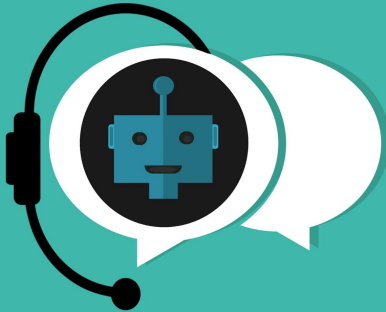


An Evaluation of ChatGPT's ability on generating code solutions for LeetCode Questions



Group 5

- Shiyu Xiu
- Yuangan Zou
- Qian Tang



Motivation

- Explore the possibility of applying ChatGPT as a code-writing assistant during development
- Evaluate the quality of code generated by ChatGPT using Python, Java, C
- Code generation for Leetcode questions provide insights for ChatGPT's ability in solving real-world questions



Research Questions

- *How does the ChatGPT perform in terms of generating correct code solutions to leetcode questions with different languages?*
 - *Out of 100 questions, How many of the generated solutions are correct and accepted?*
- *How does the ChatGPT perform in terms of running time and memory consumed when solving leetcode questions with different languages?*



Research Questions Cont

- *How does the ChatGPT perform in terms of **Cyclomatic and Cognitive complexity** of the code when solving leetcode questions with different languages?*
 - **Cognitive Complexity:** a measure of how difficult a unit of code is to intuitively understand.
 - **Cyclomatic complexity:** a count of the linearly independent paths through source code.
- *Under what circumstances does ChatGPT **produce an incorrect answer** to Leetcode questions more often?*
 - *hard questions, new questions, complex questions*



Method

A case study with quantitative analysis:

- % of correct answers
- Comparison of running time and memory consumed among languages
- Analysis of readability with a human-assigned score (1 ~5)

Quantitative analysis applied since:

- involve numerical data (# of (in)correct answers)
- statistical methods required



Data collection

Data mining technique: using python selenium library to automatically crawl the leetcode website to obtain the list of problems and its associated metadata such as difficulty, problem description.



Data Analysis

Quantitative analysis:

- Quantitative analysis on the correctness of code solutions
- Evaluate Cyclomatic and Cognitive Complexity with metrics
- Conduct a within-group survey-like study for the readability with human-assigned score (1 ~ 5)



Expected Results

- Table 1 showing : problem info

% passed test cases

Time & Space Complexity

reasons for failure

- Table 2 showing: Readability scores

Cyclomatic and Cognitive Complexity scores

Initial results

Python:

Problem number	Problem title	Problem Context	Acceptance	Difficulty	Answers by ChatGPT	Succeeded	Runtime	runtime_b	memory	memory_b	error_type	error_message	Total_testcases	Test_cases_passed
1	Two Sum	Given an array of integers, return indices of the two numbers such that they add up to a specific target.	49.6%	Easy	class Solution: def twoSum(self, nums: List[int], target: int) -> List[int]:	TRUE	61	79.46	15.1	37.43	None	None	None	None
2	Add Two Numbers	You are given two non-empty linked lists representing two non-negative integers. The digits are stored in the reverse order of the list.	40.3%	Medium	# Definition for singly-linked list. class ListNode: def __init__(self, x): self.val = x self.next = None	TRUE	64	85.11	14	32.77	None	None	None	None
3	Longest Substring Without Repeating Characters	Given a string, find the length of the longest substring without repeating characters.	33.8%	Medium	class Solution: def lengthOfLongestSubstring(self, s: str) -> int:	TRUE	61	76.93	13.9	88.64	None	None	None	None

C:

Problem number	Problem title	Problem Context	Acceptance	Difficulty	Answers by ChatGPT	Succeeded	Runtime	runtime_b	memory	memory_b	error_type	error_message	Total_testcases	Test_cases_passed
1	Two Sum	Given an array of integers, return indices of the two numbers such that they add up to a specific target.	49.7%	Easy	/** Note: You are not allowed to modify the input array. */	TRUE	132	65.42	6.3	91.3	None	None	None	None
2	Add Two Numbers	You are given two non-empty linked lists representing two non-negative integers. The digits are stored in the reverse order of the list.	40.3%	Medium	/** Definition for singly-linked list. struct ListNode { int val; struct ListNode *next; };	TRUE	19	24.99	7.9	32.88	None	None	None	None
3	Longest Substring Without Repeating Characters	Given a string, find the length of the longest substring without repeating characters.	33.8%	Medium	int lengthOfLongestSubstring(char* s)	TRUE	3	87.48	5.9	46.69	None	None	None	None

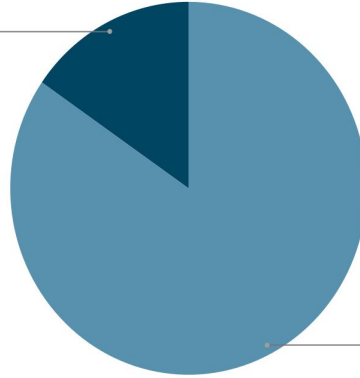
Java:

