

 **Important:** Make a copy of this template (File > Make a copy) and use it while reviewing

Rewritten response:

All of these must be true for a correct response rewrite:

General for Reviewing Rewrites:

The rewritten code corrects any errors from previous turns and fully satisfies the prompt's requirements. Making the necessary changes to fix and improve the original model response

Any explanations provided alongside the code are accurate and clarify the changes made or the solution overall.

The rewrite should be classified as “no issues” based on all of the per-turn ratings described above.

Code Testing:

If no rewrite is required (in the case of one good and one bad response), does the good response's code run?

Yes, then proceed

No, then change the toggle to say a rewrite is required, and proceed to rewrite the response and proceed

(When a rewrite is present) Does the code run in the rewrites?

Yes, then proceed

No, then fix the code

Is the code optimal and if code efficiency (i.e. O(n^2) vs O(1), etc) is described in the response, is it accurately described? Make sure to verify.

Yes, then proceed

No, then update accordingly

Does the code contain sufficient error handling and input validation?

Yes, then proceed

No, then update accordingly

Style and Presentation:

Does the code contain sufficient comments?

Yes, then proceed

No, then add comments

Does the rewritten response replace any paragraphs (especially those with at least 3 points) with bullet points instead?

Yes, then proceed

No, then re-word into three bullet points

Is all repetitive phrasing/wording removed and made more concise?

Yes, then proceed

No, then make the necessary edits

Does the rewritten response satisfy all of the requirements of the prompt?

Act like a lawyer when it comes to the prompt: for every ask or request in the prompt, does the rewritten response satisfy it?

Yes, then proceed

No, then edit the rewrite further to make sure both the text and code address the prompt

When should you do a rewrite?:

For any turn:

The attempter will be required to **always** rewrite the **preferred** response to achieve the goal of the prompt and fix any issues that were identified.

The only scenario where a rewrite would **not** be necessary is if one of the model responses is completely **perfect**, meeting all the requirements and specifications of the prompt. In this case, indicate that a **response rewrite is not needed**. **This is very rare and requires a thorough justification, in most cases, the response will require a rewrite.**

Specification List

- The re-written response should **fully achieve the most recent prompt**.
- The re-written response should **address any secondary objectives implied by the prompt**.
- The re-written response should **correct all errors (major or minor)**
- The re-written response should be **coherent and logically connected to the prior conversation**.
- The re-written response should **fulfill all the dimension for "No Issues" in the rating rubrics**.
- The rewrite should have **input validation** if the code accepts input from the user
- The re-written response should **follow the formatting and styling specifications**.
- The code of the re-written response should **contain enough comments**.
- The code of the re-written response MUST run properly and be optimal and efficient. Failure to properly test code results in removal from the project.**

In short, they are doing a rewrite step, this response is required to be perfect! All identified issues must be fixed, regardless of their severity.

Goal-Oriented Multi-Turn (MT) Coding

Task Specifications

Evaluating Responses & Continuing the Task

In this step, you will assess the model's responses based on specific criteria and provide follow-up prompts to improve its output.

3a. Producing at least one bad response

Every step should involve **at least one bad model response resulting from the prompt you provide**. If both responses are good, please write a more challenging prompt and try again.

How to quickly identify a good vs bad response: If there are any areas in the rating dimensions where a response has **at least one issue in the P0 or P1 level, see below**, we can count that response as "bad."

If a response only has issues in the presentation criteria (P2) that is **not** sufficient enough to be rated as "bad"!

3b. Rating the Responses based on each dimension

For each turn, after specifying how the two models did overall, each individual response **should** be rated on performance according to the following dimensions along with a brief explanation.

Quick snapshot of the dimensions (more detailed table below) and the level of priority (meaning these dimensions are most important vs less important):

Level of priority should be used to **determine which response is better with regards to your rating for the turn**. For example, if R1 has an instruction following issue while R2 has efficiency issues, R2 would be the better response.

Instruction Following: The response answers all requests in the prompt

This checks if the model understood **all** of the requests of the prompt and is addressing each one.

Make sure to check the entire implementation. If the model *tries* to address the request but *fails to*, it is an **accuracy** issue.

Priority Level: P0 - Highest Importance

Accuracy: All claims and code in the response is accurate and fully correct

You will need to Google the claims made by the model or execute code to check this

Priority Level: P0 - Highest Importance

Sub category - Input Validation/Error Handling: The response covers all meaningful input edge cases and handles them correctly.

The response validates all inputs, handling invalid data.

Priority Level: P0 - Highest Importance

Optimality and Efficiency: The response presents the most optimal and efficient solutions

The response is using common practices and standards

Priority Level: P1 - Second Highest Importance

Up-to-Date: The response uses only the most recent APIs, functions, or libraries available.

APIs, functions, or libraries used aren't causing compilation or runtime errors due to deprecation.

Priority Level: P1 - Second Highest Importance

Presentation: The response format follows the Style and Presentation guidelines

It should follow the presentation rubric, such as:

Enough comments in the code.

A professional tone (no pleasantries / fluff).

An answer that is concise, without repetitive statements.

Explanations that use bullet points.

Priority Level: P2 - Third Highest Importance

Each dimension rating **should** include a brief explanation as to why the given rating was chosen (e.g., why "Major Issues" was selected for Accuracy). The explanations/justifications

you write for dimension ratings should be **specific** and **clear**. If there are any bugs, or there are issues, specify **what the issues are** and **what the causes of the issues are**.

Avoid generalizations like “No Issue”, “N/A”, “This has an error”, etc. If there are no issues, give a very brief explanation as to why.

Examples of Bad Justifications:

“There is a syntax error”: ✗ This is only half of the story - you should briefly explain what’s causing the syntax error.

“The response doesn’t satisfy all of the prompt requirements”: ✗ This is very vague. If you find that the response doesn’t satisfy all of the requirements, you should be specifying which requirements it doesn’t satisfy and why it doesn’t satisfy them.

Response Rating Rubric

Dimensions	NA (0)	No Issue (1)	Minor Issue (2)	Major Issue (3)
Instruction following	<i>This dimension cannot be NA.</i>	The response meets the main request and all constraints, showing a strong understanding of the prompt, even if there are minor implementation errors. It handles any ambiguities well and stays within the specified requirements.	The response fulfills the primary request but does not entirely adhere to all the constraints. The response could have better handled the ambiguity of the prompt. Common errors: Fails some but not ALL constraints	The response fails to fulfill the primary request OR fulfills the primary request but does not adhere to <i>any</i> constraints. Common errors: Fails to do the primary request
Accuracy and [NEW] Input Validation / Error Handling	Can only be NA if the response contains no code or factual claims, and does not rely on prior context.	The code runs error-free , produces the correct output, and follows best practices. All text and comments are accurate, and the response is contextually appropriate with any previous errors fixed. — All meaningful edge cases are covered. The response includes thorough validation for expected inputs and error handling for invalid data, ensuring robustness and resilience to common input errors. Example: The function <code>calculate_discount(price, discount_percentage)</code>	The code runs but has minor warnings or low-risk security issues. The content is mostly accurate, but some statements are unclear or make unproven claims. Previous errors remain but don't affect the current response. — Some edge cases are covered. The response handles most expected inputs but misses certain edge cases, which could lead to potential errors or exceptions under specific conditions. The function <code>calculate_discount(price, discount_percentage)</code>	The code doesn't run due to logic errors, produces incorrect output, or has major security flaws. The response includes false claims, lacks context, and previous errors were not fixed, making the issues worse. — No edge cases are covered. The response lacks validation for all inputs, making it vulnerable to errors when faced with unexpected or invalid data inputs. The function <code>calculate_discount(price, discount_percentage)</code> performs no validation on its inputs, assuming <code>price</code> and <code>discount_percentage</code> are valid.

		<code>discount_percent</code> age) validates all inputs. It checks that <code>price</code> and <code>discount_percent</code> age are positive numbers, ensures <code>discount_percent</code> age does not exceed 100%, and returns a clear error message if values are out of expected ranges or of incorrect types (e.g., strings).	<code>discount_percent</code> age) includes basic validation, such as checking that <code>price</code> and <code>discount_percent</code> age are positive numbers. However, it lacks checks for certain edge cases, such as ensuring <code>discount_percent</code> age does not exceed 100% or verifying that inputs are numeric.	<code>discount_percentage</code> are always valid and within expected ranges. This could cause runtime errors or incorrect results if given invalid inputs, such as negative numbers, <code>discount_percentage</code> over 100%, or non-numeric types, making the function unreliable.
Optimality and Efficiency	Can only be NA if the response contains no code using functions or statements aside from the assignment	The code is well-optimized, handles edge cases, and follows standard best practices. If top performance isn't required, it still performs efficiently without adding unnecessary complexity.	The code performs well but could use minor optimizations. It generally follows best practices but may not scale for large datasets.	The code exhibits severe performance and efficiency issues. The code does not adhere to common practices and standards.

<p>Presentat ion</p> <p>WHEN REWRITI NG, YOU MUST FIX ALL PRESEN TATION ISSUES</p>	<p>This dimen sion canno t be NA.</p>	<p>The code is well-documented, with clear comments and explanations for any modifications. Code included in the prompt that did not originally have comments should have comments if included in the response. The response is concise, well-organized, and uses readable variable and function names. Complex processes are broken down with bullets, and Markdown is correctly formatted with clear hierarchies.</p> <p>Formatting is neat, with triple backticks for code blocks, and proper use of bold and italics for emphasis. White space and line breaks improve readability, and tables are correctly aligned. Functions are modular and follow standard patterns, such as using <code>if __name__ == "__main__":</code> blocks for structure. There are no redundant solutions provided for the same problem.</p>	<p>The documentation is generally clear but could use more detail. There are minor language errors that don't affect readability, and formatting could be improved for clarity. Variable and function names are understandable, but some structural changes—like adding bullets or logical sections—would help. Functions are present but may need more modularity, and some explanations are missing, making the code harder to follow in parts.</p> <p>Common Errors</p> <ul style="list-style-type: none"> Uses backticks inconsistently Uses camelCase and snake_case inconsistently 	<p>The documentation is missing or inadequate, or lacking code comments entirely, making the code hard to understand. The response is poorly formatted and lacks structure, with unclear variable and function names. The logic is disorganized, and there are no explanations for key decisions, making it difficult to follow, integrate, or reuse. Programming language tags are also missing.</p>
<p>Up-to-Dat e</p>	<p>NA (0)</p>	<p>Up-to-Date</p>	<p>Out-of-Date</p>	
		<p>The code does not call on any libraries or functions.</p>	<p>The code uses the most fresh API, libraries, or functions available to solve problems efficiently. The code uses a maintained library or function which is an older version that still works (even if it is less efficient).</p>	<p>The code uses a deprecated API, library or function, causing a runtime or compile-time error.</p>

Rewriting Each Preferred Response

In this step, once you have selected the better response from each turn, you will be required to **always** rewrite the **preferred** response to achieve the goal of the prompt and fix any issues that were identified.

The only scenario where a rewrite would **not** be necessary is if one of the model responses is completely **perfect**, meeting all the requirements and specifications of the prompt. In this case, indicate that a **response rewrite is not needed**. This is very rare and requires a **thorough justification, in most cases, the response will require a rewrite**.

Specification List

- The re-written response should **fully achieve the most recent prompt**.
- The re-written response should **address any secondary objectives implied by the prompt**.
- The re-written response should **correct all errors (major or minor)**
- The re-written response should **be coherent and logically connected to the prior conversation**.
- The re-written response should **fulfill all the dimensions for "No Issues" in the rating rubrics**.
- The re-written response should **follow the formatting specifications**.
- The code of the re-written response should **contain enough comments**.
- The code of the re-written response MUST run properly and be optimal and efficient. Failure to properly test code results in removal from the project.**

Formatting/Presentation Requirements

- Key terms should be highlighted in bold, whereas titles, articles, etc. are italicized.
- Remove pleasantries such as "Sure," "Certainly," "I can help with that," etc.
- Make responses more concise, remove all fluff and unnecessary phrases (i.e. "*Welcome to the world of VS Code!*")
- Tone should be straightforward/professional.
- Code should be well-commented
- Test outputs include a comment with the expected response.
- Explanations should use bullet points.
 - Rewritten response replaces any paragraphs (especially those with at least 3 points) with bullet points instead
- All repetitive phrasing/wording must be removed.

Here are examples of [formatting/presentation requirements](#)
Make sure to keep track of the changes that you made as you'll have to write them out in the next step.

Stylistic Guidance & Changes for Rewrites

Stylistic Change	Description	Bad Example	Better Example
------------------	-------------	-------------	----------------

Descriptive Code Comments	<p>Provide comments that explain the code's logic and intent rather than stating the obvious.</p>	<pre>Python # Loop through the list for item in items: process(item)</pre>	<pre>Python # Iterate over the list of items to apply a discount to each for item in items: apply_discount(item)</pre>
Concise Natural Language Explanations	<p>Avoid repetitive phrasing (e.g., "I added," "We created"). Combine related actions into single sentences to improve clarity.</p>	<p><i>Repetitive phrasing:</i> "I added a guessedWords array to keep track of words that have already been guessed." "I updated the checkSecretWord function to check if an input word has already been guessed." "I added a checkIfAllGuessed function."</p>	<p>A guessedWords array was added to track guessed words. The checkSecretWord function displays an alert if a word has been guessed. The checkIfAllGuessed function ensures all secret words are handled and manages animations."</p>
Use Lists or Bullet Points for Clarity	<p>Use bullet points or numbered lists when describing more than three actions or steps.</p>	<p><i>Not itemized:</i> "First, we import libraries like pandas, numpy, and sklearn. Then, we create a dataset, convert categorical variables, split data, and train a model."</p>	<p>Import libraries: pandas, numpy, sklearn. Create a dummy dataset with 5 attributes and 1 target variable. 1-hot coded categorical variables.. Split data into training and testing sets. Train a Random Forest model.</p>
Consistent Formatting	<p>Maintain consistent formatting when referring to libraries, methods, or code snippets.</p>	<p><i>Inconsistent formatting:</i> Highlighting one keyword but not the other. "The code uses matplotlib and tkinter."</p>	<p><i>Consistent formatting:</i> Highlighting both keywords for clarity and consistency. "The code uses matplotlib and tkinter."</p>

Use Lists for Multiple Steps	Organize steps or instructions using numbered lists to clearly present each step.	"Check the input, validate user data against the database, and store data in the database."	Use a list for clarity: 1. Check the input for errors. 2. Validate user data against the database. 3. Store the data in the database.
-------------------------------------	---	---	---

Lemur Astrologer Coding — Attempter's Resource Masterlist

Congratulations on pre-qualifying for our Lemur Astrologer Coding project! 🎉

We're so excited to have you on board. Our goal is to ensure you can work on this project with confidence and ease. This guide provides **all the resources** you need to get started.

Bookmark this page!

What You'll Be Doing:

The goal of this project is to train AI models for coding tasks. Your job is to write high-quality prompts that challenge the model, ensuring at least one response contains a substantial error — this is referred to as a "**deviation**".

Resources:

Documentation (*must read*):

[Lemur Astrologer Tasking Specifications](#)

You'll find examples of **good vs. bad prompts** at the bottom of the document.

[Lemur Astrologer Common Errors](#)

You'll find the most common pitfalls that contributors make on the project, so that you can avoid them when you begin tasking.

Live Tasking Session (*optional but recommended*):

We host multiple live tasking sessions for this project each day for those who are in the limited production phase (see **Limitations** below). You'll find the schedule in your dashboard if you're in this phase.

A lead will be available on these calls to give you live feedback as you task. Please join these sessions with a live task available to review.

Community Support (*optional but recommended*):

Lemur Astrologer forum:

You will be automatically invited to one of the Lemur Astrologer forums on discourse.

Reach out to us if you encounter any issues during this project & connect with fellow taskers.

Project Onboarding:

As you begin working on Lemur Astrologer, you'll complete a few courses and benchmark tasks *before* moving on to production tasks. Our goal is to ensure you're well-prepared to succeed and understand how to effectively challenge the model.

Courses:

Lemur Astrologer Language Survey

In this brief course, you will be asked to select all the coding languages that you are familiar with so that you will be served tasks you can complete.

Lemur Astrologer Intro Module

This instructive course teaches the step-by-step of how to successfully complete a task on this project, including:

1. Goal Setting
2. Prompt Writing and Tagging
3. Response Evaluation
4. Code Execution
5. Ranking Responses
6. Rewriting Responses
7. Continuing the Conversation

Good and Bad Prompts Walkthrough Course

This quick course covers examples of both good and bad prompts and the reasons why they are good or bad.

Attempter's Resource Masterlist

This one-page course lists all of the resources needed for this project.

Lemur Astrologer Quiz

This certification course asks questions that are answerable using the materials covered in the previous onboarding steps. You must answer at least 8/10 of the questions correctly to proceed with tasking on this project.

Production Tasking:

Video Walkthrough: [video1369749349.mp4](#)

Limitations:

We want to be transparent about the limitations we put in place in our current project setup. This helps ensure the quality and consistency of the work being submitted.

Task Limit: You may be limited to submitting 2 tasks within a 24-hour period.

Evaluation Criteria: These tasks will be closely evaluated based on quality, accuracy, and overall performance.

Trusted Status: If your tasks demonstrate high standards and reliability, the submission limit will be lifted for you.

Key takeaway: The fastest way to remove the task limit is by submitting 2 exceptional tasks.



Lemur Astrologer Coding

Goal-Oriented Multi-Turn (MT) Coding Reviewer Checklist

Overview:

[How to Use this Checklist:](#)

[When to fix vs. SBQ a task:](#)

[Prompt Quality](#)

[Ratings Per model response](#)

[Rewritten response:](#)

[How to score a task](#)

! Changelog

Date:	Change Description
10/31/24	<ul style="list-style-type: none">! MAJOR update - New sub-category under accuracy: Input Validation! MAJOR update - Rewrites are now required for any issues in the response (including 1 or more minor issues)1 minor issue is enough to classify a response as “bad”
10/29/24	<ul style="list-style-type: none">- Clarification about major issues for presentation - this includes repetition and use of paragraphs for 3+ points<ul style="list-style-type: none">- Any scenario where a model response doesn't follow the stylistic guidelines is considered to be a major issue

If a task is in an incorrect category or difficulty, do NOT SBQ it, instead select the correct category in the multiple choice field within the task.

How to Use this Checklist:

This checklist is designed to guide you in reviewing tasks. It gives criteria for assessing the prompt quality and the attempter's ratings.

- Make a copy of this template (File > Make a copy) and use it while reviewing

When to FIX vs. SBQ a task:

1. **Terrible prompt?**
 - If the prompt is low quality, spam, or not sufficient for the project guidelines, SBQ the task. ✗
2. **No Significant Failures?**
 - If there are no substantial issues or failures in either model response for a turn, SBQ the task. ✗
3. **Disagree with Ratings?**
 - If you disagree with selected ratings (accuracy, instruction following, etc.) adjust them in the per-turn steps and update the justifications. ✓
4. **Goal or Category Adjustments?**
 - If the goal, category, difficulty or similar elements need updates, make those adjustments. ✓
5. **SxS Rating Correction:**
 - If you believe the “**better**” model response remains the same, but want to adjust the degree (e.g., “slightly better” to “much better”), you may change it without affecting future turns. ✓
 - If only the **justification needs updating**, please adjust as needed. ✓
 - **Changing which response is “better”** requires re-rolling the prompt for future turns to reflect the updated context. ● Move to bullet point 7
 - Please watch this [quick 3 minute video](#) for more information
6. **Rewrite Correction:**
 - For **minor fixes** (grammar, wording, small presentation changes), edit the rewrite without re-rolling future turns. ✓
 - For example, changing the natural explanation underneath a code block to use bullet points rather than a paragraph.
 - For **major adjustments** (fixing code bugs, adding missing features), re-roll the prompt to ensure future turns reflect the new context.
 - If this is the last turn in a task, changing the rewrite will have no affect
 - For any other turns ● Move to bullet point 7.
 - Major adjustments are changes that will noticeably impact the response in future turns. For instance, if you edit the rewrite to include comments in the code, but the remaining turns are based on a version without comments, this creates a clear inconsistency.
 - Anytime you need to edit/modify the code in a rewrite, this will be a **major adjustment**
7. **What happens if I need to re-roll the next or current turn?**
 - If you can finish the task (including all re-rolled turns) within a reasonable time, go ahead and make the fix. ✓
 - If the prompt is too complex or there are too many turns to complete feasibly, SBQ the task. ✗

Prompt Quality

This section helps you evaluate the quality of the prompt. Review each point to make sure that the prompt aligns with the assigned task's requirements. All of these must be true for a correct prompt

1st turn only: Prompt fits the task category:

- The prompt is correctly labeled and fits the category (e.g., writing code, fixing bugs).
 - Does the prompt match the **difficulty level** at the top of the task?
 - Yes, then proceed
 - No, then edit the difficulty field in the prompt categorization section
 - Does the subcategory of the prompt meet the requirements?
 - Yes, then proceed
 - No, then edit the sub-category field - we have to make sure it's correctly labeled

Prompt is clear and detailed:

- The prompt has all necessary information and isn't confusing.
- There are no contradicting requests in the prompt

Correct language and tools used:

- The prompt uses the right programming language or library as instructed.

For multi-turn tasks only, each prompt builds on the previous one:

- Each prompt logically continues from the previous one
- The conversation is fluid and natural
- Does the subcategory of the prompt meet the requirements? It likely won't always be the same and may need to be adjusted.

IMPORTANT: Simplified Example of this common sub-category error:

Task Category: Code Generation/Synthesis

Turn 1 Prompt: "Make a 2d minigolf game using JS"

Turn 1 Category: Text to Code (*makes sense*)

Turn 2 Prompt: "Additionally add a stopwatch to the top right corner"

Turn 2 Category: Text to Code (~~bad~~ bad, it should be **Text to Code Edits**, because we're asking for edits to existing code)

Ratings Per model response

How to use:

This section guides you through reviewing and evaluating each model response based on the attempter's rating. For each response, use the criteria below to determine if the rating is accurate.

