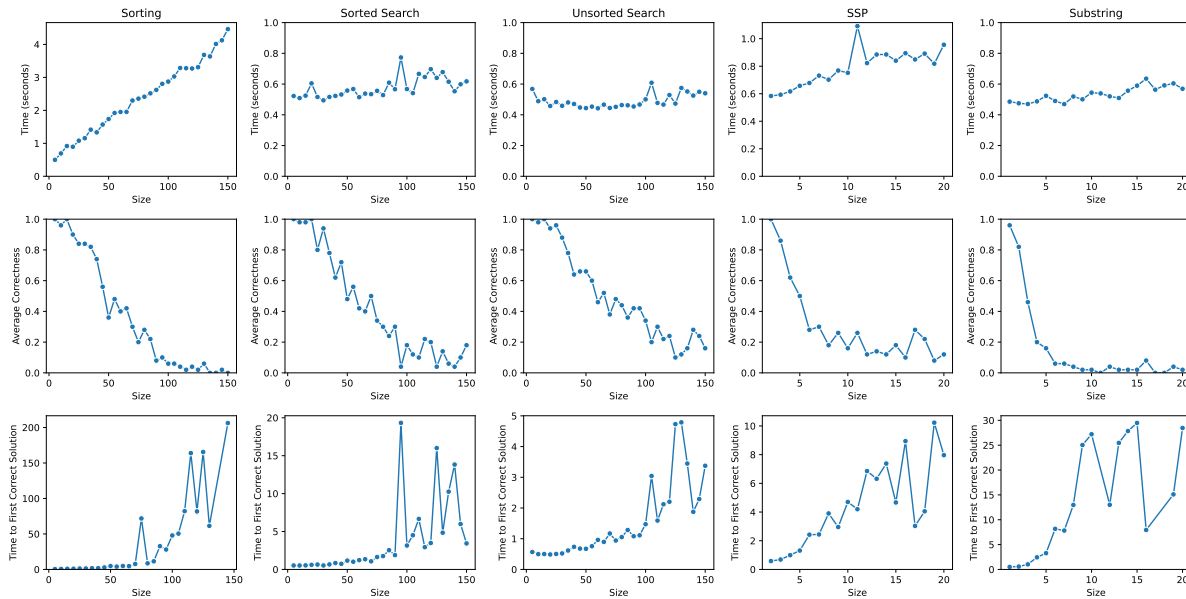


Figure 1: LLM-performance on various well-studied computational problems. *Top*: Average time to solution. *Middle*: Average proportion of correct solutions (among 30+ repetitions). *Bottom*: Expected time to generate the first correct solution.

the array. For every randomly generated input, we recorded the time taken (i.e., *efficiency*) and whether the LLM-generated result was correct (i.e., *correctness*).

For the other problems, we take a similar approach. For *searching*, we generate a random or sorted list of numbers as above, select a random element to search for and prompt the LLM to return the corresponding index or negative one (-1). For the subset sum problem (*SSP*), we generate a random list of numbers as above, find the sum S of a random subset, and prompt the LLM to return that subset whose total is S . For the longest palindromic substring problem (*substring*), we construct a random string of characters and prompt the LLM to return the result. For the latter two, we implement simple Python programs as ground truth.

2.2 Efficiency versus Correctness

Efficiency. Figure 1.top shows the efficiency of empirical computation, i.e., the average time it took the LLM to solve instances of various computational problems, as a function of instance size n . For *sorting*, time appears to be linearly increasing in list size. The increase is likely due to the corresponding increase in the number of input and output tokens [19], rather than the computational complexity of the sorting problem, i.e., $O(n \log n)$. Indeed, the execution time is relatively constant (under one second) for both *search* problems where the output is just a single token (an array index or negative one). Note that the classic theory of computation predicts a difference between the average complexity of searching sorted versus unsorted lists. For SSP, we explain the slightly increasing time with the opportunity to *output* longer subsets.