Problem number	Problem title	Problem Context	Acceptance	Difficulty	Answers by ChatGPT	Succeeded	Runtime	runtime_b	memory	memory_b	error_type	error_message	Total_testcases	Test_cases_passed
1	Two Sum	Given an array of integers, return indices of the two numbers such that they add up to a specific target.	49.7%	Easy	class Solution { public int[] twoSum(int[] nums, int target) {	TRUE	1	99.31	43	22.77	None	None	None	None
2	Add Two Numbers	You are given two non-empty linked lists representing two non-negative integers. The digits are stored in the reverse order of the list.	40.3%	Medium	Solution: /** Definition for singly-linked list. public class ListNode {	TRUE					compile error	compile error	None	None
3	Longest Substring Without Repeating Characters	Given a string, find the length of the longest substring without repeating characters.	33.8%	Medium	class Solution { public int lengthOfLongestSubstring(String s) {	TRUE	2	99.94	42.7	48.66	None	None	None	None

Analysis of Initial Results

LeetCode Result for Python

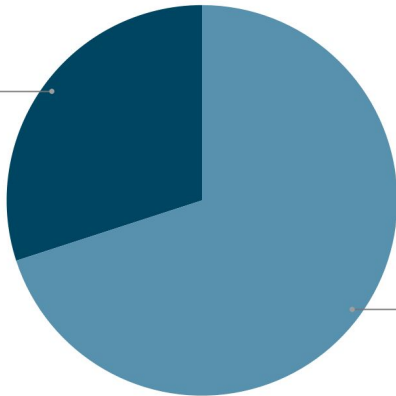
Wrong Answer
15.4%



Correct Answer
84.6%

LeetCode Result for C

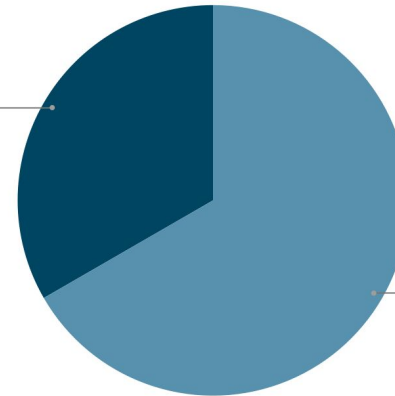
Wrong Answers
30.0%



Correct Answers
70.0%

LeetCode Result for Java

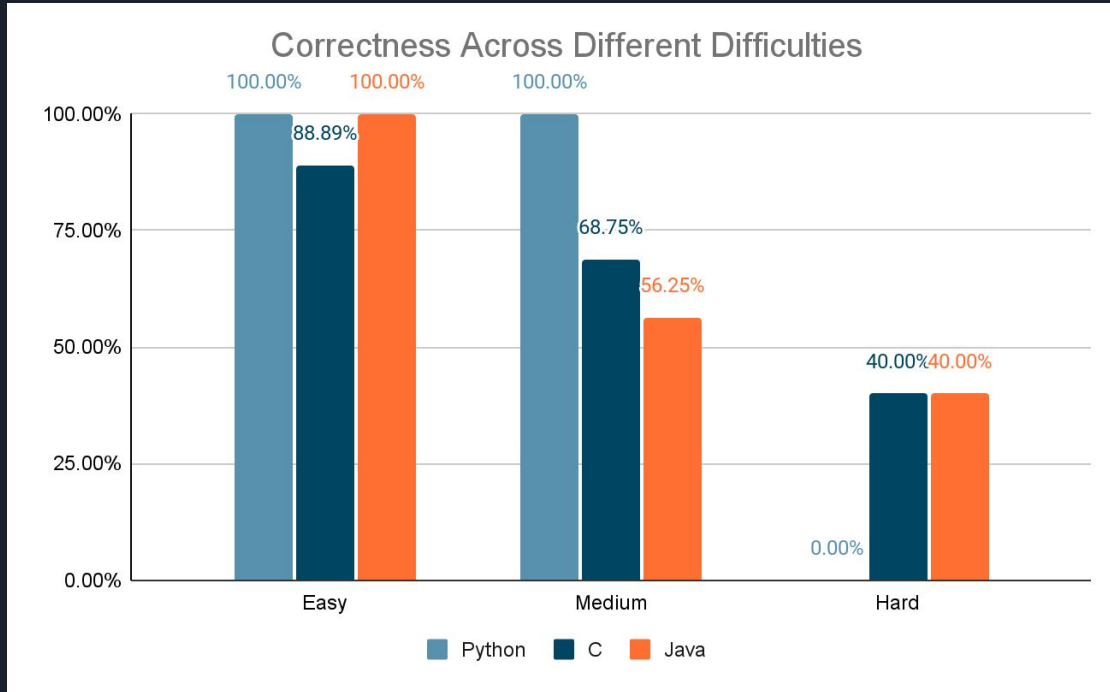
Wrong Answer
33.3%



Correct Answer
66.7%

Analysis of Initial Results cont.

% correct answers across different difficulties.





Implications & limitation

- Better understand the ability of cutting edge AI in code generation.
- Help software developers decide whether to use AI tools for writing code.
- Ways of improvement for ChatGPT as a code writer
- Compatibility
- Design of functions



Thank you !