- **For ratings of 3 (Major Issue):** **At least one** of the conditions listed under "Major Issue" in the category must be true for a rating of 3.
- **For ratings of 2 (Minor Issue):** **At least one** of the conditions under "Minor Issue" in the category must be true for a rating of 2.
- **For ratings of 1 (No Issue):** **All conditions** listed under "No Issue" in the category must be true for a rating of 1.

Instruction Following

- **If the attempter rated 3 (major issues):**
 - The response fails to fulfill the primary request OR fulfills the primary request but fails to adhere to any constraints.
- **If the attempter rated 2 (minor issues):**
 - The response fulfills the primary request but does not entirely adhere to all the constraints.
 - The response could have better handled the ambiguity of the prompt.
- **If the attempter rated 1 (no issues):**
 - The response fulfills the primary request and all of the constraints.
 - The response may have incorrect implementation (e.g., mistakes in code) if it shows a deep understanding of the prompt request and constraints in the written text.
 - The response adeptly handles any potential ambiguity in the prompt.
 - The response fully adheres to all the prompt's requirements/constraints and doesn't do more than the user requested.

Accuracy of Response (including input validation)

- **If the attempter rated 3 (major issues):**
 - The code fails to execute due to flawed logic.
 - The code output does not align with the expected program input.
 - The response contains false or inaccurate claims.
 - The code contains severe security vulnerabilities.
 - Relative to the context of the conversation, the response in the most recent turn does not make sense.
 - Errors from the prior turn were not corrected and compounded.
 - No edge cases are covered. The response lacks validation for all inputs, making it vulnerable to errors when faced with unexpected or invalid data inputs.
- **If the attempter rated 2 (minor issues):**
 - The code executes but contains non-failing minor errors; errors that would show up as a linter error but not a compilation error or would cause a warning upon running the code. An issue such as bad type setting in a dynamic

programming language like Python or JS (not TS) would also fall under this category.

- The code contains some low-risk security vulnerabilities.
- The response is factual and accurate.
- Relative to the conversation context, coherence is ambiguous or incomplete.
- The response falsely asserts claims that are not fully proven or controversial as fact.
- Some edge cases are covered. The response handles most expected inputs but misses certain edge cases, which could lead to potential errors or exceptions under specific conditions.
- **If the attempter rated 1 (no issues):**
 - The code executes without errors and generates an output that aligns precisely with the model intent.
 - The code is safe, free of vulnerabilities and follows best practices.
 - All written text (including inline code comments) is factual and accurate.
 - Relative to the context of the conversation, the most recent turn makes sense.
 - Prior turn errors were fully corrected, or there were no prior errors.
 - All meaningful edge cases are covered. The response includes thorough validation for expected inputs and error handling for invalid data, ensuring robustness and resilience to common input errors.

Optimality and Efficiency

- **If the attempter rated 3 (major issues):**
 - The model's response contains avoidable inefficiencies—such as poor algorithm design, unnecessary complexity, or incorrect performance.
 - The code is inefficient, especially where performance matters.
- **If the attempter rated 2 (minor issues):**
 - The code is relatively performant, but some low-lift optimizations could still be done.
 - The code mostly adheres to common practices and standards.
 - The code may not be scalable in real-world large-dataset use cases.
- **If the attempter rated 1 (no issues):**
 - The code is highly performant and is optimized even for edge cases.
 - The code is sufficiently performant if the prompt does not request the best performance and if a more efficient option would require much more complexity.
 - The code adheres to common practices and standards.

Presentation

- **If the attempter rated 3 (major issues):**
 - Documentation is missing or insufficient to understand the code.
 - Not Including Any Code Comments

- The response has poor readability due to a lack of structure or formatting.
- The code is missing programming language tags.
- The response includes repetitive content
- The response does not replace any paragraphs **with at least 3 points** with **bullet points** instead
- The code has lousy variable or function names.
- The logic in the code is disorganized and difficult to follow.
- Explanations are absent, leaving the reader unclear about key decisions or steps in the code.
- Code is poorly structured, making it challenging to integrate or reuse.
- **If the attempter rated 2 (minor issues):**
 - The documentation is sufficient to understand the code, but additional details would be helpful.
 - The response contains a few language mechanics errors, but they do not impact readability.
 - The response readability can be improved with better formatting.
 - Some adjustments to formatting and structure could improve clarity, such as adding more bullet points or logical sections.
- **If the attempter rated 1 (no issues):**
 - The code has an adequate amount of documentation, including in-code comments.
 - The response is clear, concise, and well-structured.
 - Variables and functions are named with readability in mind.
 - All modifications are documented in the code or in a written explanation outside the code.

Out of Date:

- **If the attempter rated 2 (major issues):**
 - The code uses a deprecated API, library or function or ones with known inefficiencies or vulnerabilities, or are generally not recommended AND stop the code from compiling
- **If the attempter rated 1 (no issues):**
 - The code uses a maintained library or function that is an older version or less efficient but still allows the code to run.

Rewritten response:

All of these must be true for a correct response rewrite:

General for Reviewing Rewrites:

- The rewritten code corrects any errors from previous turns and fully satisfies the prompt's requirements. Making the necessary changes to fix and improve the original model response
- Any explanations provided alongside the code are accurate and clarify the changes made or the solution overall.

Code Testing:

- If no rewrite is required (in the case of one good and one bad response), does the good response's code run?
 - Yes, then proceed
 - No, then change the toggle to say a rewrite is required, and proceed to rewrite the response and proceed (or SBQ)
- (When a rewrite is present) Does the code run in the rewrites?
 - Yes, then proceed
 - No, then fix the code + re-roll the next prompt or SBQ if you are time constrained
- Is the code optimal and if code efficiency (i.e. $O(n^2)$ vs $O(1)$, etc) is described in the response, is it accurately described? Make sure to verify.
 - Yes, then proceed
 - No, then update accordingly
- Does the code contain sufficient error handling and input validation?
 - Yes, then proceed
 - No, then update accordingly

Style and Presentation:

- Does the code contain sufficient comments?
 - Yes, then proceed
 - No, then add comments + re-roll the next turn or SBQ
- Does the rewritten response replace any paragraphs (especially those with at least 3 points) with bullet points instead?
 - Yes, then proceed
 - No, then re-word into three bullet points (do not need to re-roll)
- Is all repetitive phrasing/wording removed and made more concise?
 - Yes, then proceed
 - No, then make the necessary edits
- Does the rewritten response satisfy all of the requirements of the prompt?
 - Act like a lawyer when it comes to the prompt: for every ask or request in the prompt, does the rewritten response satisfy it?
 - Yes, then proceed
 - No, then edit the rewrite further to make sure both the text and code address the prompt

When should they do a rewrite?:

For any turn:

The attempter will be required to **always** rewrite the **preferred** response to achieve the goal of the prompt and fix any issues that were identified.

The only scenario where a rewrite would **not** be necessary is if one of the model responses is completely **perfect**, meeting all the requirements and specifications of the prompt. In this

case, indicate that a response rewrite is not needed. This is very rare and requires a thorough justification, in most cases, the response will require a rewrite.

Specification List

1. The re-written response should **fully achieve the most recent prompt**.
2. The re-written response should **address any secondary objectives implied by the prompt**.
3. The re-written response should **correct all errors (major or minor)**
4. The re-written response should be **coherent and logically connected to the prior conversation**.
5. The re-written response should **fulfill all the dimension for "No Issues" in the rating rubrics**.
6. The rewrite should have **input validation** if the code accepts input from the user
7. The re-written response should **follow the formatting and styling specifications**.
8. The code of the re-written response should **contain enough comments**.
9. **The code of the re-written response MUST run properly and be optimal and efficient. Failure to properly test code results in removal from the project.**

In short, they are doing a rewrite step, this response is required to be perfect! All identified issues must be fixed, regardless of their severity.

How to score a task

This rubric is designed to help you accurately score tasks based on the quality of both the prompt and the model response ratings.

- In this step, you are grading the quality of the **attempter** who's task you are reviewing.
- These are just suggestions to help guide you, if any other issues in the task that are not mentioned below come up, give the task the score that you think it deserves

1. Spam

- The task is irrelevant, nonsensical, or entirely inappropriate. The response doesn't relate to the task category, language, or requirements at all.
- The prompt uses code or questions from the internet
- **Task needs to be SBQd**

2. The Prompt or Rewrite is Bad

- The prompt is incorrect, misaligned, or does not match the project requirements. The response may not follow key requirements due to the poor prompt setup. In this case the task needs to be redone
- The rewrite provides a flawed solution and doesn't meet all of the requirements in the [rewrite requirements](#) section
- **Task needs to be SBQd or requires major fixes**

3. 3+ Incorrect Ratings /

- There are three or more incorrect ratings for instruction following, accuracy, efficiency, etc. The response contains significant issues in multiple areas that the attempter did not catch, or the attempter flagged issues that were not true.

4. 2 Incorrect Ratings ✓

- The task is almost perfect but contains one or two incorrect ratings in key areas (instruction following, accuracy, efficiency, etc.).

5. Perfect Task 🎉🏆

- The task is perfect, with a good prompt and all ratings accurately reflecting the response.

Some helpful rubrics:

Rewrite Response Rubric

Field	1-2 (Fail)	3 (Okay)	4-5 (Good/ Perfect)	Additional Notes
Accuracy	<ul style="list-style-type: none"> Major Factual Errors: Response has 1+ major factual errors or misleading points. Minor Factual Errors: Response has 2+ minor factual errors. 	<ul style="list-style-type: none"> Contains only 1 minor factual error or misleading statement. 	<ul style="list-style-type: none"> No factual errors or misleading statements. 	<ul style="list-style-type: none"> A major error involves incorrect/misleading data central to the request. A minor error is near the subject matter but doesn't affect the main point.
Instruction Following / Response Fulfillment	<ul style="list-style-type: none"> Main Goal Miss: Does not achieve the main goal or make progress in the conversation. Explicit Instruction Miss: Misses 1+ key instructions. Not Fulfilled: Fails to answer the question. 	<ul style="list-style-type: none"> All explicit instructions are followed. Secondary Objectives Miss: Some secondary objectives not addressed. Subjective Instruction Miss: Subjectively misses some parts. 	<ul style="list-style-type: none"> Fully achieves the main goal or makes clear progress. Follows all instructions and fully answers the question. 	<ul style="list-style-type: none"> Rule of thumb: for word count, ±10% is acceptable. Example: If a question asks for a historical event year, even if implied, the year should be provided.
Unnecessary Greetings / Pleasantries	<ul style="list-style-type: none"> Contains greetings/pleasantries such as "Sure, I'd love to help." 	N/A	<ul style="list-style-type: none"> No unnecessary greetings or pleasantries. 	<ul style="list-style-type: none"> Only flag unnecessary phrases at the beginning or end of the response. Phrases like "Here is..." are not considered pleasantries.
Depth / Nuance	<ul style="list-style-type: none"> Little to No Detail: Response is superficial and lacks meaningful Too Much Detail: 	<ul style="list-style-type: none"> Has enough detail but may need more depth or nuance. 	<ul style="list-style-type: none"> Balanced, insightful, and focused without going overboard. 	<ul style="list-style-type: none"> Too much detail can lead to confusion. Be clear about what points are most important.

	<p>depth or insight.</p> <ul style="list-style-type: none"> Excessive Detail: Overly complex and obscures key points. 	Slightly too much, but doesn't obscure key points.		
[Rewrite/SxS] Clearly Worse Than Model Response	<ul style="list-style-type: none"> Worse Than Original: Clearly performs worse than the original across rubric categories. Worse Than Side by Side Model: Performs worse overall. 	<ul style="list-style-type: none"> Performs similarly to the original or side-by-side comparison. 	<ul style="list-style-type: none"> Performs better overall than the original or side-by-side comparison. 	<ul style="list-style-type: none"> Don't penalize for minimal/no changes if the original response was acceptable. Compare against state-of-the-art models.



Lemur Astrologer Coding Goal-Oriented Multi-Turn (MT) Coding

[Task Specifications](#)

[Table of Contents](#)

[Project Overview](#)

[Task Overview](#)

[Task Specifications](#)

[Step 1: Goal Setting](#)

[Step 2: Prompt Writing & Tagging](#)

[Step 3: Response Evaluation](#)

[Step 4: Execution](#)

[Step 5: Ranking Responses](#)

[Step 6: SUPER IMPORTANT Rewriting Each Preferred Response](#)

[Step 7: Continuing the Conversation](#)

[Appendix](#)

[Prompt Examples](#)

[Response Examples](#)

[Grading Rubrics](#)

[Prompt Evaluation Rubric](#)

[Response Grading Rubric](#)

[Response Rewrite Rubric](#)

! Important Announcement:

NEW - Common gaming-related prompts are no longer allowed ("Build a tictactoe game", "Build a tetris game"). [Here is a list of prompts](#) that are commonly-seen and will be SBQ.

Read this [doc](#) for more info![Common Prompts](#)

There's an additional consideration for the Accuracy Dimension:

Input Validation. Read the changes to align your work.

Prompts must be labeled accurately by difficulty + sub-category

Code must be fully tested - bad code will result in SBQ + removal from project

All code in rewritten responses or preferred responses are REQUIRED to have sufficient code comments!

For rewrites only: Paragraphs describing processes in multiple sentences should be converted to bullet points (or numbered lists, if the process is ordered).

Project Overview

The goal of this project is to train AI models to handle a variety of coding-specific tasks. You'll do this by having a multi-turn ("MT") conversation with the model to guide it to fulfilling a specific goal. We call this "Goal-Oriented Multi-Turn." You'll start by having a specific goal

in mind and prompting the model to try to get as close to the goal as possible.

Your work will identify strengths and weaknesses in the model, help identify where/when the model gets something wrong, and help improve the model's reasoning skills with each prompt and rewritten response. In a nutshell, the **main objective** is to craft realistic prompts that cause the model to fail.

IMPORTANT!

If you've already worked on this project, please review these **Recent Changes**:

10/31/24 - [NEW Rating Criteria: Input Validation.](#)

10/31/24 - [Changes to rewrite guidelines](#) - now required at every turn

10/31/24 - Presentation failures cannot be the only mistake/issue

[Stylistic Guidance & Changes when Rewriting Responses](#)

[Execution Instructions](#)

Task Overview

Here is a high-level overview of the task:

Step 1: Goal Setting

Establish a clear and achievable goal tailored to a specific coding task and for a target programming language

Step 2: Prompt Writing & Tagging

Craft an effective prompt and properly tag its task category and difficulty level

Step 3: Response Evaluation

Ensure your prompt results in either one or both responses to your prompt failing, and then rating each model responses for accuracy, efficiency, and adherence to instructions

Step 4: Execution

Test the code in each of the model responses by uploading a screenshot of your code's output, list any setup and run commands, and show the input and actual result when executed

Step 5: Ranking Responses

Compare and rank the two responses, determining which one best aligns with the task criteria

Step 6: Rewriting the Preferred Response

Rewrite the *preferred* response to ensure it fully achieves the goal of the task

Step 7: Continuing the Conversation

Prompt the model again until reaching the minimum turn requirement to guide the model closer to the goal

Task Specifications

Step 1: Goal Setting

In this step, you'll **create a goal that matches the specified task category, target programming language, and difficulty level.**

Goals should

Be specific, clear, and easy to understand, leading the model to fulfilling the goal

Example Goal: "Build a pong game"

If you're unfamiliar with the task category or target programming language, it's best to skip the task and select another one - you will set your preferences in the courses.

After writing your goal, select the goal type and sub-category. There are **2 GOAL TYPES**, with sub-categories as follows:

Error Correction

Fixing Model Errors

Identifies and corrects specific errors in model responses.

Rectifies inaccuracies, bugs, or other mistakes in code or text.

Handling Ambiguous Requests

Trains the model to interpret and clarify unclear instructions.

Encourages reasonable assumptions or requests for clarification to fulfill tasks properly.

Multi-Step

Building Progressively

Guides model through a series of steps to achieve a complex task. Encourages structured and logical problem-solving through each step.

Deep Dive

Explores a concept in greater depth. Encourages nuanced responses and deeper understanding.

Refining Requests

Refines model's responses with more precise requests.

Adjusts and narrows prompts to improve quality and accuracy.

Step 2: Prompt Writing & Tagging

In this step, you will write a prompt that **aligns with the chosen category and difficulty level**, **guiding the model toward achieving the task's goal**. Additionally, you will **label the prompt** to clearly indicate the selected **category** and **difficulty level**. If the prompt you write covers more than one category, the **main request** must be of the **chosen category**.

Prompt Creation:

Write the initial prompt based on the task category and programming language.

Ensure each prompt is clear, specific, and challenging to the model.

Design the prompts to reflect the set difficulty level.

Use only English for all prompts.

Avoid low-effort prompts, such as single-sentence prompts or incredibly generic prompts without any constraints.

Special Notes:

Prompts specifically not allowed for this project include “counting” prompts (asking the model to do some kind of counting operation like “generate N of something”, “do something in N steps”, “give me X amount of lines”, “have this on line number x” - the customer does not want these as they are already aware of them, so this will be rejected if you submit a counting related prompt). **This includes referencing line numbers in the prompt, such as “debug the function on line 40”.**

Avoid asking the model to interpret scenarios or formulas from topics such as math, physics, etc - this is a **Coding project** and the model failures/deviations should be failures on coding logic and understanding. **Failures/deviations that rely on this approach will be SBQ'd.** It's alright if

a task includes formulas, as long as the model isn't expected to interpret them and penalized for not doing so correctly.

Prompts should not be related to common games (i.e. **Build a tic-tac-toe game**, **Build a tetris game**, [see this sheet for more examples](#)), these tasks are no longer being accepted and will be rejected.

Prompt Tagging:

Label your prompt with the appropriate [task category](#) (make sure to pay close attention for subsequent turns - this is a common error!)

IMPORTANT: Simplified for Educational Purposes - Example of common sub-category error:

Task Category: Code Generation/Synthesis

Turn 1 Prompt: "Make a 2d minigolf game using JS"

Turn 1 Category: Text to Code (*makes sense*)

Turn 2 Prompt: "Additionally add a stopwatch to the top right corner"

Turn 2 Category: Text to Code (~~X~~ bad, it should be **Text to Code Edits**, because we're asking for edits to existing code)

Next, tag the [difficulty level](#) (your prompt should match the difficulty, if you think it doesn't then tag it differently!)

Enhancing Testability:

Recommend code that runs on Replit (or any convenient IDE of your choice): This tool simplifies testing by setting up 80%+ of coding environments.

Add setup instructions: In subsequent prompts, ask the model to include setup steps for testing the code (e.g., using [pip](#) for Python or [npm](#) for Node.js).

Limit dependencies: Avoid prompts that require external dependencies like local databases or APIs, unless you can provide them and they are easily accessible for a reviewer - otherwise, you may be penalized.

Use a staged approach: For complex tasks like app development, break up the prompts into steps (e.g., pull requests) for easier testing. You should **not** have a laundry list of requirements in your prompts.

Final Considerations:

Make your prompts **realistic** and diverse.

The goal is to simulate real-world tasks that are both complex and testable, pushing the model's capabilities while maintaining clarity.

A Prompt Error we commonly see on this project:

Not thoroughly testing code generated by the models in response to your prompt.

Solution: Test each part of the code to ensure it runs correctly and meets your requirements.

Your prompt is not actually failing the model in either response! This will result in your task being rejected - you must ensure your prompt fails **at least one of the model responses!**

Your prompt must include all of the requirements you expect and/or desire, if one of the model responses does something that is not to your preference (i.e. it solves something iteratively instead of recursively), you cannot penalize the model for that unless you specifically required in your prompt that the solution follow a specific approach.

Tips for Avoiding this Error:

Write Clear and Testable Prompts

Create prompts that result in outputs you can easily test and review.

Use Tools Like Replit

Ask the model to generate code that runs directly on platforms like Replit.

Replit supports many libraries and can automatically set up testing environments.

Minimize External Dependencies

Avoid requiring libraries or databases that are difficult to set up unless you provide clear instructions or alternative solutions.

⚠ You must review the [Prompt Examples](#) (<- click here) below to get a good sense of what is good and bad for this project. ⚡

Task Categories

Instruction Type	Category Type	Description	Code required in prompt?
Generation/Synthesis	Code completion Text to code Text-to-SQL	Generate immediately executable code from natural language text description or existing code	No

		snippet, infilling, etc.	
Editing/Rewriting	Code summarization/compression Text to Code edits Code translation Code refactoring	Make changes or adjustments to existing code to meet new requirements or conditions, such as altering functionality or updating or enhancing features.	Yes (code must work properly)
Debugging	Debugging and troubleshooting Testing Security Review	Identify and correct errors in existing code, such as debugging, resolving syntax errors, and fixing logical mistakes.	Yes (code must contain a bug)
Documentation	Codebase documentation Comment generation Commit text generation API documentation Create example usages of this function Document this function	Generating or updating documentation related to code to help developers understand the code, its functionality, and its usage	Yes (code must work properly)

Review/Critique	Code review Log analysis (text -> text) Quality assurance	Code review & best practices is the process of reviewing and improving code quality, security, and maintainability by applying best practices, standards, and guidelines, and ensuring compliance with coding standards and regulations	Yes (code must work properly)
Code Ecosystem	IDEs or development workflows CLI (command line interface) Version control	Interacting with coding environments, such as IDEs, version control systems, or build tools, to automate tasks, integrate code, or manage development workflows	No

Difficulty Level

Here are the difficulty levels that you may encounter:

Please note that all prompts must be original and copying anything that already exists online on sources such as Leetcode and HackerRank etc. are strictly prohibited and will result in automatic removal from the project. The goal of having human contributors writing prompts is that we want to challenge the models in ways that the models can't already learn from the internet!