Correctness. Figure 1.middle shows the correctness of empirical computation, i.e., the proportion of empirically solved instances that are actually correct, as a function of n . For all problems, we can

see an obvious decrease in correctness. Interestingly, for sorting and both searching problems, the reduction in correctness in n is roughly equivalent: For Python lists with 50 integers,

- the probability that empirical sorting returns the correctly sorted list is equivalent to a coin flip (0.5), and
- the probability that empirical searching returns the correct index is also equivalent to a coin flip (0.5).

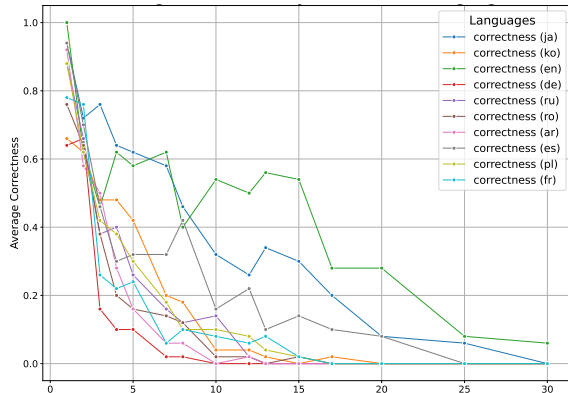
In contrast, for the substring problem, correctness reduces to close to zero (0) already for strings with five characters likely due to tokenization. For the subset sum problem, correctness is close to zero (0) already for lists of length 10. If we combine the efficiency and correctness results into the expected time until the LLM returns the first correct solution (Fig. 1.bottom), we can clearly see the consequence of the decrease in correctness on the performance of empirical computation for these problems.

There might be several (non-traditional) mechanisms to improve the correctness of empirical computation, e.g., with methods like majority voting, prompt engineering [16], fine-tuning [5], reinforcement learning [7], or test-time-training [9]. Currently, for our sorting problem, ChatGPT 3.5 Turbo would sometimes add numbers not present in the input or remove numbers that were present. Some numbers would appear multiple times, and sometimes the list would not be sorted entirely. Large inputs (e.g., $n > 200$) would sometimes be truncated or produce repeated sequences.

Familiarity. We conjectured that the empirical availability of similar problem instances in the training data improves the correctness of the empirical computation. As preliminary study of this hypothesis, we *prompted the LLM* to generate a random sequence of numbers of a specific length, and then ask another LLM (instance) to sort this sequence. We assume that the LLM would be *more familiar* with the self-generated “random” sequence of numbers. Table 1

Table 1: Correctness of sorting “familiar” instances. Number of instances actually returned with that size in parenthesis.

| Array Size | Random | Familiar |
|------------|--------|------------|
| 10 | 0.99 | 1.00 (100) |
| 20 | 0.95 | 1.00 (100) |
| 30 | 0.81 | 0.95 (100) |
| 40 | 0.66 | 0.67 (122) |
| 50 | 0.56 | 0.70 (161) |

**Figure 2: Correctness of empirical sorting for various languages (measured as the proportion of correctly sorted lists).**

shows the corresponding correctness for various array sizes. Indeed, the correctness on familiar instances significantly improved.

2.3 Inputs in Natural Language

A fundamental difference between the formal approach to computation and the empirical approach is that no “formal contract” is required for the latter on how the inputs are provided. In contrast, a program requires inputs to be passed in a certain structure that is formally agreed on beforehand. That input structure could be a specific file format, like this PDF file, a specific data structure, like a linked list, graph, tree, or object, or it could be a primitive data type, like a (signed or unsigned) integer or float.

In contrast, an LLM takes inputs flexibly, either informally in natural language or formally, e.g., as structured JSON file. The LLM infers from context how the provided input is interpreted. This flexibility enables broader applicability but also introduces ambiguities. There is no certainty about how an input is interpreted.