	Medium (Undergrad)	Hard (Graduate)	Challenger (SME)
Knowledge Required	Limited domain/algorithms knowledge or implementation context (e.g., architecture, libraries, pre-existing code)	Knowledge of standard algorithms and data structures, common libraries and concepts or additional code context may be	Expert domain knowledge or information on the specific application or deployment scenario, including substantial specific API/code context

		required to achieve an optimal solution	
Prompt Ambiguity	There is little ambiguity in the question (in the case of underspecification, good default behaviors are easy to come up with or not essential) and limited complexity of specifications (in instructions)	Medium ambiguity in the prompt (e.g., needs to come up with reasonable ad-hoc data representation or class structure without explicit guidance), multiple requirements should be satisfied, or multiple bugs should be found	Finding good solutions needs non-trivial design decisions regarding data structures, algorithms or code architecture/design patterns
Solution Complexity	The solution is easy to explain (e.g., code doesn't need comments to be understood) and to test for/debug (limited corner cases)	Involves corner cases that should be dealt with separately; an explanation of the solution requires some abstraction or decomposition of the problem into a few subproblems	Finding a solution requires solving several non-trivial subproblems or finding non-trivial bugs; a problem involves tricky corner cases, and explaining the solution to a non-expert requires adding context

Prompt Example	<p>I have a folder that contains MIDI covers of songs that I created myself. Here is what the song metadata looks like:</p> <pre>[{"songData": {"title": "", "artist": "", "album": "", "duration": ""}, "midiFilePath": ""}]</pre> <p>I want a React app that displays my MIDI music library with a play button next to each song. Clicking the play button should play the MIDI file. Use the midi-player-js library. Create some visualization based on the MIDI sequence while a song is playing.</p>	<p>A requirement dictated by management is that files are NOT allowed to use import statements that are erroneous, i.e. a package is imported but none of the methods are utilized in the file it is included inside of. The reasoning being that this can be confusing, wastes space, and is just sloppy programming. We have a quite large Python codebase and it would be unfeasible for someone to manually comb through all the files and check for this condition. Please help write a script in Python which accepts a root folder as an input and then recursively searches through all files and directories, checking the import statements used within each file and ensuring that they are not erroneous. Any files found should be recorded in a text file for manual review.</p>	<p>I work for an oil company in the Midwest, and we are drilling in some newly acquired fields that supposedly contain oil and other gas deposits five layers deep into the Earth. However, it has been brought to our attention recently that there are ore and precious metal deposits where we planned on drilling. This is problematic because our equipment for drilling for oil and gas will be damaged if it comes into contact with the precious metals. Furthermore, the precious metals will be destroyed in the process. We are trying to create a system where we can analyze the depth of where we plan to drill and determine where the metals are placed in the holes amongst the oil and gas so we can plan how we will drill without damaging anything. If we can do this programmatically, we hope to implement it with our equipment so they can drill automatically by analyzing the soil at each depth. How can we generate an artificial dataset and implement this system in Python, TensorFlow?</p>
-----------------------	---	--	---

! Prompt Examples (Highly recommended)

Please review the examples here in the: [Prompt Examples](#) which is in the appendix of these instructions

For more examples, please reference the [prompt example bank](#).

Step 3: Evaluating Responses & Continuing the Task

In this step, you will assess the model's responses based on specific criteria and provide follow-up prompts to improve its output.

3a. Producing at least one bad response

Every step should involve **at least one bad model response resulting from the prompt you provide**. If both responses are good, please write a more challenging prompt and try again.

How to quickly identify a good vs bad response: If there are any areas in the rating dimensions where a response has **at least one issue in the P0 or P1 level, see below**, we can count that response as "bad."

If a response only has issues in the presentation criteria (P2) that is **not** sufficient enough to be rated as "bad"!

3b. Rating the Responses based on each dimension

For each turn, after specifying how the two models did overall, each individual response should be rated on performance according to the following dimensions along with a brief explanation.

Quick snapshot of the dimensions (more detailed table below) and the level of priority (meaning these dimensions are most important vs less important):

Level of priority should be used to **determine which response is better with regards to your rating for the turn**. For example, if R1 has an instruction following issue while R2 has efficiency issues, R2 would be the better response.

Instruction Following: The response answers all requests in the prompt
This checks if the model understood **all** of the requests of the prompt and is addressing each one.

Make sure to check the entire implementation. If the model *tries* to address the request but *fails to*, it is an **accuracy** issue.

Priority Level: **P0 - Highest Importance**

Accuracy: All claims and code in the response is accurate and fully correct
You will need to Google the claims made by the model or execute code to check this

Priority Level: **P0 - Highest Importance**

Sub category - Input Validation/Error Handling: The response covers all meaningful input edge cases and handles them correctly.
The response validates all inputs, handling invalid data.

Priority Level: **P0 - Highest Importance**

Optimality and Efficiency: The response presents the most optimal and efficient solutions
The response is using common practices and standards

Priority Level: P1 - Second Highest Importance

Up-to-Date: The response uses only the most recent APIs, functions, or libraries available.

APIs, functions, or libraries used aren't causing compilation or runtime errors due to deprecation.

Priority Level: P1 - Second Highest Importance

Presentation: The response format follows the Style and Presentation guidelines

It should follow the presentation rubric, such as:

Enough comments in the code.

A professional tone (no pleasantries / fluff).

An answer that is concise, without repetitive statements.

Explanations that use bullet points.

Priority Level: P2 - Third Highest Importance

Each dimension rating should include a brief explanation as to why the given rating was chosen (e.g., why "Major Issues" was selected for Accuracy). The explanations/justifications you write for dimension ratings should be **specific** and **clear**. If there are any bugs, or there are issues, specify **what the issues are** and **what the causes of the issues are**.

Avoid generalizations like "No Issue", "N/A", "This has an error", etc. If there are no issues, give a very brief explanation as to why.

Examples of Bad Justifications:

"There is a syntax error": ✗ This is only half of the story - you should briefly explain what's causing the syntax error.

"The response doesn't satisfy all of the prompt requirements": ✗ This is very vague. If you find that the response doesn't satisfy all of the requirements, you should be specifying which requirements it doesn't satisfy and why it doesn't satisfy them.

Response Rating Rubric

Dimensions	NA (0)	No Issue (1)	Minor Issue (2)	Major Issue (3)
Instruction following	<i>This dimension cannot be NA.</i>	The response meets the main request and all constraints, showing a strong understanding of the prompt, even if there are minor implementation errors. It handles any ambiguities well and stays within the specified requirements.	The response fulfills the primary request but does not entirely adhere to all the constraints. The response could have better handled the ambiguity of the prompt. Common errors: Fails some but not ALL constraints	The response fails to fulfill the primary request OR fulfills the primary request but does not adhere to <i>any</i> constraints. Common errors: Fails to do the primary request

Accuracy and [NEW] Input Validation / Error Handling	Can only be NA if the response contains no code or factual claims, and does not rely on prior context.	<p>The code runs error-free, produces the correct output, and follows best practices. All text and comments are accurate, and the response is contextually appropriate with any previous errors fixed.</p> <p>—</p> <p>All meaningful edge cases are covered. The response includes thorough validation for expected inputs and error handling for invalid data, ensuring robustness and resilience to common input errors.</p> <p>Example:</p> <p>The function <code>calculate_discount(price, discount_percent)</code> validates all inputs. It checks that <code>price</code> and <code>discount_percent</code> are positive numbers, ensures <code>discount_percent</code> does not exceed 100%, and returns a clear error message if values are out of expected ranges or of incorrect types (e.g., strings).</p>	<p>The code runs but has minor warnings or low-risk security issues. The content is mostly accurate, but some statements are unclear or make unproven claims.</p> <p>Previous errors remain but don't affect the current response.</p> <p>—</p> <p>Some edge cases are covered. The response handles most expected inputs but misses certain edge cases, which could lead to potential errors or exceptions under specific conditions.</p> <p>The function <code>calculate_discount(price, discount_percent)</code> includes basic validation, such as checking that <code>price</code> and <code>discount_percent</code> are positive numbers. However, it lacks checks for certain edge cases, such as ensuring <code>discount_percent</code> does not exceed 100% or verifying that inputs are numeric.</p>	<p>The code doesn't run due to logic errors, produces incorrect output, or has major security flaws. The response includes false claims, lacks context, and previous errors were not fixed, making the issues worse.</p> <p>—</p> <p>No edge cases are covered. The response lacks validation for all inputs, making it vulnerable to errors when faced with unexpected or invalid data inputs.</p> <p>The function <code>calculate_discount(price, discount_percentage)</code> performs no validation on its inputs, assuming <code>price</code> and <code>discount_percentage</code> are always valid and within expected ranges. This could cause runtime errors or incorrect results if given invalid inputs, such as negative numbers, <code>discount_percentage</code> over 100%, or non-numeric types, making the function unreliable.</p>
Optimality and Efficiency	Can only be NA if the response contains no code using functions or statements aside from	The code is well-optimized, handles edge cases, and follows standard best practices. If top performance isn't required, it still performs efficiently without adding unnecessary complexity.	The code performs well but could use minor optimizations. It generally follows best practices but may not scale for large datasets.	The code exhibits severe performance and efficiency issues. The code does not adhere to common practices and standards.

	the assignment			
Presentat ion WHEN REWRITI NG, YOU MUST FIX ALL PRESEN TATION ISSUES	<i>This dimension cannot be NA.</i>	<p>The code is well-documented, with clear comments and explanations for any modifications. Code included in the prompt that did not originally have comments should have comments if included in the response. The response is concise, well-organized, and uses readable variable and function names. Complex processes are broken down with bullets, and Markdown is correctly formatted with clear hierarchies.</p> <p>Formatting is neat, with triple backticks for code blocks, and proper use of bold and italics for emphasis. White space and line breaks improve readability, and tables are correctly aligned. Functions are modular and follow standard patterns, such as using <code>if __name__ == "__main__":</code> blocks for structure. There are no redundant solutions provided for the same problem.</p>	<p>The documentation is generally clear but could use more detail. There are minor language errors that don't affect readability, and formatting could be improved for clarity. Variable and function names are understandable, but some structural changes—like adding bullets or logical sections—would help. Functions are present but may need more modularity, and some explanations are missing, making the code harder to follow in parts.</p> <p>Common Errors</p> <ul style="list-style-type: none"> Uses backticks inconsistently Uses camelCase and snake_case inconsistently 	<p>The documentation is missing or inadequate, or lacking code comments entirely, making the code hard to understand. The response is poorly formatted and lacks structure, with unclear variable and function names. The logic is disorganized, and there are no explanations for key decisions, making it difficult to follow, integrate, or reuse. Programming language tags are also missing.</p>
Up-to-Dat e	NA (0)	Up-to-Date	Out-of-Date	
	The code does not call on any libraries or	<p>The code uses the most fresh API, libraries, or functions available to solve problems efficiently. The code uses a maintained library or function which is an older version that still</p>	<p>The code uses a deprecated API, library or function, causing a runtime or compile-time error.</p>	

functions.	works (even if it is less efficient).	
------------	---------------------------------------	--

Step 4: Execution (YOU MUST TEST YOUR CODE!)

For any response that changes the executable code, you will also be required to execute the code in the response. Note: *This doesn't apply to responses that have only added/modified code comments.*

Below are the steps that you will encounter for running code!

Execution Guidelines	
<p>Program ming Language :</p>	<p>What additional languages were involved in the response? (Select all that apply) <small> ⓘ </small></p> <p>Exclude the programming language assigned to the task! Applies to both code and code-related discussion in the model response.</p> <p>Hover over the hint tooltip for the list of possible selections</p> <p>Choice Paths:</p> <div style="border: 1px solid #ccc; padding: 5px; width: fit-content;"><input type="text" value="html"/> <input type="checkbox"/> HTML/CSS</div>



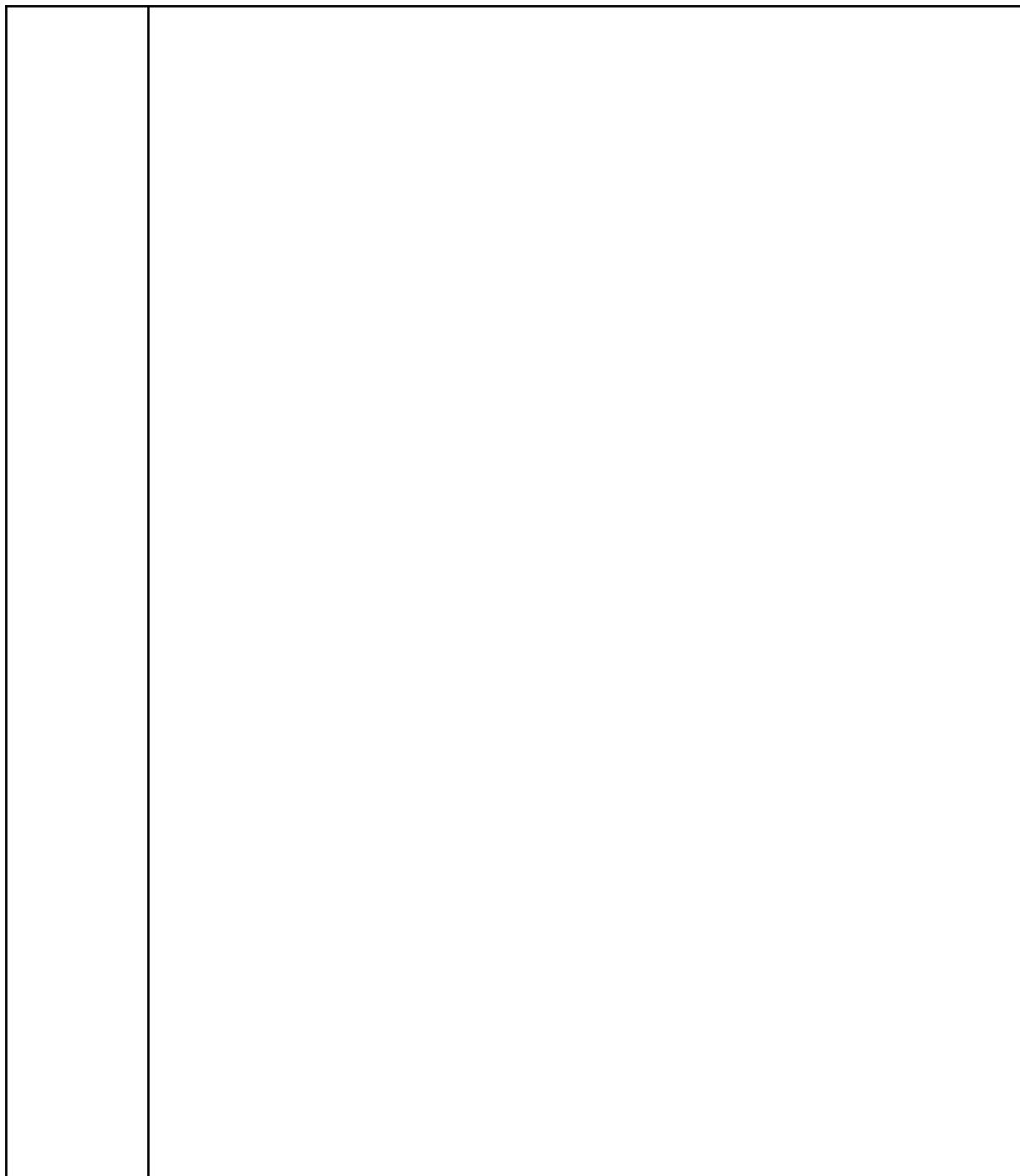




Executable Code:	<p>Does the model response edit or add any executable code? *</p> <p>Mark "Yes" if there is any new or altered code present in the response, regardless it is an entire program or just a code snippet. Does not apply to new comments!</p> <p><input type="radio"/> Yes <input type="radio"/> No</p>







**Degree of
Execution**

:

Rate the degree of execution of the code *

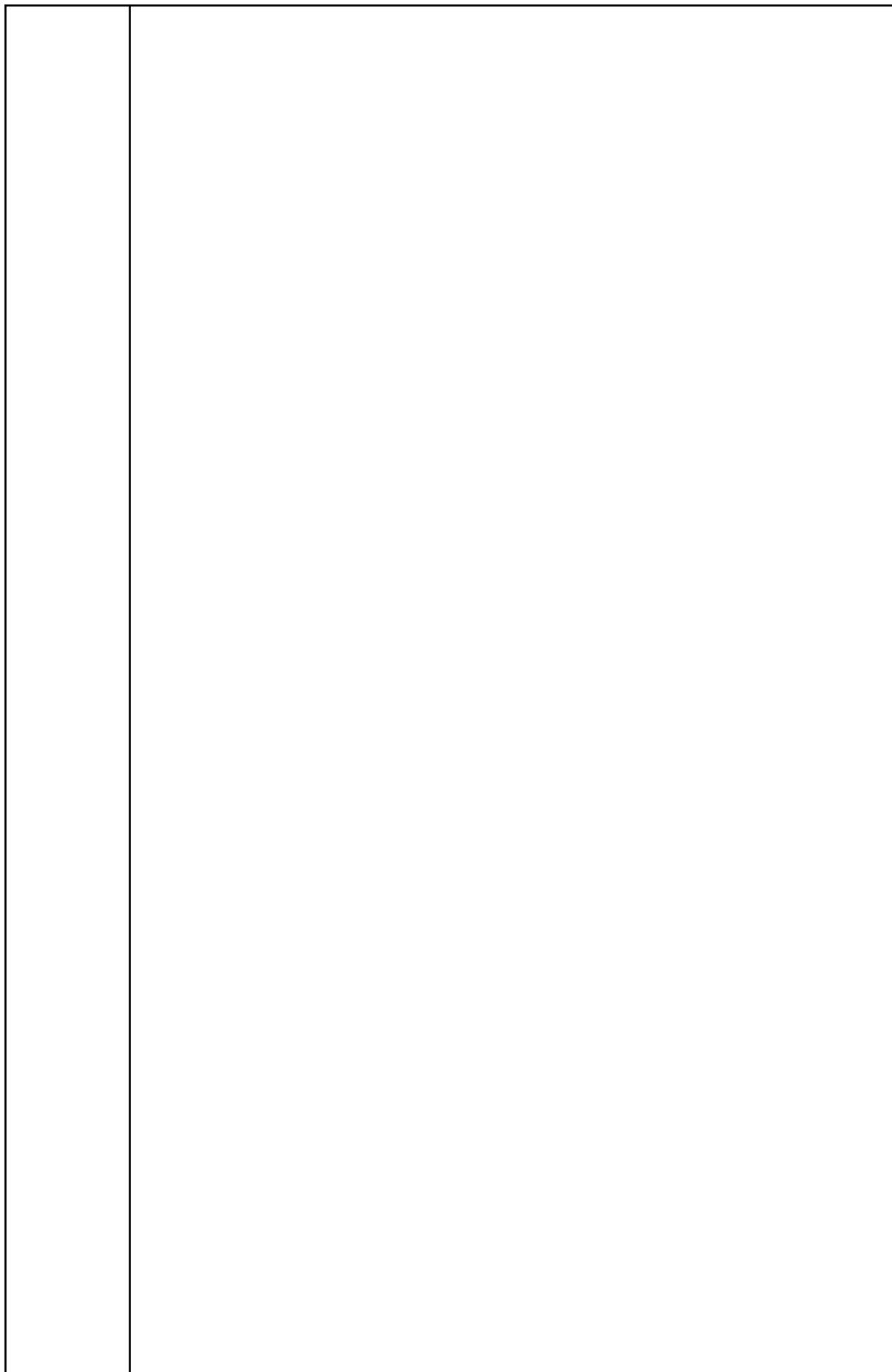
(Only if there is executable code in the response)

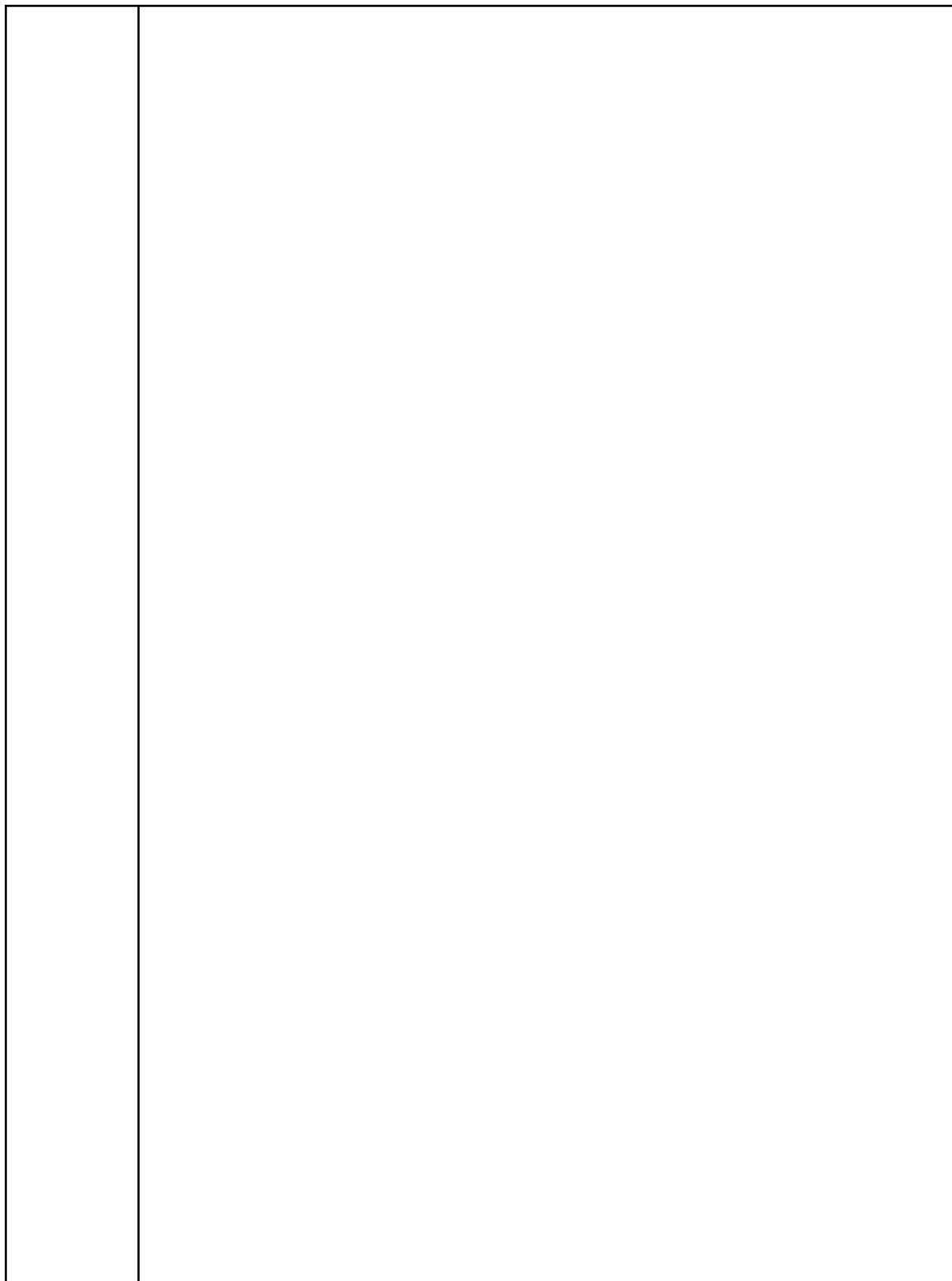
- Template ⓘ
- Partial Update ⓘ
- Function Update ⓘ
- Out-of-the-Box ⓘ
- NA