Methodology. We used a Python library called num2words to translate the random numbers in the list into words or characters in different languages. The list of languages is shown on the top right in Figure 2. Then, we executed LLM-GPT-Sort on the resulting list as input and used the num2words library to parse the LLM output.

Results. Figure 2 shows the correctness of empirical sorting, i.e., the proportion of empirically solved instances that are actually correct, as a function of n , when words or characters are used for numbers instead of digits. In comparison to Figure 1.middle, we can immediately see that the correctness of empirical sorting reduces more quickly. For languages that are well represented in the internet, correctness of empirical sorting is significantly better

than for underrepresented languages. For instance, sorting a list of 15 numbers represented in the english language is correct with more than 50% probability, while when represented in Korean is never correct in 50 repetitions (0%).

3 Perspective and Vision

Everyone talks about trustworthy artificial intelligence. Yet, no one knows how to define, test, improve, or guarantee the correctness of the response of a Large Language Model (LLM) to a given prompt. Our classic theory of computation offers no help in answering such questions; nor does the powerful program analysis machinery that we have developed in the software engineering community (e.g., static or dynamic analysis, verification, or model checking).

It is our vision that the analysis of the properties of empirical computation will emerge as a new area in software engineering that is both timely and rich with interesting problems.

Given an empirical computer, like an LLM, what type of statements can we make about the correctness on various computational problems (whether or not formally specified [17])? Can we make absolute or only statistical statements [1, 3, 8]? Can we rely on existing formal methods, like analysis or proof, or would empirical methods be more suitable, like statistical, counterfactual, or causal reasoning?

What are the most effective encodings of a problem statement in natural language [20]? Can we predict the correctness for a given problem instance (e.g., based on an estimate of “familiarity”)? How can we (reactively—or better proactively) improve correctness on specific computational problems or their instances in the absence of a (problem-specific) program?

How can we maximize the generality of our statements about the correctness of empirical computation? Can we make statements about correctness across *all instances* of a problem? Can we make general statements about the correctness of an empirical computer, i.e., across all problems, whether known or unknown?

The analysis of correctness and other properties is a classic problem in software engineering. Hence, we call on the software engineering research community to develop the tools and techniques required to analyze the correctness of empirical computation as an important step toward the systematic analysis of the trustworthiness of artificial intelligence.

Acknowledgments

This research was conducted as part of the CS@max planck internship program.

References

- [1] Marcel Böhme. 2022. Statistical Reasoning About Programs. In *Proceedings of the 44th International Conference on Software Engineering (Pittsburgh, USA) (ICSE 2022)*. 5 pages. <https://doi.org/10.1145/3510455.3512796>
- [2] Marcel Böhme, Eric Bodden, Tevfik Bultan, Cristian Cadar, Yang Liu, and Giuseppe Scanniello. 2025. Software Security Analysis in 2030 and Beyond: A Research Roadmap. *ACM Transactions on Software Engineering and Methodology* (2025), 25 pages. <https://doi.org/10.1145/3708533>
- [3] Marcel Böhme, Danushka Liyanage, and Valentin Wüstholtz. 2021. Estimating Residual Risk in Greybox Fuzzing. In *Proceedings of the 15th Joint meeting of the European Software Engineering Conference and the ACM SIGSOFT Symposium on the Foundations of Software Engineering (ESEC/FSE)*. 494–504. <https://doi.org/10.1145/3468264.3468570>