**Installatio
n
Comman
ds:**

What installation commands are necessary to test the code? *

Please separate each individual command with a comma (,) and write 'N/A' if there are no extra install commands required.

ex: pip install tkinter, sudo apt get

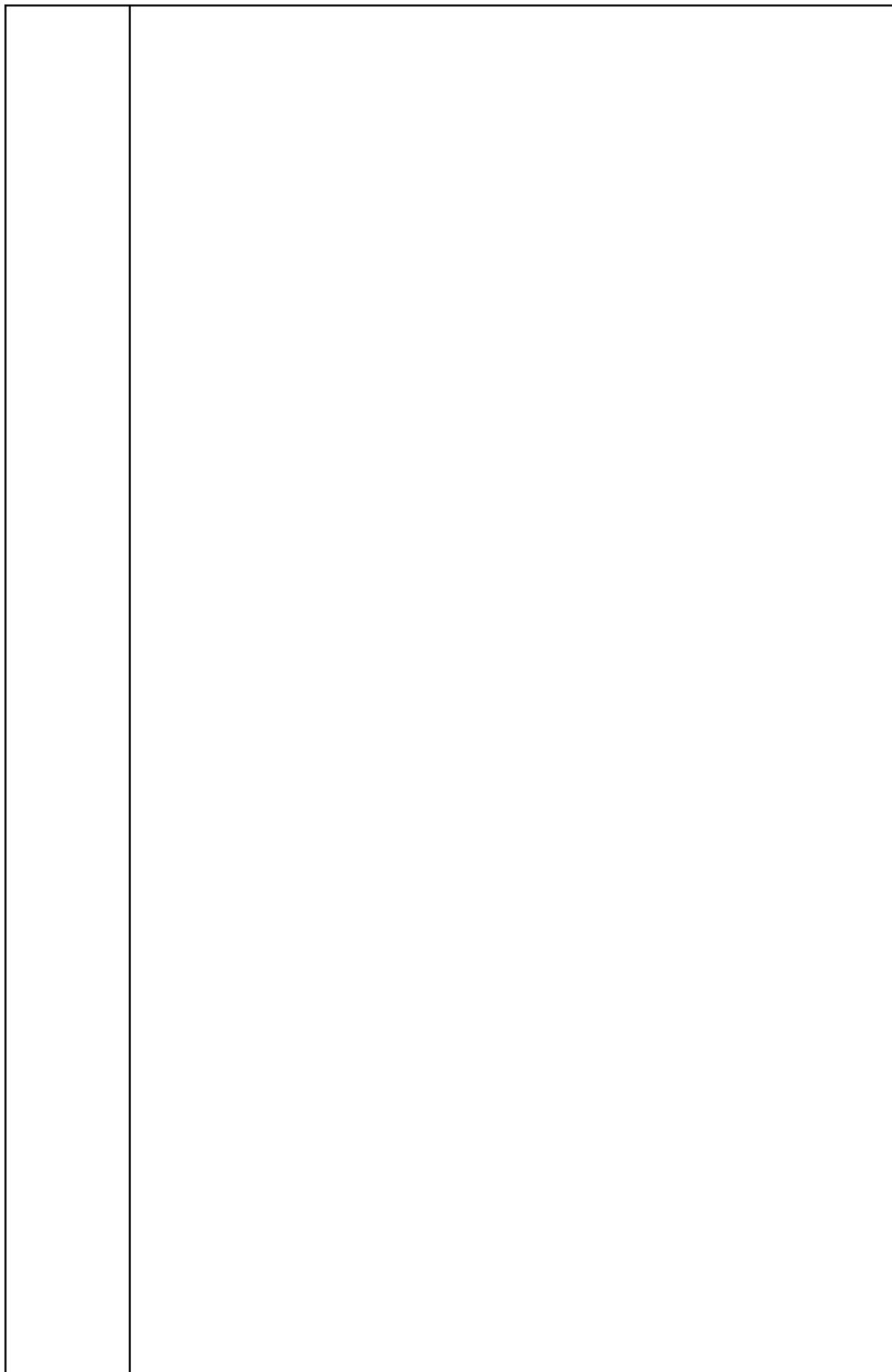






Run Command ds:	<p>What commands are necessary to run the code? Provide a comma separated list if there are multiple commands. *</p> <p>ex: unicorn main:app --host 0.0.0.0 --port 8000, python app.py</p> <div style="border: 1px solid #ccc; height: 100px; width: 100%;"></div>





Output Expectati on:	<p>Is the output of the code as expected? *</p> <div style="border: 1px solid #ccc; padding: 5px; width: fit-content;"><input checked="" type="radio"/> Yes <input type="radio"/> No</div>



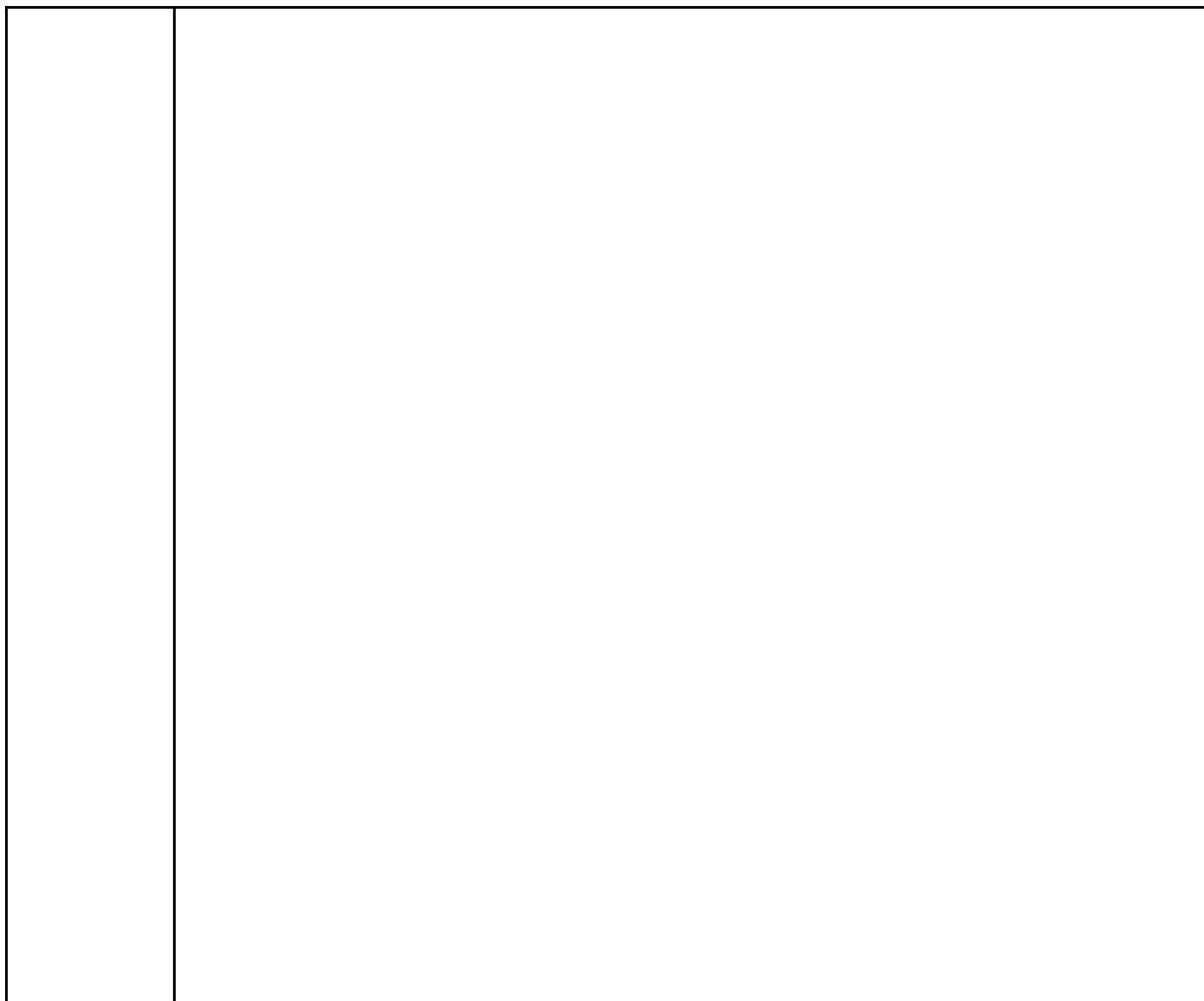


Output MIME Type:	<p>What is the mime type of the output? *</p> <div style="border: 1px solid #ccc; padding: 5px; width: fit-content;"><input type="radio"/> text/plain <input type="radio"/> Other (specify below)</div>



Record any errors:	<p>Did the code produce an error? *</p> <p><input checked="" type="radio"/> Yes</p> <p><input type="radio"/> No</p> <p>Code Execution: Error</p> <p>Please paste the error from the code execution.</p> <div style="border: 1px solid #ccc; height: 100px; width: 100%;"></div>





Upload Screenshot/ Recordin g of Output

Upload Screenshot/Screen Recording of Output (required)

Upload a screenshot of your code output (including the execution call). If the output of the code is not static, upload a screen recording instead.



Drag and drop here or click to upload your data





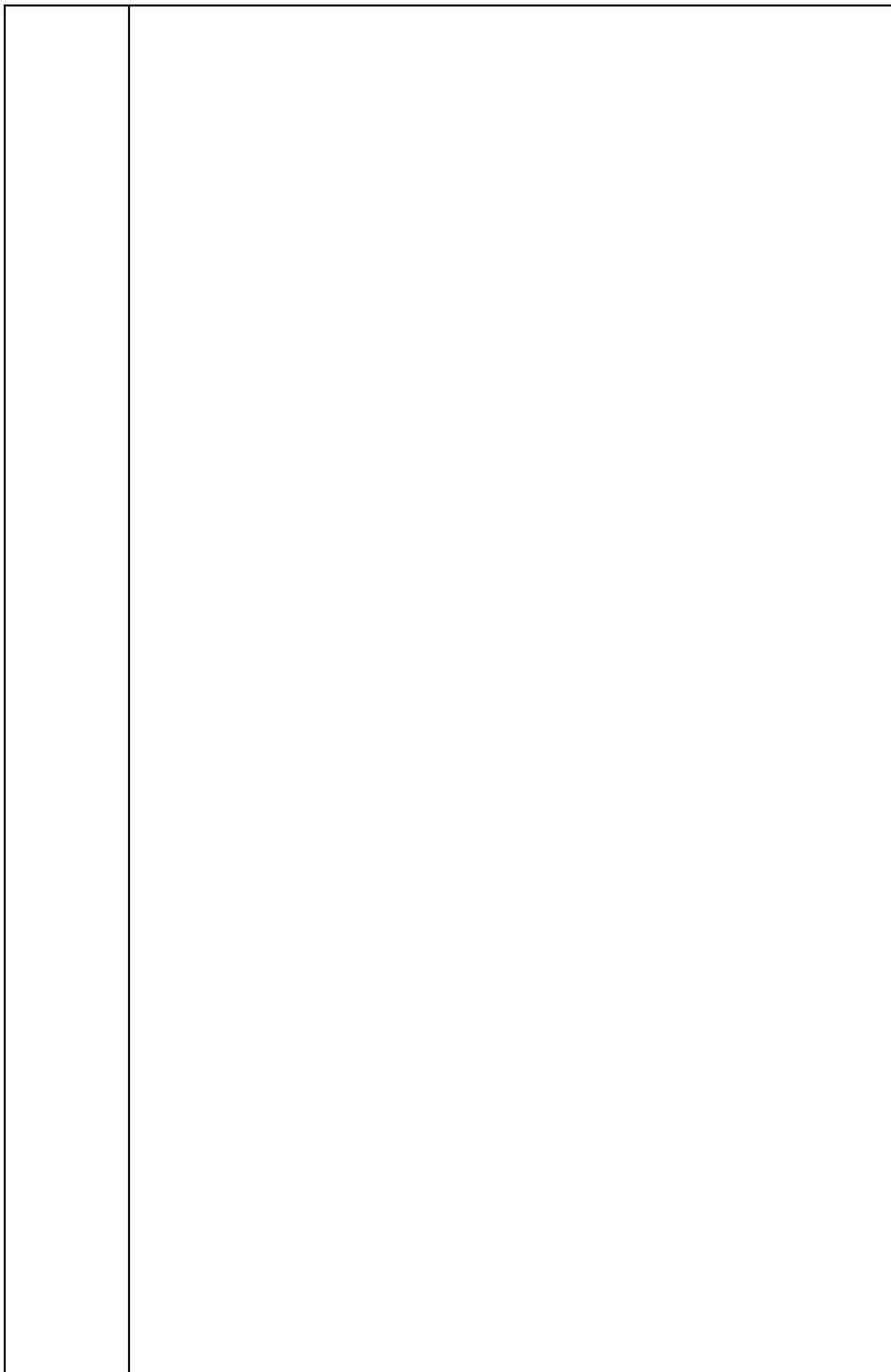








Upload Files	<p>Is your code a single file or multiple files? *</p> <p>Ignore any third party dependencies and libraries</p> <p><input checked="" type="radio"/> Single file <input type="radio"/> Multiple files</p> <p>File Name *</p> <p>Write the name of the single file containing your code. This should be the same as the one used in your run commands. (ex: main.py)</p> <p>main.py</p> <p>Paste Code (required) *</p> <p>Paste your entire codefile here. DO NOT USE IF YOU HAVE MULTIPLE FILES (see previous question)</p> <pre>def read_file(file_path): try: with open(file_path, 'r') as file: data = file.read()</pre> <p>Upload Code (required)</p> <p>Upload a zip file containing the code you used to test the response.</p> <p>Remember: The code you upload should be executable (even if it throws an error), and should correspond to the Execution Instructions you provided above.</p> <p>Drag and drop here or click to upload your data</p>







Output:	



<p>Output</p> <p>Expected:</p>	





Output	
MIME	
Type:	



JS	<p>JSON Template</p> <p>Please leave blank!</p> <p>DO NOT FILL!</p>

--	--

Step 5: Ranking Responses

In this step you will compare the two model responses and rate which one is better on a scale of 1-8, indicating whether Response 1 or Response 2 performed better. This should match how you rated the models on each dimension (i.e. if Response 2 performed worse on every dimension, then the score should be a 1).

Ranking Justification Guidelines:

Base your ranking on the rating dimensions in the rubric, listed by importance.

Provide helpful and factual feedback without unnecessary comments or conversation.

Clear Documentation: Use docstrings for code explanations and modify code directly in your response, avoiding extra copy-pasting.

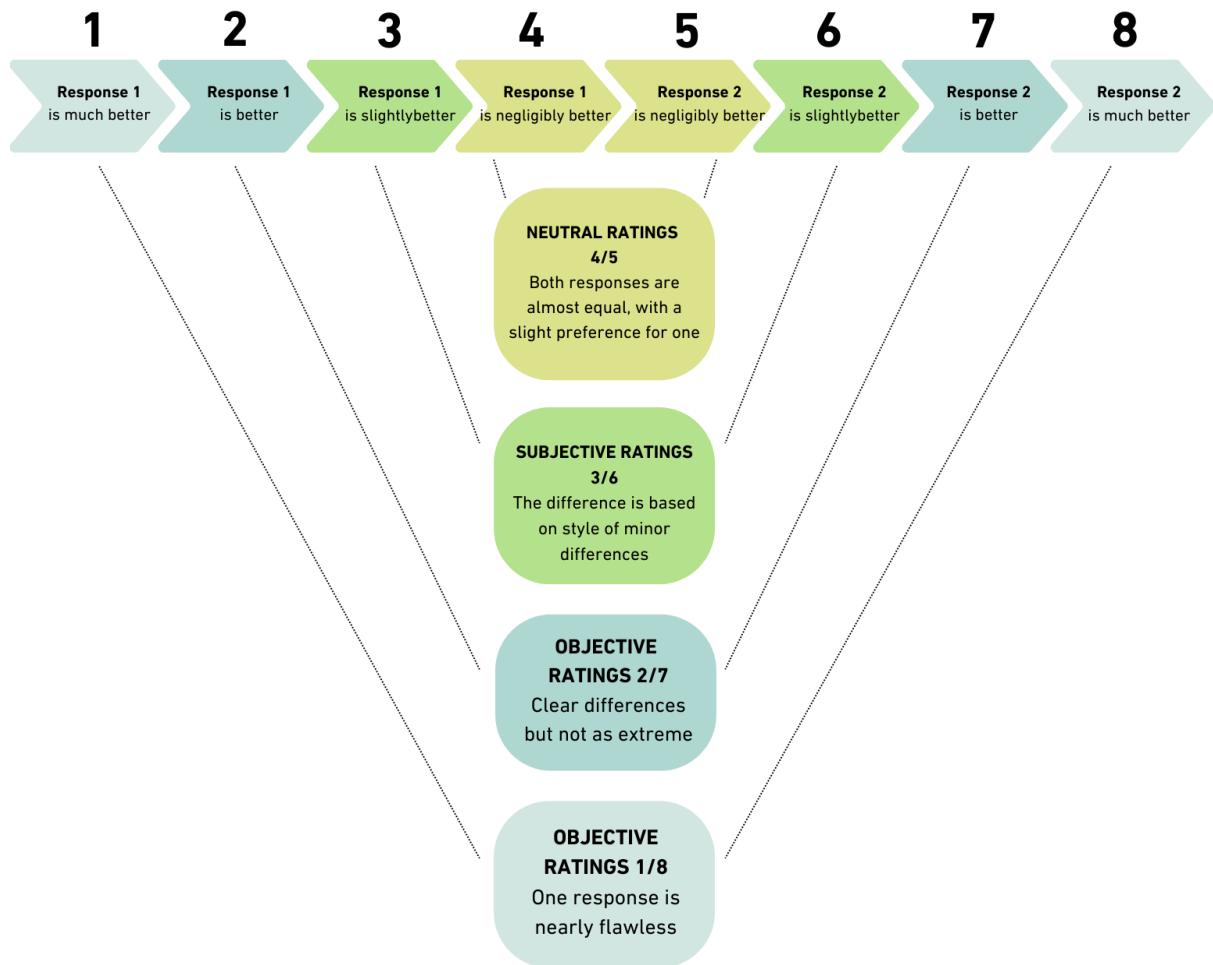
Readability and Structure: Ensure code is easy to read, with a focus on structured and modular design (e.g., use functions and `if __name__ == "__main__":` blocks).

Use bullet points or numbered lists for clear and organized explanations.

Choose the response that is **easier to understand and integrate** with minimal user effort.

Rating Scale

Scores from 1-8 will demonstrate either a strong or neutral preference between the two responses



Use the priorities mentioned here to determine which response is better based on the severity of the issues.

Step 6: Rewriting Each Preferred Response

In this step, once you have selected the better response from each turn, you will be required to **always** rewrite the **preferred** response to achieve the goal of the prompt and fix any issues that were identified.

The only scenario where a rewrite would **not** be necessary is if one of the model responses is completely **perfect**, meeting all the requirements and specifications of the prompt. In this case, indicate that a **response rewrite is not needed**. This is very rare and requires a **thorough justification**, in most cases, the response will require a rewrite.

Specification List

The re-written response should **fully achieve the most recent prompt**.

The re-written response should **address any secondary objectives implied by the prompt**.

The re-written response should **correct all errors (major or minor)**

The re-written response should **be coherent and logically connected to the prior conversation**.

The re-written response should **fulfill all the dimensions for "No Issues" in the rating rubrics**.

The re-written response should **follow the formatting specifications**.
The code of the re-written response should **contain enough comments**.
The code of the re-written response MUST run properly and be optimal and efficient. Failure to properly test code results in removal from the project.

Formatting/Presentation Requirements

Key terms **should** be highlighted in bold, whereas titles, articles, etc. are italicized.

Remove pleasantries such as "Sure," "Certainly," "I can help with that," etc.
Make responses more concise, remove all fluff and unnecessary phrases (i.e. "Welcome to the world of VS Code!")

Tone should be straightforward/professional.

Code should be well-commented

Test outputs include a comment with the expected response.

Explanations should use bullet points.

Rewritten response replaces any paragraphs (especially those with at least 3 points) with bullet points instead

All repetitive phrasing/wording must be removed.

Here are examples of [formatting/presentation requirements](#)

Make sure to keep track of the changes that you made as you'll have to write them out in the next step.

Stylistic Guidance & Changes for Rewrites

Stylistic Change	Description	Bad Example	Better Example
Use descriptive variable names that make the logic evident rather than obscure it.		<pre>for (int i = 0; i < items.length; i++) { items[i].apply_discount(); }</pre>	<pre>for (Item item : items) { item.apply_discount(); }</pre>

Use of language	Repetitive phrasing: (e.g., "We"). Check if words have been guessed: The checkSecretWord function checks if an input word has been guessed. Single vs double quotes to indicate clarity: The checkIfAllGuessed function.	A guessedWords array was added to track guessed words. The checkSecretWord function displays an alert if a word has been guessed. The checkIfAllGuessed function ensures all secret words are handled and manages animations.
Items or hints for clarity	Not itemized: let or lists more see or . we import libraries like numpy, and sklearn. we create a dataset, categorical variables, split and train a model."	Itemized: Import libraries: pandas, numpy, sklearn. Create a dummy dataset with 5 attributes and 1 target variable. 1-hot coded categorical variables.. Split data into training and testing sets. Train a Random Forest model.
Consistency	Inconsistent formatting: when using one keyword but not the other. It uses matplotlib and tkinter."	Consistent formatting: Highlighting both keywords for clarity and consistency. "The code uses matplotlib and tkinter ."
Use for steps	steps to validate user input, validate user input against the database, and store data in the database." lists to each	Use a list for clarity: 1. Check the input for errors. 2. Validate user data against the database. 3. Store the data in the database.

Step 7: Continuing the Conversation

Multi-Turn Conversation Guidelines:

For goal-oriented multi-turn tasks, structure the conversation to build logically toward the final goal:

Build Progressively

Each prompt should add value and guide the model closer to completing the task.

Encourage Depthng

As the conversation advances, focus on enhancing the model's understanding and pushing for more detailed responses.

Evaluate and Iterate

Assess each response, rewrite as needed, and ensure every turn contributes meaningfully to the final solution.

If the required number of turns has been met, you can end the task, even if the goal hasn't been fully achieved. This approach maintains focus and prevents unnecessary iterations.

PRO TIPS!

For Subsequent Turns, you may choose any category or difficulty but they must be aligned with your prompt

THE OVERALL IDEA of this project is to test the model's ability to recall context from prior turns

Subsequent Prompts **should not be disconnected from previous prompts**

It's easier to go back and modify the goal, if the model is doing well in recalling and building on previous context

You don't need to "fully complete" a goal, as it is simply an alignment tool

EXAMPLE: If the goal was to build a game, having parts of that game built is fine

Appendix

Grading Rubrics

Prompt Evaluation Rubric

This is the rubric we use to measure the quality of your prompts.

Criteria	1-2 (Fail)	3 (Okay)	4-5 (Good/ Perfect)	Additional Notes
Prompt Instruction Type	<ul style="list-style-type: none"> The initial prompt does not follow the specified instruction type. Subsequent prompts are not labeled correctly. 	N/A	<ul style="list-style-type: none"> The prompt follows the specified instruction type. Subsequent prompts are labeled correctly. 	<ul style="list-style-type: none"> Instruction types include: <ul style="list-style-type: none"> Generation/Synthesis Editing/Rewriting Debugging Documentation Review/Critique Code Coherence
Prompt Difficulty	<ul style="list-style-type: none"> The difficulty does not match the specified level. 	N/A	<ul style="list-style-type: none"> The prompt reflects the specified difficulty level. 	<ul style="list-style-type: none"> Difficulty levels include: <ul style="list-style-type: none"> Medium (Undergrad) Hard (Graduate) Challenger (SME)
Prompt Continuity	<ul style="list-style-type: none"> The prompt does not logically follow the conversation. It fails to achieve the conversation goals. 	<ul style="list-style-type: none"> The prompt may have slight shifts but still furthers the conversation. 	<ul style="list-style-type: none"> The prompt logically continues the conversation. Topic shifts are allowed if they help achieve the conversation goals. 	<ul style="list-style-type: none"> Applies if not the first turn.
Clarity & Specificity	<ul style="list-style-type: none"> The prompt is vague and unclear. Critical details are missing. 	<ul style="list-style-type: none"> The prompt is mostly clear but may be interpreted in multiple ways. 	<ul style="list-style-type: none"> The prompt is clear and specific. No assumptions are needed to answer. 	<ul style="list-style-type: none"> The prompt should be in English.
Feasibility	<ul style="list-style-type: none"> The request is impractical or impossible for an AI in a single response. Conflicting instructions are given. 	<ul style="list-style-type: none"> The request is feasible but requires effort or compromises. 	<ul style="list-style-type: none"> The request is fully actionable and realistic. No conflicting instructions. 	<ul style="list-style-type: none"> Requests should not contain conflicting or impractical elements. Example: Asking for a complex algorithm in one step.

Response Grading Rubric

This is the rubric we use to grade your response ratings.

Field	1-2 (Fail)	3 (Okay)	4-5 (Good/ Perfect)
-------	------------	----------	---------------------

Model Responses Rating	<ul style="list-style-type: none"> Major Rating Disagreement: (e.g., marked "Good" vs. "Bad") and poor justification. 	N/A	Rating aligns with expectations.
Prompt Adherence and Instruction following	<p>Major Rating Disagreement: The prompt adherence or instruction following is incorrect (e.g., marked "N/A" instead of "Major Issues").</p> <p>Minor Rating Disagreement (2+): Minor disagreements occur on 2 or more occasions.</p>	<p>Minor Rating Disagreement (1): Minor variance on 1 occasion, and justification does not fully support this variance.</p>	Rating is correct, or "N/A" is appropriately used when the response cannot be assessed.
Correctness and Accuracy [input validation]	<p>Major Rating Disagreement: Significant issues with correctness (e.g., marked "Major Issues") or minor disagreements on 2+ occasions.</p> <p>No edge cases are covered. The response lacks validation for all inputs, making it vulnerable to errors when faced with unexpected or invalid data inputs.</p>	<p>Minor Rating Disagreement (1): Slight variance on 1 occasion, justification lacks support.</p> <p>Some edge cases are covered. The response handles most expected inputs but misses certain edge cases, which could lead to potential errors or exceptions under specific conditions.</p>	Rating matches expectations, or "N/A" is used when applicable. All meaningful edge cases are covered. The response includes thorough validation for expected inputs and error handling for invalid data, ensuring robustness and resilience to common input errors.
Performance, Optimality, and Efficiency	<p>Egregious Rating Disagreement: Issues with performance were missed (e.g., not marking "Major Issues"), with flawed justification.</p>	<p>Major Issues: Minor performance issues were missed, but the justification somewhat supports the rating.</p>	No major issues missed in performance.
Readability, Documentation, and Presentation	<p>Egregious Rating Disagreement: Significant documentation or readability issues were not marked.</p>	<p>Major Issues: Some readability or presentation issues were missed, but the justification is somewhat valid.</p>	No major readability or documentation issues were missed.

Up-to-Date	Egregious Rating Disagreement: Incorrectly marked as "N/A" or "Up-to-Date" when it should have been "Out-of-Date."	N/A	Rating is accurate or correctly identifies outdated information.
MT Conversation Quality	<p>Task with multiple turns does not consider holistic, conversation-level factors. Conversation is incoherent or turns don't follow logically. Task doesn't require multi-turn structure (MT) or lacks context from prior turns.</p>	<p>Task is somewhat long-winded but follows logically, considering conversational factors. Conversation accomplishes its goal but could be completed in one turn.</p>	<p>Task fulfills the main goal of the conversation. Turns follow logically and handle topic shifts effectively. Each turn builds on the prior, maintaining consistency and context. New turns incorporate older turns for elaboration and coherence. Multi-turn context</p>

				is crucial for generating appropriate responses.
--	--	--	--	--

Response Rewrite Rubric (if applicable)

This is the rubric we use to measure the quality of your rewritten responses

Field	1-2 (Fail)	3 (Okay)	4-5 (Good/Perfect)	Additional Notes
Accuracy	<ul style="list-style-type: none"> Major Factual Errors: Response has 1+ major factual errors or misleading points. Minor Factual Errors: Response has 2+ minor factual errors. 	<ul style="list-style-type: none"> Contains only 1 minor factual error or misleading statement. 	<ul style="list-style-type: none"> No factual errors or misleading statements. 	<ul style="list-style-type: none"> A major error involves incorrect/misleading data central to the request. A minor error is near the subject matter but doesn't affect the main point.
Instruction Following / Response Fulfillment	<ul style="list-style-type: none"> Main Goal Miss: Does not achieve the main goal or make progress in the conversation. Explicit Instruction Miss: Misses 1+ key instructions. Not Fulfilled: Fails to answer the question. 	<ul style="list-style-type: none"> All explicit instructions are followed. Secondary Objectives Miss: Some secondary objectives not addressed. Subjective Instruction Miss: Subjectively misses some parts. 	<ul style="list-style-type: none"> Fully achieves the main goal or makes clear progress. Follows all instructions and fully answers the question. 	<ul style="list-style-type: none"> Rule of thumb: for word count, $\pm 10\%$ is acceptable. Example: If a question asks for a historical event year, even if implied, the year should be provided.
Unnecessary Greetings / Pleasantries	<ul style="list-style-type: none"> Contains greetings/pleasantries such as "Sure, I'd love to help." 	N/A	<ul style="list-style-type: none"> No unnecessary greetings or pleasantries. 	<ul style="list-style-type: none"> Only flag unnecessary phrases at the beginning or end of the response. Phrases like "Here is..." are not

				considered pleasantries.
Depth / Nuance	<ul style="list-style-type: none"> Little to No Detail: Response is superficial and lacks meaningful depth or insight. Excessive Detail: Overly complex and obscures key points. 	<ul style="list-style-type: none"> Has enough detail but may need more depth or nuance. Too Much Detail: Slightly too much, but doesn't obscure key points. 	<ul style="list-style-type: none"> Balanced, insightful, and focused without going overboard. 	<ul style="list-style-type: none"> Too much detail can lead to confusion. Be clear about what points are most important.
[Rewrite/SxS] Clearly Worse Than Model Response	<ul style="list-style-type: none"> Worse Than Original: Clearly performs worse than the original across rubric categories. Worse Than Side by Side Model: Performs worse overall. 	<ul style="list-style-type: none"> Performs similarly to the original or side-by-side comparison. 	<ul style="list-style-type: none"> Performs better overall than the original or side-by-side comparison. 	<ul style="list-style-type: none"> Don't penalize for minimal/no changes if the original response was acceptable. Compare against state-of-the-art models.

Prompt Examples

Prompt Category	Good Examples (specific)	Bad Examples (vague)
Code generation	<p>I am a software engineer analyzing our CDN service performance. I want to write a Bash script to extract the related data from the cache statistics report, count the hit ratio, and filter out the cache items' ID with a hit ratio of less than 0.85.</p> <ol style="list-style-type: none"> The cache statistics report is a CSV file. Its header row contains the Item ID, Name, Viewer Location, Time, Request Count, Hit Count, Miss Count, and Error Count. And the Item ID is unique. The input to the script is the name of the report. Before calculating, first check whether the count is valid. If any Request Count or Hit Count is missing or Hit Count is greater than Request Count, the script 	<p>[Asking for a Sudoku solver is too generic of a prompt, and it lacks complexity.]</p> <p>Hello, I want you to create an automatic sudoku solver in Bash that will take a .txt file that includes 81 numbers from the top row to the bottom with empty cells replaced by zeros in the first line.</p>

	<p>will log the missing or wrong data row and skip it. The hit ratio formula is Hit Count divided by Request Count. The hit ratio should be rounded to two decimal places.</p> <p>4. The output file is a text file containing unique item IDs and their hit ratio, sorted by ratio in descending order.</p>	
--	--	--

Code debugging

As a friend group, we really love to play DnD and other roleplaying games. For the upcoming weekend, I wanted to surprise my friends by developing a multiplayer text-based adventure game with C++. In this game players can explore a fantasy world, interact with characters, and solve puzzles. The game uses a complex system of pointers and dynamic memory allocation to manage player actions, game state, and world events. However, there are several bugs in the code that leads me to crashes, memory leaks, and incorrect game states. I cannot test the game because it outputs nothing. Could you help me to identify the bugs in the code and fix them?

```
#include <iostream>
#include <string>
#include <map>
using namespace std;
class GameEvent {
public:
    string description;
    bool completed;
    GameEvent(string desc) : description(desc),
    completed(false) {}
};
class Character {
public:
    string name;
    int health;
    Character(string charName)
: name(charName), health(100)
    {}
    void takeDamage(int amount) {
        health -= amount;
        if (health < 0) health
            = 0;
    }
};
class Player : public
Character {
public:
    vector<GameEvent*> events;
    Player(string name) :
    Character(name) {}
```

[This is bad because it is too vague for a debugging prompt, being “concerned” over code is not enough to warrant a debugging prompt, a debugging prompt points to a very specific issue like an error or crash; additionally the prompt veers off into being a code editing prompt in the second part which is bad]

Hello, I'm currently working on a JavaScript project which is aimed to actively monitor system performance statistics, as well as network activity and performance. It will generate logs every hour and properly handle data so that memory doesn't become an issue. Currently though, I'm suspicious of the log file handling capabilities of the code, and the asynchronous os.cpuUsage has me concerned. Please look over my code and confirm or deny if these concerns are warranted, make any changes you feel would improve the code. Also, I intend to include an active alert system for critical performance issues, as well as a text-based report which will compare logs and display a comparison for the user. Please implement these. Also keep in mind this is a multi-person project, please do not include OS specific libraries and functions.

```
// Import the os-utils library
const os = require('os-utils');
const fs = require('fs');
// Function to monitor network
// performance
function monitorNetworkPerformance() {
    let networkData = {
        downloadSpeed: 0,
        uploadSpeed: 0,
        latency: 0,
        packetLoss: 0
    };
    // Use the navigator.connection
    // API to get network performance data
    if (navigator.connection) {
        networkData.downloadSpeed =
        navigator.connection.downlink; // in
        // Mbps
        networkData.uploadSpeed =
        navigator.connection.uplink; // in
        // Mbps
        networkData.latency =
        navigator.connection.rtt; // in ms
```

<pre> void addEvent(GameEvent* event) { events.push_back(event); } void completeEvent(string eventDescription) { for (auto& event : events) { if (event->description == eventDescription && !event->completed) { event->completed = true; cout << name << " completed: " << event->description << endl; return; } } cout << "Event not found or already completed." << endl; } class Game { private: map<string, Character*> characters; Player* currentPlayer; public: Game() : currentPlayer(nullptr) {} ~Game() { for (auto& pair : characters) { delete pair.second; } } void addCharacter(string name) { characters[name] = new Character(name); } void setCurrentPlayer(string name) { if (characters.find(name) != characters.end()) { currentPlayer = static_cast<Player*>(characters[name]); } else { </pre>	<pre> networkData.packetLoss = navigator.connection.effectiveType; // effective network type } else { console.error("Network Information API not supported by this browser."); } return networkData; } // Function to monitor system performance function monitorSystemPerformance() { let systemData = { cpuUsage: 0, memoryUsage: 0 }; // Use os-utils to get system performance data os.cpuUsage(function(v) { systemData.cpuUsage = v * 100; // in % }); systemData.memoryUsage = (1 - os.freememPercentage()) * 100; // in % return systemData; } // Function to log data function logData() { const networkData = monitorNetworkPerformance(); const systemData = monitorSystemPerformance(); const logEntry = { timestamp: new Date().toISOString(), networkData: networkData, systemData: systemData }; fs.appendFile('performance_logs.json', JSON.stringify(logEntry) + '\n', (err) => { if (err) { console.error('Error writing to log file:', err); } else { console.log('Log entry recorded:', logEntry); } }); } // Function to start monitoring function startMonitoring() { setInterval(logData, 3600000); // Log data every hour } </pre>
---	---

```

        cout << "Character
not found." << endl;
    }
}
void displayStatus() {
    for (auto& pair :
characters) {
        cout <<
"Character: " << pair.first <<
        ", Health: " <<
pair.second->health << endl;
    }
}
void simulateEvent(string
description) {
    if (currentPlayer ==
nullptr) {
        cout << "No player
set." << endl;
        return;
    }
    GameEvent* newEvent =
new GameEvent(description);

currentPlayer->addEvent(newEve
nt);
    cout << "New event
added: " << description <<
        endl;
}
void resolveConflict(int
damage) {
    if (currentPlayer !=
nullptr) {

currentPlayer->takeDamage(dama
ge);
    cout <<
currentPlayer->name << " took
" << damage << " damage!" <<
        endl;
    }
}
int main() {
    Game myGame;

myGame.addCharacter("Alice");
    myGame.addCharacter("Bob");

myGame.setCurrentPlayer("Alice
");
    myGame.simulateEvent("Find
the lost treasure");

myGame.resolveConflict(20);
}

// Start monitoring
startMonitoring();

```

	<pre> myGame.setCurrentPlayer("Bob") ; myGame.simulateEvent("Save the village"); myGame.displayStatus(); return 0; } </pre>	
Code review	<p>As part of our cybersecurity module, I'm tasked with creating a Bash script that implements a basic file encryption and decryption service. The service should allow users to encrypt files using OpenSSL and decrypt them using a password. This is my initial implementation:</p> <pre> #!/bin/bash # Simple file encryption and decryption service using OpenSSL encrypt_file() { local file_to_encrypt=\$1 echo "Encrypting \$file_to_encrypt..." openssl enc -aes-256-cbc -salt -in "\$file_to_encrypt" -out "\${file_to_encrypt}.enc" -k "\$2" echo "File encrypted: \${file_to_encrypt}.enc" } decrypt_file() { local encrypted_file=\$1 echo "Decrypting \$encrypted_file..." openssl enc -aes-256-cbc -d -in "\$encrypted_file" -out "\${encrypted_file%.enc}" -k "\$2" echo "File decrypted: \${encrypted_file%.enc}" } # Check if user wants to encrypt or decrypt if [[\$1 == "encrypt" && -f \$2]]; then encrypt_file "\$2" "\$3" elif [[\$1 == "decrypt" && -f \$2]]; then decrypt_file "\$2" "\$3" else echo "Usage: \$0 [encrypt decrypt] [filename] [password]" fi </pre>	<p>[This is bad because it's too vague!]</p> <p>I'm working on the "Living the Social Life" web page. Please review the HTML and CSS code below and do the following:</p> <ul style="list-style-type: none"> Enhance the code structure and organization. Update the HTML code to use appropriate HTML5 semantic elements. <p>Here is the HTML code:</p> <pre>html placeholder</pre> <p>Here is the CSS code:</p> <pre>css placeholder</pre>

	<p>However, I have concerns about some potential security risks. For example, the current implementation accepts the encryption password in plaintext as a command-line argument, which can expose sensitive information. I also want to ensure the script handles errors more gracefully and is efficient for large file sizes. Please help review the code and suggest improvements to make it more secure and scalable. How can I avoid exposing the password and improve the script's error handling and efficiency for large files?</p>	
--	--	--

<h3>Code Editing/Writing (Modification)</h3>	<p>I wrote HTML and JavaScript codes. I aimed to print a car on the screen and move it with the "WASD" keys. However, the car doesn't look like a car. There's only a red rectangle, and I expect you to change it to a more realistic car. The point of view should be from above, so I expect to see the car from the upper perspective. I should be able to see the 4 tires, the windscreens, and the front lamps. Also, add a feature that allows the car to explode when the user presses the "space" button, and a "game over" message should be seen on the screen. There's no need to add any buttons or functionalities in the "game over" screen. Here is my code:</p> <pre> HTML: <!DOCTYPE html> <html lang="en"> <head> <meta charset="UTF-8"> <meta name="viewport" content="width=device-width, initial-scale=1.0"> <title>Move Car with WASD Keys</title> <style> body { margin: 0; overflow: hidden; background-color: lightgray; } #gameArea { width: 100vw; height: 100vh; position: relative; } #car { width: 50px; height: 100px; background-color: red; position: absolute; top: 50%; left: 50%; transform: translate(-50%, -50%); } </style> </head> <body></pre>	<p>[Prompt would have been better if contributor had given the model an 'index.html' and/or some more specific instructions about how the contributor planned to setup and execute the code. Additionally, this is a strange thing to ask for; the contributor essentially having the model build a back end on the front end. In terms of prompts, this is a rather unrealistic/unfeasible scenario for a model to address.]</p> <p>I am working on upgrading the Task Manager class in JavaScript into a sophisticated task management application. The following enhancements are required: first, implement a User model to manage multiple users, allowing each user to have their own set of tasks. Incorporate a method for adding users, ensuring tasks are associated with the specific user. Next, modify the Task model to include a tags property that allows tasks to be categorized by multiple tags. The addTask method should ensure validation for the task. Include methods for filtering tasks based on priority and sorting tasks by due date. Finally, local file storage for persistent task management should be implemented to save and load tasks automatically from the file. It should not duplicate users and tasks. Users should be loaded into the users' list, and tasks should be loaded into the tasks list by default.</p> <p>Here is the code:</p> <pre> // Task Model class Task { constructor(id, title, completed = false) { this.id = id; this.title = title; this.completed = completed; } } // Task Manager Class class TaskManager { constructor() { this.tasks = []; } addTask(task) { if (!task.id !task.title) { throw new Error('Invalid task. Ensure it has an id and title.'); } } }</pre>
--	--	--

```

<div id="gameArea">
  <div id="car"></div>
    </div>
    <script
      src="script.js"></script>
    </body>
  </html>

```

JavaScript:

```

window.onload = function() {
  const car =
document.getElementById('car')
  ;
  let carX = window.innerWidth
    / 2 - 25;
  let carY =
window.innerHeight / 2 - 50;
  const speed = 10;
  car.style.left =
    `${carX}px`;
  car.style.top = `${carY}px`;

document.addEventListener('key
  down', (event) => {
  switch
(event.key.toLowerCase()) {
    case 'w':
      carY -= speed;
      break;
    case 's':
      carY += speed;
      break;
    case 'a':
      carX -= speed;
      break;
    case 'd':
      carX += speed;
      break;
    }
  carX = Math.max(0,
    Math.min(carX,
      window.innerWidth - 50));
  carY = Math.max(0,
    Math.min(carY,
      window.innerHeight - 100));
  car.style.left =
    `${carX}px`;
  car.style.top =
    `${carY}px`;
  });
};

```

```

}
  this.tasks.push(task);
}

removeTask(taskId) {
  const index =
this.tasks.findIndex(task => task.id
  === taskId);
  if (index === -1) throw new
Error('Task ID not found.');
  this.tasks.splice(index, 1);
}

getTasks() {
  return this.tasks;
}

// Sample Usage
const taskManager = new TaskManager();
const task1 = new Task(1, "Learn
  JavaScript");
const task2 = new Task(2, "Refactor
  Task Manager");

taskManager.addTask(task1);
taskManager.addTask(task2);
console.log(taskManager.getTasks());

```

Code Ecosystem	<p>I am a solutions architect at a tech company and I need an auto-scaling mechanism for a Python-based web application running in Docker containers, deployed across a hybrid cloud infrastructure (Microsoft Azure, AWS, and GCP). The system should scale based on real-time CPU load. The auto-scaling mechanism should use a custom load monitoring tool that will run within the containers and report the CPU load directly to a central monitoring system. When the CPU load exceeds a certain threshold, the system should automatically scale across all three cloud providers. Explain how this monitoring tool should be designed and how the communication between containers running in different cloud environments would take place even when affected by network inconsistencies and varying latencies while maintaining high availability and load balancing and considering data consistency, fault tolerance, and potential for cross-cloud issues. Explain each step.</p>	
-----------------------	--	--

Lemur Astrologer Coding Goal-Oriented Multi-Turn (MT) Coding

Task Specifications

Table of Contents

[Project Overview](#)

[Task Overview](#)

[Task Specifications](#)

[Step 1: Goal Setting](#)

[Step 2: Prompt Writing & Tagging](#)

[Step 3: Response Evaluation](#)

[Step 4: Execution](#)

[Step 5: Ranking Responses](#)

[Step 6: SUPER IMPORTANT Rewriting Each Preferred Response](#)

[Step 7: Continuing the Conversation](#)

[Appendix](#)

[Prompt Examples](#)

[Response Examples](#)

[Grading Rubrics](#)

[Prompt Evaluation Rubric](#)
[Response Grading Rubric](#)
[Response Rewrite Rubric](#)

! Important Announcement:

NEW - Common gaming-related prompts are no longer allowed ("Build a tictactoe game", "Build a tetris game"). [Here is a list of prompts](#) that are commonly-seen and will be SBQ.