- [4] Benji Edwards. 2024. Google CEO says over 25% of new Google code is generated by AI. <https://arstechnica.com/ai/2024/10/google-ceo-says-over-25-of-new-google-code-is-generated-by-ai/> (Accessed on 15/01/2025).
- [5] Zeyu Han, Chao Gao, Jinyang Liu, Jeff Zhang, and Sai Qian Zhang. 2024. Parameter-Efficient Fine-Tuning for Large Models: A Comprehensive Survey. arXiv:2403.14608 [cs.LG] <https://arxiv.org/abs/2403.14608>
- [6] Ellis Horowitz and Sartaj Sahni. 1974. Computing Partitions with Applications to the Knapsack Problem. *Journal of the ACM* 21, 2 (1974), 277–292.
- [7] L. P. Kaelbling, M. L. Littman, and A. W. Moore. 1996. Reinforcement Learning: A Survey. arXiv:cs/9605103 [cs.AI] <https://arxiv.org/abs/cs/9605103>
- [8] Seongmin Lee and Marcel Böhme. 2023. Statistical Reachability Analysis. In *Proceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE 2023)*. 12. <https://doi.org/10.1145/3611643.3616268>
- [9] Jian Liang, Ran He, and Tieniu Tan. 2024. A Comprehensive Survey on Test-Time Adaptation Under Distribution Shifts. *International Journal of Computer Vision* 133, 1 (July 2024), 31–64. <https://doi.org/10.1007/s11263-024-02181-w>
- [10] Lezhi Ma, Shangqing Liu, Yi Li, Xiaofei Xie, and Lei Bu. 2024. SpecGen: Automated Generation of Formal Program Specifications via Large Language Models. arXiv:2401.08807 [cs.SE] <https://arxiv.org/abs/2401.08807>
- [11] Glenn Manacher. 1975. A New Linear-Time “On-Line” Algorithm for Finding the Smallest Initial Palindrome of a String. *J. ACM* 22, 3 (July 1975), 346–351. <https://doi.org/10.1145/321892.321896>
- [12] Ruijie Meng, Martin Mirchev, Marcel Böhme, and Abhik Roychoudhury. 2024. Large Language Model guided Protocol Fuzzing. In *Proceedings of the Network and Distributed System Security Symposium (NDSS'24)*. 15 pages.
- [13] Microsoft. [n. d.]. Copilot. <https://copilot.microsoft.com/>. (Accessed on 15/01/2025).
- [14] Hammond Pearce, Baleegh Ahmad, Benjamin Tan, Brendan Dolan-Gavitt, and Ramesh Karri. 2022. Asleep at the Keyboard? Assessing the Security of GitHub Copilot’s Code Contributions. In *SP. IEEE*, 754–768.
- [15] Neil Perry, Megha Srivastava, Deepak Kumar, and Dan Boneh. 2023. Do Users Write More Insecure Code with AI Assistants?. In *Proceedings of the 2023 ACM SIGSAC Conference on Computer and Communications Security (CCS '23)*. Association for Computing Machinery, New York, NY, USA, 2785–2799. <https://doi.org/10.1145/3576915.3623157>
- [16] Pranab Sahoo, Ayush Kumar Singh, Sriparna Saha, Vinija Jain, Samrat Mondal, and Aman Chadha. 2024. A Systematic Survey of Prompt Engineering in Large Language Models: Techniques and Applications. arXiv:2402.07927 [cs.AI] <https://arxiv.org/abs/2402.07927>
- [17] Ion Stoica, Matei Zaharia, Joseph Gonzalez, Ken Goldberg, Koushik Sen, Hao Zhang, Anastasios Angelopoulos, Shishir G. Patil, Lingjiao Chen, Wei-Lin Chiang, and Jared Q. Davis. 2024. Specifications: The missing link to making the development of LLM systems an engineering discipline. arXiv:2412.05299 [cs.SE] <https://arxiv.org/abs/2412.05299>
- [18] L. G. Valiant. 1984. A theory of the learnable. *Commun. ACM* 27, 11 (Nov. 1984), 1134–1142. <https://doi.org/10.1145/1968.1972>
- [19] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Proceedings of the 31st International Conference on Neural Information Processing Systems (Long Beach, California, USA) (NIPS'17)*. Curran Associates Inc., Red Hook, NY, USA, 6000–6010.
- [20] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. arXiv:2201.11903 [cs.CL] <https://arxiv.org/abs/2201.11903>