Read this [doc](#) for more info! [Common Prompts](#)

There's an additional consideration for the Accuracy Dimension:

Input Validation. Read the changes to align your work.

Prompts must be labeled accurately by difficulty + sub-category

Code must be fully tested - bad code will result in SBQ + removal from project

All code in rewritten responses or preferred responses are REQUIRED to have sufficient code comments!

For rewrites only: Paragraphs describing processes in multiple sentences should be converted to bullet points (or numbered lists, if the process is ordered).

Project Overview

The goal of this project is to train AI models to handle a variety of coding-specific tasks. You'll do this by having a multi-turn ("MT") conversation with the model to guide it to fulfilling a specific goal. We call this "Goal-Oriented Multi-Turn." You'll start by having a specific goal in mind and prompting the model to try to get as close to the goal as possible. Your work will identify strengths and weaknesses in the model, help identify where/when the model gets something wrong, and help improve the model's reasoning skills with each prompt and rewritten response. In a nutshell, the main objective is to craft realistic prompts that cause the model to fail.

IMPORTANT!

If you've already worked on this project, please review these [Recent Changes](#):

10/31/24 - [NEW Rating Criteria: Input Validation](#).

10/31/24 - [Changes to rewrite guidelines](#) - now required at every turn

10/31/24 - Presentation failures cannot be the only mistake/issue

[Stylistic Guidance & Changes when Rewriting Responses](#)

[Execution Instructions](#)

Task Overview

[Here is a high-level overview of the task](#):

Step 1: Goal Setting

Establish a clear and achievable goal tailored to a specific coding task and for a target programming language

Step 2: Prompt Writing & Tagging

Craft an effective prompt and properly tag its task category and difficulty level

Step 3: Response Evaluation

Ensure your prompt results in either one or both responses to your prompt failing, and then rating each model responses for accuracy, efficiency, and adherence to instructions

Step 4: Execution

Test the code in each of the model responses by uploading a screenshot of your code's output, list any setup and run commands, and show the input and actual result when executed

Step 5: Ranking Responses

Compare and rank the two responses, determining which one best aligns with the task criteria

Step 6: Rewriting the Preferred Response

Rewrite the *preferred* response to ensure it fully achieves the goal of the task

Step 7: Continuing the Conversation

Prompt the model again until reaching the minimum turn requirement to guide the model closer to the goal

Task Specifications

Step 1: Goal Setting

In this step, you'll **create a goal that matches the specified task category, target programming language, and difficulty level.**

Goals should

Be specific, clear, and easy to understand, leading the model to fulfilling the goal

Example Goal: "Build a pong game"

If you're unfamiliar with the task category or target programming language, it's best to skip the task and select another one - you will set your preferences in the courses.

After writing your goal, select the goal type and sub-category. There are **2 GOAL TYPES**, with sub-categories as follows:

Error Correction

Fixing Model Errors

Identifies and corrects specific errors in model responses.

Rectifies inaccuracies, bugs, or other mistakes in code or text.

Handling Ambiguous Requests

Trains the model to interpret and clarify unclear instructions.

Encourages reasonable assumptions or requests for clarification to fulfill tasks properly.

Multi-Step

Building Progressively

Guides model through a series of steps to achieve a complex task. Encourages structured and logical problem-solving through each step.

Deep Dive

Explores a concept in greater depth. Encourages nuanced responses and deeper understanding.

Refining Requests

Refines model's responses with more precise requests.

Adjusts and narrows prompts to improve quality and accuracy.

Step 2: Prompt Writing & Tagging

In this step, you will write a prompt that **aligns with the chosen category and difficulty level, guiding the model toward achieving the task's goal**. Additionally, you will **label the prompt** to clearly indicate the selected **category and difficulty level**. If the prompt you write covers more than one category, the **main request** must be of the **chosen category**.

Prompt Creation:

Write the initial prompt based on the task category and programming language.

Ensure each prompt is clear, specific, and challenging to the model.

Design the prompts to reflect the set difficulty level.

Use only English for all prompts.

Avoid low-effort prompts, such as single-sentence prompts or incredibly generic prompts without any constraints.

Special Notes:

Prompts specifically not allowed for this project include “counting” prompts (asking the model to do some kind of counting operation like “generate N of something”, “do something in N steps”, “give me X amount of lines”, “have this on line number x” - the customer does not want these as they are already aware of them, so this will be rejected if you submit a counting related prompt). **This includes referencing line numbers in the prompt, such as “debug the function on line 40”.**

Avoid asking the model to interpret scenarios or formulas from topics such as math, physics, etc - this is a **Coding project** and the model failures/deviations should be failures on coding logic and understanding.

Failures/deviations that rely on this approach will be SBQ'd. It's alright if a task includes formulas, as long as the model isn't expected to interpret them and penalized for not doing so correctly.

Prompts should not be related to common games (i.e. **Build a tic-tac-toe game, Build a tetris game, see this sheet for more examples**), these tasks are no longer being accepted and will be rejected.

Prompt Tagging:

Label your prompt with the appropriate [task category](#) (make sure to pay close attention for subsequent turns - this is a common error!)

IMPORTANT: Simplified for Educational Purposes - Example of common sub-category error:

Task Category: Code Generation/Synthesis

Turn 1 Prompt: "Make a 2d minigolf game using JS"

Turn 1 Category: Text to Code (*makes sense*)

Turn 2 Prompt: "Additionally add a stopwatch to the top right corner"

Turn 2 Category: Text to Code (~~X~~ bad, it should be **Text to Code Edits**, because we're asking for edits to existing code)

Next, tag the [difficulty level](#) (your prompt should match the difficulty, if you think it doesn't then tag it differently!)

Enhancing Testability:

Recommend code that runs on Replit (or any convenient IDE of your choice): This tool simplifies testing by setting up 80%+ of coding environments.

Add setup instructions: In subsequent prompts, ask the model to include setup steps for testing the code (e.g., using [pip](#) for Python or [npm](#) for Node.js).

Limit dependencies: Avoid prompts that require external dependencies like local databases or APIs, unless you can provide them and they are easily accessible for a reviewer - otherwise, you may be penalized.

Use a staged approach: For complex tasks like app development, break up the prompts into steps (e.g., pull requests) for easier testing.
You should **not** have a laundry list of requirements in your prompts.

Final Considerations:

Make your prompts **realistic** and diverse.

The goal is to simulate real-world tasks that are both complex and testable, pushing the model's capabilities while maintaining clarity.

A Prompt Error we commonly see on this project:

Not thoroughly testing code generated by the models in response to your prompt.

Solution: Test each part of the code to ensure it runs correctly and meets your requirements.

Your prompt is not actually failing the model in either response! This will result in your task being rejected - you must ensure your prompt fails **at least one of the model responses!**

Your prompt must include all of the requirements you expect and/or desire, if one of the model responses does something that is not to your preference (i.e. it solves something iteratively instead of recursively), you cannot penalize the model for that unless you specifically required in your prompt that the solution follow a specific approach.

Tips for Avoiding this Error:

Write Clear and Testable Prompts

Create prompts that result in outputs you can easily test and review.

Use Tools Like Replit

Ask the model to generate code that runs directly on platforms like Replit.

Replit supports many libraries and can automatically set up testing environments.

Minimize External Dependencies

Avoid requiring libraries or databases that are difficult to set up unless you provide clear instructions or alternative solutions.

⚠ You must review the [Prompt Examples](#) (<- click here) below to get a good sense of what is good and bad for this project. ⚡

Task Categories

Instruction Type	Category Type	Description	Code required in prompt?
------------------	---------------	-------------	--------------------------

Generation/Synthesis	Code completion Text to code Text-to-SQL	Generate immediately executable code from natural language text description or existing code snippet, infilling, etc.	No
Editing/Rewriting	Code summarization/compression Text to Code edits Code translation Code refactoring	Make changes or adjustments to existing code to meet new requirements or conditions, such as altering functionality or updating or enhancing features.	Yes (code must work properly)
Debugging	Debugging and troubleshooting Testing Security Review	Identify and correct errors in existing code, such as debugging, resolving syntax errors, and fixing logical mistakes.	Yes (code must contain a bug)
Documentation	Codebase documentation Comment generation Commit text generation API documentation Create example usages of this function Document this function	Generating or updating documentation related to code to help developers understand the code, its functionality, and its usage	Yes (code must work properly)

Review/Critique	Code review Log analysis (text -> text) Quality assurance	Code review & best practices is the process of reviewing and improving code quality, security, and maintainability by applying best practices, standards, and guidelines, and ensuring compliance with coding standards and regulations	Yes (code must work properly)
Code Ecosystem	IDEs or development workflows CLI (command line interface) Version control	Interacting with coding environments, such as IDEs, version control systems, or build tools, to automate tasks, integrate code, or manage development workflows	No

Difficulty Level

Here are the difficulty levels that you may encounter:

Please note that all prompts must be original and copying anything that already exists online on sources such as Leetcode and HackerRank etc. are strictly prohibited and will result in automatic removal from the project. The goal of having human contributors writing prompts is that we want to challenge the models in ways that the models can't already learn from the internet!

	Medium (Undergrad)	Hard (Graduate)	Challenger (SME)
Knowledge Required	Limited domain/algorithms knowledge or implementation context (e.g., architecture, libraries, pre-existing code)	Knowledge of standard algorithms and data structures, common libraries and concepts or additional code context may be	Expert domain knowledge or information on the specific application or deployment scenario, including substantial specific API/code context

		required to achieve an optimal solution	
Prompt Ambiguity	There is little ambiguity in the question (in the case of underspecification, good default behaviors are easy to come up with or not essential) and limited complexity of specifications (in instructions)	Medium ambiguity in the prompt (e.g., needs to come up with reasonable ad-hoc data representation or class structure without explicit guidance), multiple requirements should be satisfied, or multiple bugs should be found	Finding good solutions needs non-trivial design decisions regarding data structures, algorithms or code architecture/design patterns
Solution Complexity	The solution is easy to explain (e.g., code doesn't need comments to be understood) and to test for/debug (limited corner cases)	Involves corner cases that should be dealt with separately; an explanation of the solution requires some abstraction or decomposition of the problem into a few subproblems	Finding a solution requires solving several non-trivial subproblems or finding non-trivial bugs; a problem involves tricky corner cases, and explaining the solution to a non-expert requires adding context

Prompt Example	<p><i>I have a folder that contains MIDI covers of songs that I created myself. Here is what the song metadata looks like:</i></p> <pre>[{"songData": {"title": "", "artist": "", "album": "", "duration": ""}, "midiFilePath": ""}]</pre> <p><i>I want a React app that displays my MIDI music library with a play button next to each song. Clicking the play button should play the MIDI file. Use the midi-player-js library. Create some visualization based on the MIDI sequence while a song is playing.</i></p>	<p>A requirement dictated by management is that files are NOT allowed to use import statements that are erroneous, i.e. a package is imported but none of the methods are utilized in the file it is included inside of. The reasoning being that this can be confusing, wastes space, and is just sloppy programming. We have a quite large Python codebase and it would be unfeasible for someone to manually comb through all the files and check for this condition. Please help write a script in Python which accepts a root folder as an input and then recursively searches through all files and directories, checking the import statements used within each file and ensuring that they are not erroneous. Any files found should be recorded in a text file for manual review.</p>	<p><i>I work for an oil company in the Midwest, and we are drilling in some newly acquired fields that supposedly contain oil and other gas deposits five layers deep into the Earth. However, it has been brought to our attention recently that there are ore and precious metal deposits where we planned on drilling. This is problematic because our equipment for drilling for oil and gas will be damaged if it comes into contact with the precious metals. Furthermore, the precious metals will be destroyed in the process. We are trying to create a system where we can analyze the depth of where we plan to drill and determine where the metals are placed in the holes amongst the oil and gas so we can plan how we will drill without damaging anything. If we can do this programmatically, we hope to implement it with our equipment so they can drill automatically by analyzing the soil at each depth. How can we generate an artificial dataset and implement this system in Python, TensorFlow?</i></p>
-----------------------	---	--	--

! Prompt Examples (Highly recommended)

Please review the examples here in the: [Prompt Examples](#) which is in the appendix of these instructions

For more examples, please reference the [prompt example bank](#).

Step 3: Evaluating Responses & Continuing the Task

In this step, you will assess the model's responses based on specific criteria and provide follow-up prompts to improve its output.

3a. Producing at least one bad response

Every step should involve **at least one bad model response resulting from the prompt you provide**. If both responses are good, please write a more challenging prompt and try again.

How to quickly identify a good vs bad response: If there are any areas in the rating dimensions where a response has **at least one issue in the P0 or P1 level, see below**, we can count that response as "bad."

If a response only has issues in the presentation criteria (P2) that is **not** sufficient enough to be rated as "bad"!

3b. Rating the Responses based on each dimension

For each turn, after specifying how the two models did overall, each individual response **should** be rated on performance according to the following dimensions along with a brief explanation.

Quick snapshot of the dimensions (more detailed table below) and the level of priority (meaning these dimensions are most important vs less important):

Level of priority should be used to **determine which response is better with regards to your rating for the turn**. For example, if R1 has an instruction following issue while R2 has efficiency issues, R2 would be the better response.

Instruction Following: The response answers all requests in the prompt
This checks if the model understood **all** of the requests of the prompt and is addressing each one.

Make sure to check the entire implementation. If the model *tries* to address the request but *fails to*, it is an **accuracy** issue.

Priority Level: P0 - Highest Importance

Accuracy: All claims and code in the response is accurate and fully correct
You will need to Google the claims made by the model or execute code to check this

Priority Level: P0 - Highest Importance

Sub category - Input Validation/Error Handling: The response covers all meaningful input edge cases and handles them correctly.
The response validates all inputs, handling invalid data.

Priority Level: P0 - Highest Importance

Optimality and Efficiency: The response presents the most optimal and efficient solutions
The response is using common practices and standards

Priority Level: P1 - Second Highest Importance

Up-to-Date: The response uses only the most recent APIs, functions, or libraries available.

APIs, functions, or libraries used aren't causing compilation or runtime errors due to deprecation.

Priority Level: P1 - Second Highest Importance

Presentation: The response format follows the Style and Presentation guidelines

It should follow the presentation rubric, such as:

Enough comments in the code.

A professional tone (no pleasantries / fluff).

An answer that is concise, without repetitive statements.

Explanations that use bullet points.

Priority Level: P2 - Third Highest Importance

Each dimension rating should include a brief explanation as to why the given rating was chosen (e.g., why "Major Issues" was selected for Accuracy). The explanations/justifications you write for dimension ratings should be **specific** and **clear**. If there are any bugs, or there are issues, specify **what the issues are** and **what the causes of the issues are**.

Avoid generalizations like "No Issue", "N/A", "This has an error", etc. If there are no issues, give a very brief explanation as to why.

Examples of Bad Justifications:

"There is a syntax error": ✗ This is only half of the story - you should briefly explain what's causing the syntax error.

"The response doesn't satisfy all of the prompt requirements": ✗ This is very vague. If you find that the response doesn't satisfy all of the requirements, you should be specifying which requirements it doesn't satisfy and why it doesn't satisfy them.

Response Rating Rubric

Dimensions	NA (0)	No Issue (1)	Minor Issue (2)	Major Issue (3)
Instruction following	<i>This dimension cannot be NA.</i>	The response meets the main request and all constraints, showing a strong understanding of the prompt, even if there are minor implementation errors. It handles any ambiguities well and stays within the specified requirements.	The response fulfills the primary request but does not entirely adhere to all the constraints. The response could have better handled the ambiguity of the prompt. Common errors: Fails some but not ALL constraints	The response fails to fulfill the primary request OR fulfills the primary request but does not adhere to any constraints. Common errors: Fails to do the primary request

Accuracy and [NEW] Input Validation / Error Handling	<p>Can only be NA if the response contains no code or factual claims, and does not rely on prior context.</p> <p>The code runs error-free, produces the correct output, and follows best practices. All text and comments are accurate, and the response is contextually appropriate with any previous errors fixed.</p> <p>–</p> <p>All meaningful edge cases are covered. The response includes thorough validation for expected inputs and error handling for invalid data, ensuring robustness and resilience to common input errors.</p> <p>Example:</p> <p>The function <code>calculate_discount(price, discount_percent)</code> validates all inputs. It checks that <code>price</code> and <code>discount_percent</code> are positive numbers, ensures <code>discount_percent</code> does not exceed 100%, and returns a clear error message if values are out of expected ranges or of incorrect types (e.g., strings).</p>	<p>The code runs but has minor warnings or low-risk security issues. The content is mostly accurate, but some statements are unclear or make unproven claims.</p> <p>Previous errors remain but don't affect the current response.</p> <p>–</p> <p>Some edge cases are covered. The response handles most expected inputs but misses certain edge cases, which could lead to potential errors or exceptions under specific conditions.</p> <p>The function <code>calculate_discount(price, discount_percent)</code> includes basic validation, such as checking that <code>price</code> and <code>discount_percent</code> are positive numbers. However, it lacks checks for certain edge cases, such as ensuring <code>discount_percent</code> does not exceed 100% or verifying that inputs are numeric.</p>	<p>The code doesn't run due to logic errors, produces incorrect output, or has major security flaws. The response includes false claims, lacks context, and previous errors were not fixed, making the issues worse.</p> <p>–</p> <p>No edge cases are covered. The response lacks validation for all inputs, making it vulnerable to errors when faced with unexpected or invalid data inputs.</p> <p>The function <code>calculate_discount(price, discount_percentage)</code> performs no validation on its inputs, assuming <code>price</code> and <code>discount_percentage</code> are always valid and within expected ranges. This could cause runtime errors or incorrect results if given invalid inputs, such as negative numbers, <code>discount_percentage</code> over 100%, or non-numeric types, making the function unreliable.</p>
Optimality and Efficiency	<p>Can only be NA if the response contains no code using functions or statements aside from</p>	<p>The code is well-optimized, handles edge cases, and follows standard best practices. If top performance isn't required, it still performs efficiently without adding unnecessary complexity.</p>	<p>The code performs well but could use minor optimizations. It generally follows best practices but may not scale for large datasets.</p>

	the assignment			
Presentat ion WHEN REWRITI NG, YOU MUST FIX ALL PRESEN TATION ISSUES	<i>This dimension cannot be NA.</i>	<p>The code is well-documented, with clear comments and explanations for any modifications. Code included in the prompt that did not originally have comments should have comments if included in the response. The response is concise, well-organized, and uses readable variable and function names. Complex processes are broken down with bullets, and Markdown is correctly formatted with clear hierarchies.</p> <p>Formatting is neat, with triple backticks for code blocks, and proper use of bold and italics for emphasis. White space and line breaks improve readability, and tables are correctly aligned. Functions are modular and follow standard patterns, such as using <code>if __name__ == "__main__":</code> blocks for structure. There are no redundant solutions provided for the same problem.</p>	<p>The documentation is generally clear but could use more detail. There are minor language errors that don't affect readability, and formatting could be improved for clarity. Variable and function names are understandable, but some structural changes—like adding bullets or logical sections—would help. Functions are present but may need more modularity, and some explanations are missing, making the code harder to follow in parts.</p> <p>Common Errors</p> <ul style="list-style-type: none"> Uses backticks inconsistently Uses camelCase and snake_case inconsistently 	<p>The documentation is missing or inadequate, or lacking code comments entirely, making the code hard to understand. The response is poorly formatted and lacks structure, with unclear variable and function names. The logic is disorganized, and there are no explanations for key decisions, making it difficult to follow, integrate, or reuse. Programming language tags are also missing.</p>
Up-to-Dat e	NA (0)	Up-to-Date	Out-of-Date	
	The code does not call on any libraries or	The code uses the most fresh API, libraries, or functions available to solve problems efficiently. The code uses a maintained library or function which is an older version that still	The code uses a deprecated API, library or function, causing a runtime or compile-time error.	

functions.	works (even if it is less efficient).	
------------	---------------------------------------	--

Step 4: Execution (YOU MUST TEST YOUR CODE!)

For any response that changes the executable code, you will also be required to execute the code in the response. Note: *This doesn't apply to responses that have only added/modified code comments.*

Below are the steps that you will encounter for running code!

Execution Guidelines

Program
ming
Language
:

What additional languages were involved in the response? (Select all that apply)

Exclude the programming language assigned to the task! Applies to both code and related discussion in the model response.

Hover over the hint tooltip for the list of possible selections

Choice Paths:

html|

HTML/CSS



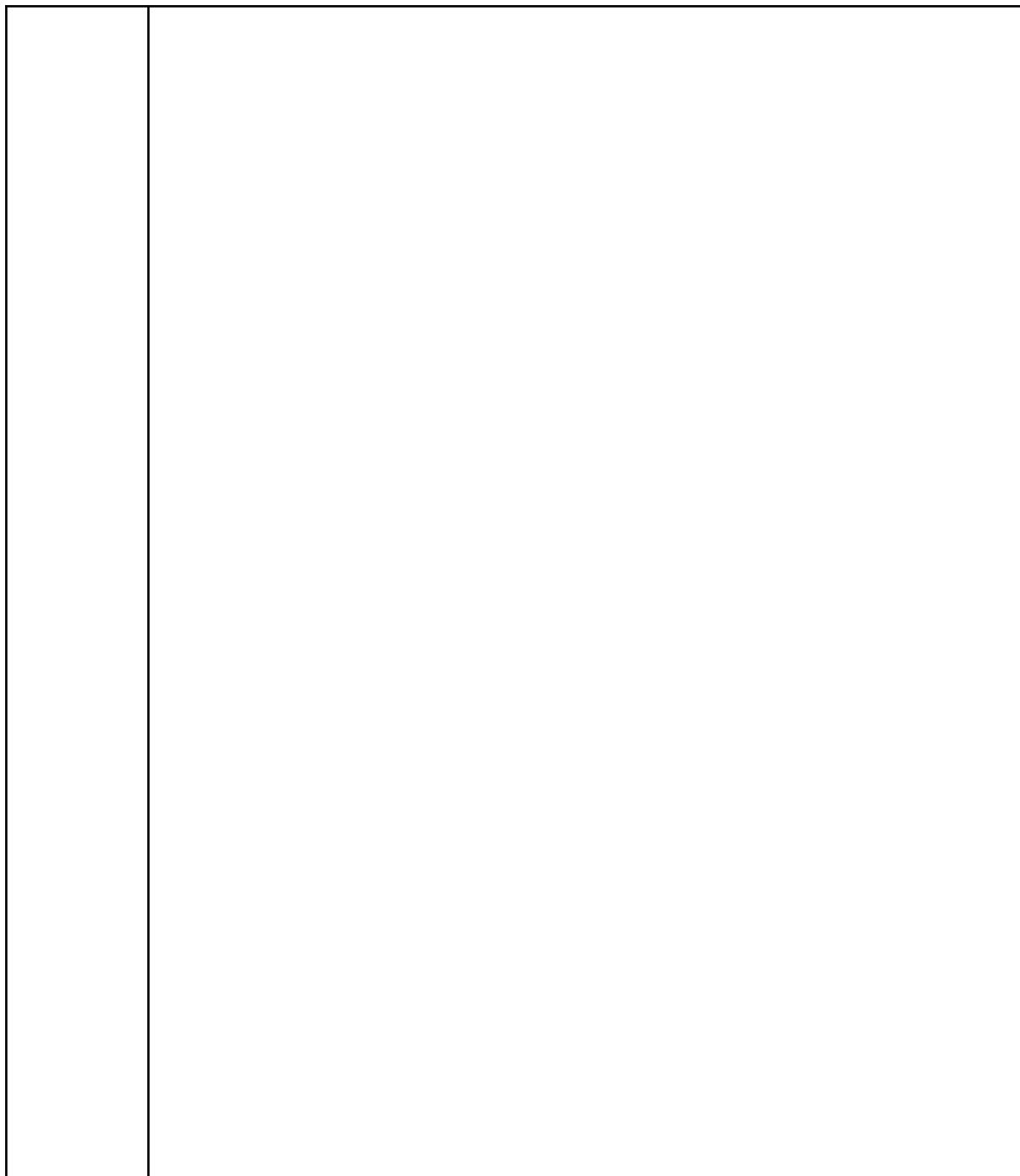




Executable Code:	<p>Does the model response edit or add any executable code? *</p> <p>Mark "Yes" if there is any new or altered code present in the response, regardless it is an entire program or just a code snippet. Does not apply to new comments!</p> <p><input type="radio"/> Yes <input type="radio"/> No</p>







Degree of Execution :

Rate the degree of execution of the code *

(Only if there is executable code in the response)

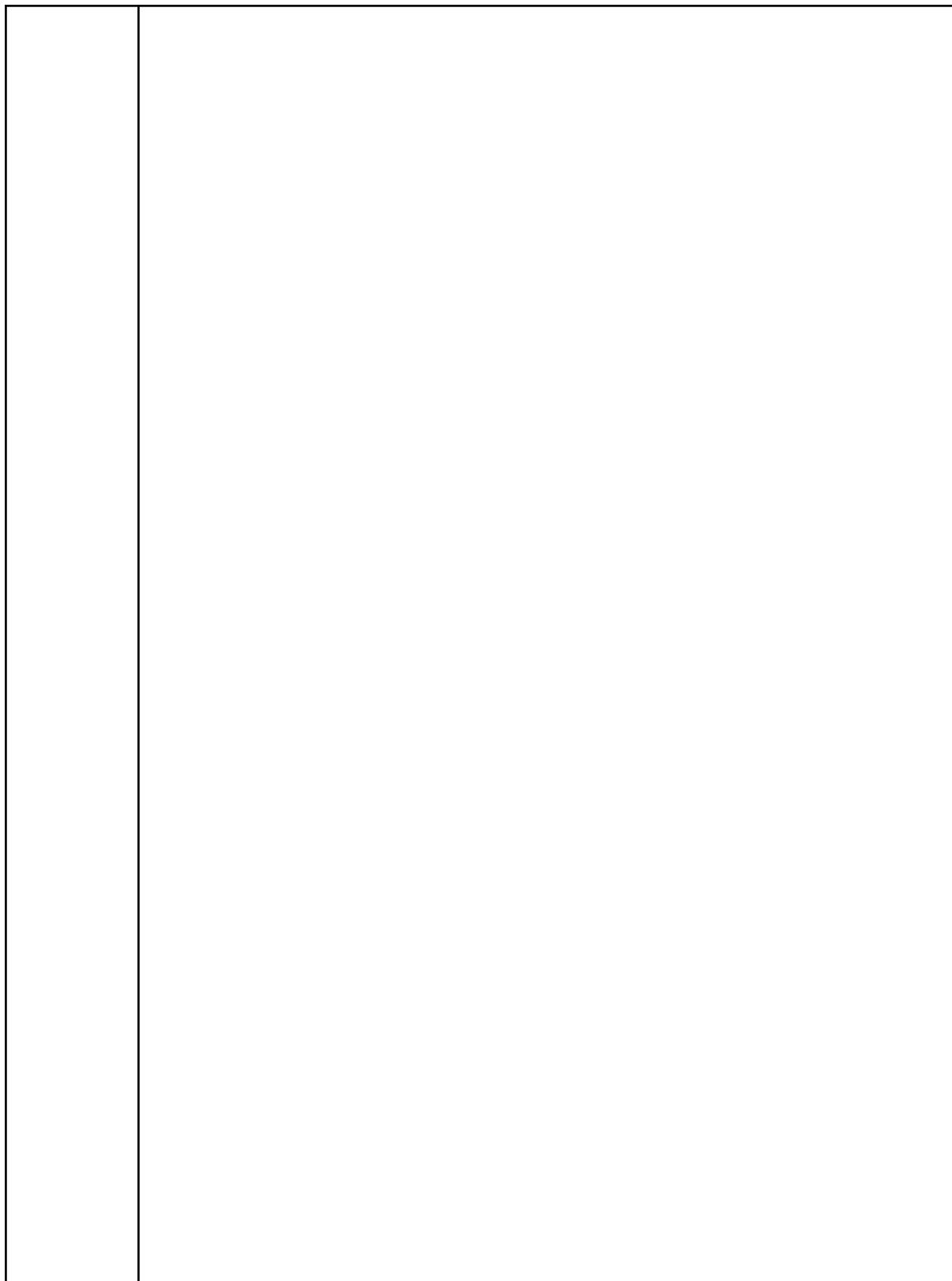
- Template ⓘ
- Partial Update ⓘ
- Function Update ⓘ
- Out-of-the-Box ⓘ
- NA











**Installatio
n
Comman
ds:**

What installation commands are necessary to test the code? *

Please separate each individual command with a comma (,) and write 'N/A' if there are no extra install commands required.

ex: pip install tkinter, sudo apt get







Run Comman ds:	<p>What commands are necessary to run the code? Provide a comma separated list if there are multiple commands. *</p> <p>ex: uvicorn main:app --host 0.0.0.0 --port 8000, python app.py</p> <div style="border: 1px solid #ccc; height: 60px; width: 100%;"></div>





Output Expectation:	<p>Is the output of the code as expected? *</p> <div style="border: 1px solid #ccc; padding: 5px; width: fit-content;"><input checked="" type="radio"/> Yes <input type="radio"/> No</div>



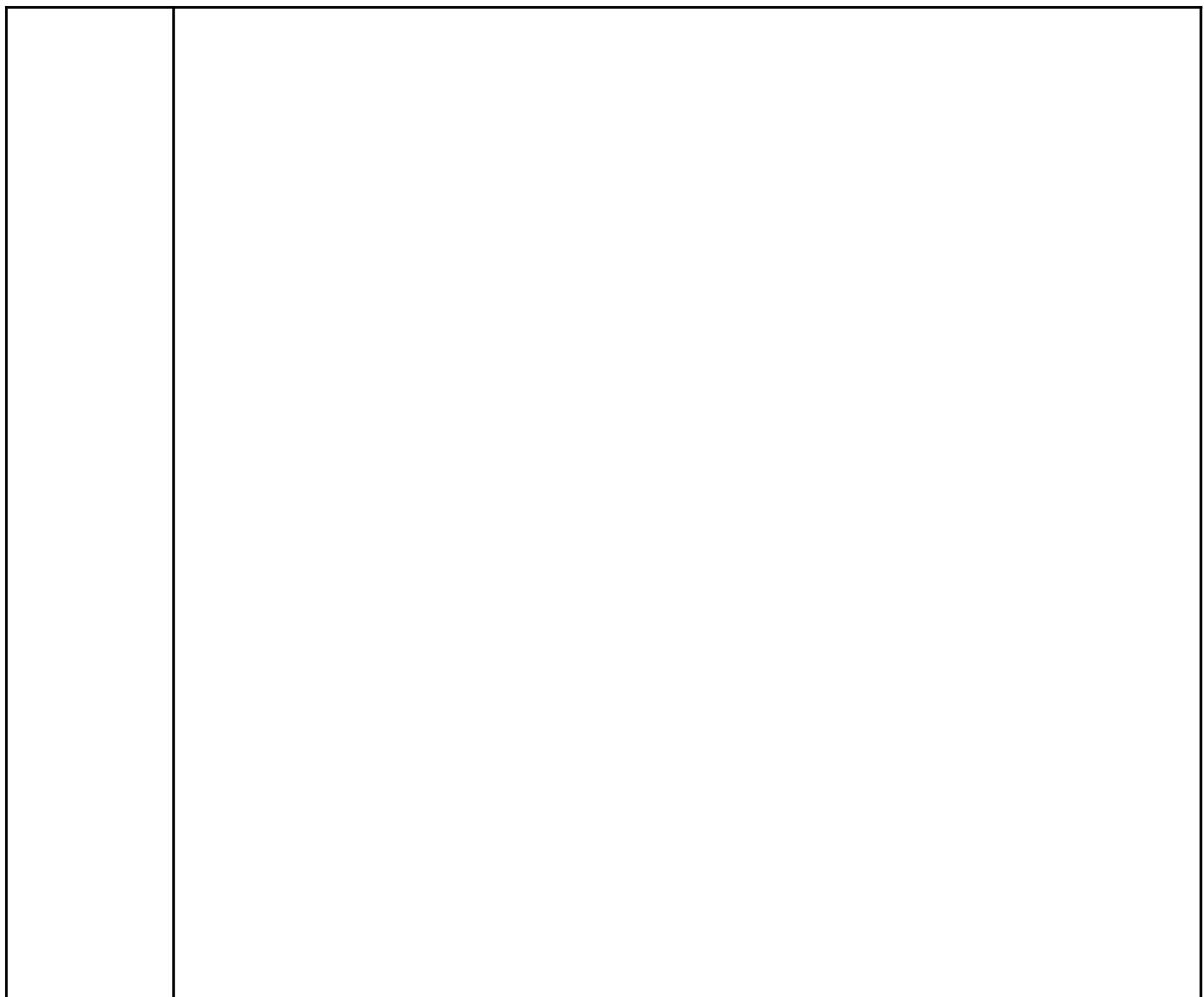


Output MIME Type:	<p>What is the mime type of the output? *</p> <div style="border: 1px solid #ccc; padding: 5px; width: fit-content;"><input type="radio"/> text/plain <input type="radio"/> Other (specify below)</div>



Record any errors:	<p>Did the code produce an error? *</p> <p><input checked="" type="radio"/> Yes</p> <p><input type="radio"/> No</p> <p>Code Execution: Error</p> <p>Please paste the error from the code execution.</p> <div style="border: 1px solid #ccc; height: 100px; width: 100%;"></div>

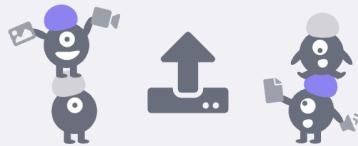




Upload Screenshot/ Screen Recording of Output

Upload Screenshot/Screen Recording of Output (required)

Upload a screenshot of your code output (including the execution call). If the output of the code is not static, upload a screen recording instead.



Drag and drop here or click to upload your data













Upload Files	<p>Is your code a single file or multiple files? *</p> <p>Ignore any third party dependencies and libraries</p> <p><input checked="" type="radio"/> Single file <input type="radio"/> Multiple files</p> <p>File Name * Write the name of the single file containing your code. This should be the same as the one used in your run commands. (ex: main.py)</p> <p>main.py</p> <p>Paste Code (required) * Paste your entire codefile here. DO NOT USE IF YOU HAVE MULTIPLE FILES (see previous question)</p> <pre>def read_file(file_path): try: with open(file_path, 'r') as file: data = file.read()</pre> <p>Upload Code (required) Upload a zip file containing the code you used to test the response. Remember: The code you upload should be executable (even if it throws an error), and should correspond to the Execution Instructions you provided above.</p> <p>Drag and drop here or click to upload your data</p>







Output:	



<p>Output</p> <p>Expected:</p>	





Output

MIME

Type:



JS	<p>JSON Template Please leave blank!</p> <div style="border: 1px solid #ccc; padding: 5px; width: fit-content;">DO NOT FILL!</div>

--	--

Step 5: Ranking Responses

In this step you will compare the two model responses and rate which one is better on a scale of 1-8, indicating whether Response 1 or Response 2 performed better.

This should match how you rated the models on each dimension (i.e. if Response 2 performed worse on every dimension, then the score should be a 1).

Ranking Justification Guidelines:

Base your ranking on the rating dimensions in the rubric, listed by importance.

Provide helpful and factual feedback without unnecessary comments or conversation.

Clear Documentation: Use docstrings for code explanations and modify code directly in your response, avoiding extra copy-pasting.

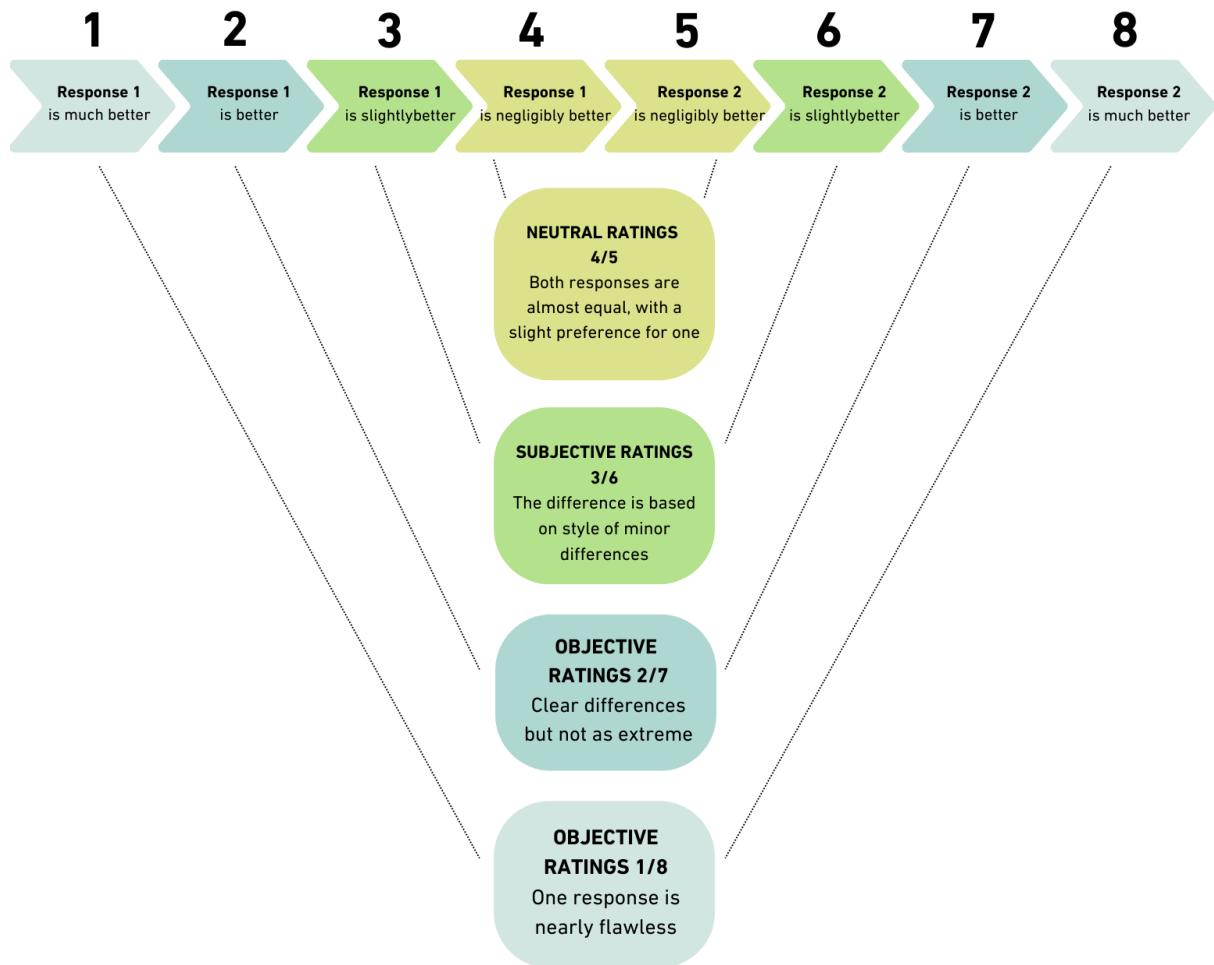
Readability and Structure: Ensure code is easy to read, with a focus on structured and modular design (e.g., use functions and `if __name__ == "__main__":` blocks).

Use bullet points or numbered lists for clear and organized explanations.

Choose the response that is **easier to understand and integrate** with minimal user effort.

Rating Scale

Scores from 1-8 will demonstrate either a strong or neutral preference between the two responses



Use the [priorities mentioned here](#) to determine which response is better based on the severity of the issues.

Step 6: Rewriting Each Preferred Response

In this step, once you have selected the better response from each turn, you will be required to **always** rewrite the **preferred** response to achieve the goal of the prompt and fix any issues that were identified.

The only scenario where a rewrite would **not** be necessary is if one of the model responses is completely **perfect**, meeting all the requirements and specifications of the prompt. In this case, indicate that a **response rewrite is not needed**. This is very rare and requires a thorough justification, in most cases, the response will require a rewrite.

Specification List

The re-written response should **fully achieve the most recent prompt**.

The re-written response should **address any secondary objectives implied by the prompt**.

The re-written response should **correct all errors (major or minor)**

The re-written response should **be coherent and logically connected to the prior conversation**.

The re-written response should **fulfill all the dimensions for "No Issues" in the rating rubrics**.

The re-written response should **follow the formatting specifications**.
The code of the re-written response should **contain enough comments**.
The code of the re-written response MUST run properly and be optimal and efficient. Failure to properly test code results in removal from the project.

Formatting/Presentation Requirements

Key terms **should** be highlighted in bold, whereas titles, articles, etc. are italicized.
Remove pleasantries such as "Sure," "Certainly," "I can help with that," etc.
Make responses more concise, remove all fluff and unnecessary phrases (i.e. "*Welcome to the world of VS Code!*")
Tone should be straightforward/professional.
Code should be well-commented
Test outputs include a comment with the expected response.
Explanations should use bullet points.
Rewritten response replaces any paragraphs (especially those with at least 3 points) with bullet points instead
All repetitive phrasing/wording must be removed.

Here are examples of [formatting/presentation requirements](#)
Make sure to keep track of the changes that you made as you'll have to write them out in the next step.

Stylistic Guidance & Changes for Rewrites

Stylistic Change	Description	Bad Example	Better Example
Use of 're' and 's' instead of 'for' and 'in'	that loops through a list of items to apply a discount to each item.	<pre>for item in items: apply_discount(item)</pre>	<pre>for item in items: apply_discount(item)</pre>

options	titive e.g., "We ions to arity.	<i>phrasing:</i> guessedWords array to of words that have en guessed." the checkSecretWord check if an input word y been guessed." checkIfAllGuessed	A guessedWords array was added to track guessed words. The checkSecretWord function displays an alert if a word has been guessed. The checkIfAllGuessed function ensures all secret words are handled and manages animations."
or nts	lists more	<i>ed:</i> import libraries like umpy, and sklearn. create a dataset, ategorical variables, split train a model."	Import libraries: pandas, numpy, sklearn. Create a dummy dataset with 5 attributes and 1 target variable. 1-hot coded categorical variables.. Split data into training and testing sets. Train a Random Forest model.
it g	when or sets.	<i>nt formatting:</i> g one keyword but not uses matplotlib and	<i>t formatting:</i> g both keywords for clarity and consistency. uses matplotlib and tkinter ."
for steps	steps lists ch	input, validate user st the database, and in the database."	for clarity: he input for errors. user data against the database. e data in the database.

Step 7: Continuing the Conversation

Multi-Turn Conversation Guidelines:

For goal-oriented multi-turn tasks, structure the conversation to build logically toward the final goal:

Build Progressively

Each prompt should add value and guide the model closer to completing the task.

Encourage Depthng

As the conversation advances, focus on enhancing the model's understanding and pushing for more detailed responses.

Evaluate and Iterate

Assess each response, rewrite as needed, and ensure every turn contributes meaningfully to the final solution.

If the required number of turns has been met, you can end the task, even if the goal hasn't been fully achieved. This approach maintains focus and prevents unnecessary iterations.

PRO TIPS!

For Subsequent Turns, you may choose any category or difficulty but they must be aligned with your prompt

THE OVERALL IDEA of this project is to test the model's ability to recall context from prior turns

Subsequent Prompts **should not be disconnected from previous prompts**

It's easier to go back and modify the goal, if the model is doing well in recalling and building on previous context

You don't need to "fully complete" a goal, as it is simply an alignment tool

EXAMPLE: If the goal was to build a game, having parts of that game built is fine

Appendix

Grading Rubrics

Prompt Evaluation Rubric

This is the rubric we use to measure the quality of your prompts.

Criteria	1-2 (Fail)	3 (Okay)	4-5 (Good/ Perfect)	Additional Notes
Prompt Instruction Type	<ul style="list-style-type: none"> The initial prompt does not follow the specified instruction type. Subsequent prompts are not labeled correctly. 	N/A	<ul style="list-style-type: none"> The prompt follows the specified instruction type. Subsequent prompts are labeled correctly. 	<ul style="list-style-type: none"> Instruction types include: Generation/Synthesis Editing/Rewriting Debugging Documentation Review/Critique Code Coherence
Prompt Difficulty	<ul style="list-style-type: none"> The difficulty does not match the specified level. 	N/A	<ul style="list-style-type: none"> The prompt reflects the specified difficulty level. 	<ul style="list-style-type: none"> Difficulty levels include: Medium (Undergrad) Hard (Graduate) Challenger (SME)
Prompt Continuity	<ul style="list-style-type: none"> The prompt does not logically follow the conversation. It fails to achieve the conversation goals. 	<ul style="list-style-type: none"> The prompt may have slight shifts but still furthers the conversation. 	<ul style="list-style-type: none"> The prompt logically continues the conversation. Topic shifts are allowed if they help achieve the conversation goals. 	<ul style="list-style-type: none"> Applies if not the first turn.
Clarity & Specificity	<ul style="list-style-type: none"> The prompt is vague and unclear. Critical details are missing. 	<ul style="list-style-type: none"> The prompt is mostly clear but may be interpreted in multiple ways. 	<ul style="list-style-type: none"> The prompt is clear and specific. No assumptions are needed to answer. 	<ul style="list-style-type: none"> The prompt should be in English.
Feasibility	<ul style="list-style-type: none"> The request is impractical or impossible for an AI in a single response. Conflicting instructions are given. 	<ul style="list-style-type: none"> The request is feasible but requires effort or compromises. 	<ul style="list-style-type: none"> The request is fully actionable and realistic. No conflicting instructions. 	<ul style="list-style-type: none"> Requests should not contain conflicting or impractical elements. Example: Asking for a complex algorithm in one step.

Response Grading Rubric

This is the rubric we use to grade your response ratings.

Field	1-2 (Fail)	3 (Okay)	4-5 (Good/ Perfect)
-------	------------	----------	---------------------

Model Responses Rating	<ul style="list-style-type: none"> Major Rating Disagreement: (e.g., marked "Good" vs. "Bad") and poor justification. 	N/A	Rating aligns with expectations.
Prompt Adherence and Instruction following	<p>Major Rating Disagreement: The prompt adherence or instruction following is incorrect (e.g., marked "N/A" instead of "Major Issues").</p> <p>Minor Rating Disagreement (2+): Minor disagreements occur on 2 or more occasions.</p>	<p>Minor Rating Disagreement (1): Minor variance on 1 occasion, and justification does not fully support this variance.</p>	Rating is correct, or "N/A" is appropriately used when the response cannot be assessed.
Correctness and Accuracy [input validation]	<p>Major Rating Disagreement: Significant issues with correctness (e.g., marked "Major Issues") or minor disagreements on 2+ occasions.</p> <p>No edge cases are covered. The response lacks validation for all inputs, making it vulnerable to errors when faced with unexpected or invalid data inputs.</p>	<p>Minor Rating Disagreement (1): Slight variance on 1 occasion, justification lacks support.</p> <p>Some edge cases are covered. The response handles most expected inputs but misses certain edge cases, which could lead to potential errors or exceptions under specific conditions.</p>	Rating matches expectations, or "N/A" is used when applicable. All meaningful edge cases are covered. The response includes thorough validation for expected inputs and error handling for invalid data, ensuring robustness and resilience to common input errors.
Performance, Optimality, and Efficiency	<p>Egregious Rating Disagreement: Issues with performance were missed (e.g., not marking "Major Issues"), with flawed justification.</p>	<p>Major Issues: Minor performance issues were missed, but the justification somewhat supports the rating.</p>	No major issues missed in performance.
Readability, Documentation, and Presentation	<p>Egregious Rating Disagreement: Significant documentation or readability issues were not marked.</p>	<p>Major Issues: Some readability or presentation issues were missed, but the justification is somewhat valid.</p>	No major readability or documentation issues were missed.

Up-to-Date	Egregious Rating Disagreement: Incorrectly marked as "N/A" or "Up-to-Date" when it should have been "Out-of-Date."	N/A	Rating is accurate or correctly identifies outdated information.
MT Conversation Quality	<p>Task with multiple turns does not consider holistic, conversation-level factors. Conversation is incoherent or turns don't follow logically. Task doesn't require multi-turn structure (MT) or lacks context from prior turns.</p>	<p>Task is somewhat long-winded but follows logically, considering conversational factors. Conversation accomplishes its goal but could be completed in one turn.</p>	<p>Task fulfills the main goal of the conversation. Turns follow logically and handle topic shifts effectively. Each turn builds on the prior, maintaining consistency and context. New turns incorporate older turns for elaboration and coherence. Multi-turn context</p>

				is crucial for generating appropriate responses.
--	--	--	--	--

Response Rewrite Rubric (if applicable)

This is the rubric we use to measure the quality of your rewritten responses

Field	1-2 (Fail)	3 (Okay)	4-5 (Good/Perfect)	Additional Notes
Accuracy	<ul style="list-style-type: none"> Major Factual Errors: Response has 1+ major factual errors or misleading points. Minor Factual Errors: Response has 2+ minor factual errors. 	<ul style="list-style-type: none"> Contains only 1 minor factual error or misleading statement. 	<ul style="list-style-type: none"> No factual errors or misleading statements. 	<ul style="list-style-type: none"> A major error involves incorrect/misleading data central to the request. A minor error is near the subject matter but doesn't affect the main point.
Instruction Following / Response Fulfillment	<ul style="list-style-type: none"> Main Goal Miss: Does not achieve the main goal or make progress in the conversation. Explicit Instruction Miss: Misses 1+ key instructions. Not Fulfilled: Fails to answer the question. 	<ul style="list-style-type: none"> All explicit instructions are followed. Secondary Objectives Miss: Some secondary objectives not addressed. Subjective Instruction Miss: Subjectively misses some parts. 	<ul style="list-style-type: none"> Fully achieves the main goal or makes clear progress. Follows all instructions and fully answers the question. 	<ul style="list-style-type: none"> Rule of thumb: for word count, $\pm 10\%$ is acceptable. Example: If a question asks for a historical event year, even if implied, the year should be provided.
Unnecessary Greetings / Pleasantries	<ul style="list-style-type: none"> Contains greetings/pleasantries such as "Sure, I'd love to help." 	N/A	<ul style="list-style-type: none"> No unnecessary greetings or pleasantries. 	<ul style="list-style-type: none"> Only flag unnecessary phrases at the beginning or end of the response. Phrases like "Here is..." are not

				considered pleasantries.
Depth / Nuance	<ul style="list-style-type: none"> Little to No Detail: Response is superficial and lacks meaningful depth or insight. Excessive Detail: Overly complex and obscures key points. 	<ul style="list-style-type: none"> Has enough detail but may need more depth or nuance. Too Much Detail: Slightly too much, but doesn't obscure key points. 	<ul style="list-style-type: none"> Balanced, insightful, and focused without going overboard. 	<ul style="list-style-type: none"> Too much detail can lead to confusion. Be clear about what points are most important.
[Rewrite/SxS] Clearly Worse Than Model Response	<ul style="list-style-type: none"> Worse Than Original: Clearly performs worse than the original across rubric categories. Worse Than Side by Side Model: Performs worse overall. 	<ul style="list-style-type: none"> Performs similarly to the original or side-by-side comparison. 	<ul style="list-style-type: none"> Performs better overall than the original or side-by-side comparison. 	<ul style="list-style-type: none"> Don't penalize for minimal/no changes if the original response was acceptable. Compare against state-of-the-art models.

Prompt Examples

Prompt Category	Good Examples (specific)	Bad Examples (vague)
Code generation	<p>I am a software engineer analyzing our CDN service performance. I want to write a Bash script to extract the related data from the cache statistics report, count the hit ratio, and filter out the cache items' ID with a hit ratio of less than 0.85.</p> <ol style="list-style-type: none"> The cache statistics report is a CSV file. Its header row contains the Item ID, Name, Viewer Location, Time, Request Count, Hit Count, Miss Count, and Error Count. And the Item ID is unique. The input to the script is the name of the report. Before calculating, first check whether the count is valid. If any Request Count or Hit Count is missing or Hit Count is greater than Request Count, the script 	<p>[Asking for a Sudoku solver is too generic of a prompt, and it lacks complexity.]</p> <p>Hello, I want you to create an automatic sudoku solver in Bash that will take a .txt file that includes 81 numbers from the top row to the bottom with empty cells replaced by zeros in the first line.</p>

	<p>will log the missing or wrong data row and skip it. The hit ratio formula is Hit Count divided by Request Count. The hit ratio should be rounded to two decimal places.</p> <p>4. The output file is a text file containing unique item IDs and their hit ratio, sorted by ratio in descending order.</p>
--	--

Code debugging

As a friend group, we really love to play DnD and other roleplaying games. For the upcoming weekend, I wanted to surprise my friends by developing a multiplayer text-based adventure game with C++. In this game players can explore a fantasy world, interact with characters, and solve puzzles. The game uses a complex system of pointers and dynamic memory allocation to manage player actions, game state, and world events.

However, there are several bugs in the code that leads me to crashes, memory leaks, and incorrect game states. I cannot test the game because it outputs nothing. Could you help me to identify the bugs in the code and fix them?

```
#include <iostream>
#include <string>
#include <map>
using namespace std;
class GameEvent {
public:
    string description;
    bool completed;
    GameEvent(string desc) : description(desc),
    completed(false) {}
};
class Character {
public:
    string name;
    int health;
    Character(string charName) : name(charName), health(100)
    {}
    void takeDamage(int amount) {
        health -= amount;
        if (health < 0) health
= 0;
    }
};
class Player : public
Character {
public:
    vector<GameEvent*> events;
    Player(string name) :
    Character(name) {}
```

[This is bad because it is too vague for a debugging prompt, being “concerned” over code is not enough to warrant a debugging prompt, a debugging prompt points to a very specific issue like an error or crash; additionally the prompt veers off into being a code editing prompt in the second part which is bad]

Hello, I'm currently working on a JavaScript project which is aimed to actively monitor system performance statistics, as well as network activity and performance. It will generate logs every hour and properly handle data so that memory doesn't become an issue. Currently though, I'm suspicious of the log file handling capabilities of the code, and the asynchronous os.cpuUsage has me concerned. Please look over my code and confirm or deny if these concerns are warranted, make any changes you feel would improve the code. Also, I intend to include an active alert system for critical performance issues, as well as a text-based report which will compare logs and display a comparison for the user. Please implement these. Also keep in mind this is a multi-person project, please do not include OS specific libraries and functions.

```
// Import the os-utils library
const os = require('os-utils');
const fs = require('fs');
// Function to monitor network
performance
function monitorNetworkPerformance() {
    let networkData = {
        downloadSpeed: 0,
        uploadSpeed: 0,
        latency: 0,
        packetLoss: 0
    };
    // Use the navigator.connection
    API to get network performance data
    if (navigator.connection) {
        networkData.downloadSpeed =
        navigator.connection.downlink; // in
        Mbps
        networkData.uploadSpeed =
        navigator.connection.uplink; // in
        Mbps
        networkData.latency =
        navigator.connection.rtt; // in ms
```

```

void addEvent(GameEvent* event) {
    events.push_back(event);
}
void completeEvent(string eventDescription) {
    for (auto& event : events) {
        if (event->description == eventDescription && !event->completed) {
            event->completed = true;
            cout << name << " completed: " << event->description << endl;
            return;
        }
    }
    cout << "Event not found or already completed." << endl;
}
class Game {
private:
    map<string, Character*> characters;
    Player* currentPlayer;
public:
    Game() : currentPlayer(nullptr) {}
    ~Game() {
        for (auto& pair : characters) {
            delete pair.second;
        }
    }
    void addCharacter(string name) {
        characters[name] = new Character(name);
    }
    void setCurrentPlayer(string name) {
        if (characters.find(name) != characters.end()) {
            currentPlayer =
static_cast<Player*>(characters[name]);
        } else {
            networkData.packetLoss =
navigator.connection.effectiveType; // effective network type
        } else {
            console.error("Network Information API not supported by this browser.");
        }
        return networkData;
    }
    // Function to monitor system performance
    function monitorSystemPerformance() {
        let systemData = {
            cpuUsage: 0,
            memoryUsage: 0
        };
        // Use os-utils to get system performance data
        os.cpuUsage(function(v) {
            systemData.cpuUsage = v * 100;
// in %
        });
        systemData.memoryUsage = (1 - os.freememPercentage()) * 100; // in %
        return systemData;
    }
    // Function to log data
    function logData() {
        const networkData =
monitorNetworkPerformance();
        const systemData =
monitorSystemPerformance();
        const logEntry = {
            timestamp: new Date().toISOString(),
            networkData: networkData,
            systemData: systemData
        };
        fs.appendFile('performance_logs.json',
JSON.stringify(logEntry) + '\n', (err) => {
            if (err) {
                console.error('Error writing to log file:', err);
            } else {
                console.log('Log entry recorded:', logEntry);
            }
        });
    }
    // Function to start monitoring
    function startMonitoring() {
        setInterval(logData, 3600000); // Log data every hour
    }
}

```

```
                cout << "Character  
not found." << endl;  
            }  
        }  
        void displayStatus() {  
            for (auto& pair :  
characters) {  
                cout <<  
"Character: " << pair.first <<  
, Health: " <<  
pair.second->health << endl;  
            }  
        }  
        void simulateEvent(string  
description) {  
            if (currentPlayer ==  
nullptr) {  
                cout << "No player  
set." << endl;  
                return;  
            }  
            GameEvent* newEvent =  
new GameEvent(description);  
  
currentPlayer->addEvent(newEve  
nt);  
            cout << "New event  
added: " << description <<  
endl;  
        }  
        void resolveConflict(int  
damage) {  
            if (currentPlayer !=  
nullptr) {  
  
currentPlayer->takeDamage(dama  
ge);  
            cout <<  
currentPlayer->name << " took  
" << damage << " damage!" <<  
endl;  
        }  
    }  
};  
int main() {  
    Game myGame;  
  
myGame.addCharacter("Alice");  
  
myGame.addCharacter("Bob");  
  
myGame.setCurrentPlayer("Alice  
");  
    myGame.simulateEvent("Find  
the lost treasure");  
  
myGame.resolveConflict(20);  
// Start monitoring  
startMonitoring();
```

	<pre> myGame.setCurrentPlayer("Bob") ; myGame.simulateEvent("Save the village"); myGame.displayStatus(); return 0; } </pre>	
Code review	<p>As part of our cybersecurity module, I'm tasked with creating a Bash script that implements a basic file encryption and decryption service. The service should allow users to encrypt files using OpenSSL and decrypt them using a password. This is my initial implementation:</p> <pre> #!/bin/bash # Simple file encryption and decryption service using OpenSSL encrypt_file() { local file_to_encrypt=\$1 echo "Encrypting \$file_to_encrypt..." openssl enc -aes-256-cbc -salt -in "\$file_to_encrypt" -out "\${file_to_encrypt}.enc" -k "\$2" echo "File encrypted: \${file_to_encrypt}.enc" } decrypt_file() { local encrypted_file=\$1 echo "Decrypting \$encrypted_file..." openssl enc -aes-256-cbc -d -in "\$encrypted_file" -out "\${encrypted_file%.enc}" -k "\$2" echo "File decrypted: \${encrypted_file%.enc}" } # Check if user wants to encrypt or decrypt if [[\$1 == "encrypt" && -f \$2]]; then encrypt_file "\$2" "\$3" elif [[\$1 == "decrypt" && -f \$2]]; then decrypt_file "\$2" "\$3" else echo "Usage: \$0 [encrypt decrypt] [filename] [password]" fi </pre>	<p>[This is bad because it's too vague!]</p> <p>I'm working on the "Living the Social Life" web page. Please review the HTML and CSS code below and do the following:</p> <ul style="list-style-type: none"> Enhance the code structure and organization. Update the HTML code to use appropriate HTML5 semantic elements. <p>Here is the HTML code:</p> <p>html placeholder</p> <p>Here is the CSS code:</p> <p>css placeholder</p>

	<p>However, I have concerns about some potential security risks. For example, the current implementation accepts the encryption password in plaintext as a command-line argument, which can expose sensitive information. I also want to ensure the script handles errors more gracefully and is efficient for large file sizes. Please help review the code and suggest improvements to make it more secure and scalable. How can I avoid exposing the password and improve the script's error handling and efficiency for large files?</p>
--	--

<h3>Code Editing/Writing (Modification)</h3>	<p>I wrote HTML and JavaScript codes. I aimed to print a car on the screen and move it with the "WASD" keys. However, the car doesn't look like a car. There's only a red rectangle, and I expect you to change it to a more realistic car. The point of view should be from above, so I expect to see the car from the upper perspective. I should be able to see the 4 tires, the windscreen, and the front lamps. Also, add a feature that allows the car to explode when the user presses the "space" button, and a "game over" message should be seen on the screen. There's no need to add any buttons or functionalities in the "game over" screen. Here is my code:</p> <p>HTML:</p> <pre><!DOCTYPE html> <html lang="en"> <head> <meta charset="UTF-8"> <meta name="viewport" content="width=device-width, initial-scale=1.0"> <title>Move Car with WASD Keys</title> <style> body { margin: 0; overflow: hidden; background-color: lightgray; } #gameArea { width: 100vw; height: 100vh; position: relative; } #car { width: 50px; height: 100px; background-color: red; position: absolute; top: 50%; left: 50%; transform: translate(-50%, -50%); } </style> </head> <body></pre>	<p>[Prompt would have been better if contributor had given the model an 'index.html' and/or some more specific instructions about how the contributor planned to setup and execute the code. Additionally, this is a strange thing to ask for; the contributor essentially having the model build a back end on the front end. In terms of prompts, this is a rather unrealistic/unfeasible scenario for a model to address.]</p> <p>I am working on upgrading the Task Manager class in JavaScript into a sophisticated task management application. The following enhancements are required: first, implement a User model to manage multiple users, allowing each user to have their own set of tasks. Incorporate a method for adding users, ensuring tasks are associated with the specific user. Next, modify the Task model to include a tags property that allows tasks to be categorized by multiple tags. The addTask method should ensure validation for the task. Include methods for filtering tasks based on priority and sorting tasks by due date. Finally, local file storage for persistent task management should be implemented to save and load tasks automatically from the file. It should not duplicate users and tasks. Users should be loaded into the users' list, and tasks should be loaded into the tasks list by default.</p> <p>Here is the code:</p> <pre>// Task Model class Task { constructor(id, title, completed = false) { this.id = id; this.title = title; this.completed = completed; } } // Task Manager Class class TaskManager { constructor() { this.tasks = []; } addTask(task) { if (!task.id !task.title) { throw new Error('Invalid task. Ensure it has an id and title.'); } } }</pre>
--	---	---

<pre> <div id="gameArea"> <div id="car"></div> </div> <script src="script.js"></script> </body> </html> </pre> <p>JavaScript:</p> <pre> window.onload = function() { const car = document.getElementById('car') ; let carX = window.innerWidth / 2 - 25; let carY = window.innerHeight / 2 - 50; const speed = 10; car.style.left = `\${carX}px`; car.style.top = `\${carY}px`; document.addEventListener('key down', (event) => { switch (event.key.toLowerCase()) { case 'w': carY -= speed; break; case 's': carY += speed; break; case 'a': carX -= speed; break; case 'd': carX += speed; break; } carX = Math.max(0, Math.min(carX, window.innerWidth - 50)); carY = Math.max(0, Math.min(carY, window.innerHeight - 100)); car.style.left = `\${carX}px`; car.style.top = `\${carY}px`; }); }; </pre>	<pre> } this.tasks.push(task); } removeTask(taskId) { const index = this.tasks.findIndex(task => task.id === taskId); if (index === -1) throw new Error('Task ID not found.'); this.tasks.splice(index, 1); } getTasks() { return this.tasks; } // Sample Usage const taskManager = new TaskManager(); const task1 = new Task(1, "Learn JavaScript"); const task2 = new Task(2, "Refactor Task Manager"); taskManager.addTask(task1); taskManager.addTask(task2); console.log(taskManager.getTasks()); </pre>
---	--

Code Ecosystem	<p>I am a solutions architect at a tech company and I need an auto-scaling mechanism for a Python-based web application running in Docker containers, deployed across a hybrid cloud infrastructure (Microsoft Azure, AWS, and GCP). The system should scale based on real-time CPU load. The auto-scaling mechanism should use a custom load monitoring tool that will run within the containers and report the CPU load directly to a central monitoring system. When the CPU load exceeds a certain threshold, the system should automatically scale across all three cloud providers. Explain how this monitoring tool should be designed and how the communication between containers running in different cloud environments would take place even when affected by network inconsistencies and varying latencies while maintaining high availability and load balancing and considering data consistency, fault tolerance, and potential for cross-cloud issues. Explain each step.</p>	
-----------------------	--	--