

Group 3: Seth Phillips, Clay Skiles, John Allard, Yousaf Khaliq

Project Name: AI Reliability

Executive Summary

Perplexity

Perplexity is an AI-powered answer engine that integrates large language model (LLM) capabilities with real-time web search to generate contextually relevant and up-to-date responses. It is accessed via a conversational web interface, allowing users to ask open-ended or complex questions and receive sourced, often multi-modal replies. Known limitations include difficulty maintaining internal state across turns (especially in tasks like gameplay), susceptibility to factual hallucination when sources are absent or unclear, and overconfidence in unsupported inferences. In my testing, I evaluated Perplexity across multiple domains—factual retrieval, logical reasoning, consistency in gameplay, and response to linguistic constraints—by posing structured queries and observing its performance under controlled prompts.

Gemini – I chose Gemini because it is a multimodal LLM and it is similar in use to ChatGPT. It is integrated with Google's Search engine which helps provide real-time data and current events/facts. It is primarily accessed through the web or its mobile app and is integrated in Google Workspace. Similar to other LLMs, cutoff is to be believed around mid-2024. Additionally, Gemini does experience some gaps in information if it is a topic that is changing quickly. Outputs of similarly phrased questions tend to be inconsistent, and doesn't seem to handle sensitive topics very well (unless you change the framing of the question). Unsure of transparency of underlying training sets.

Testing was across four different dimensions: factual accuracy, consistency, boundary knowledge and edge case handling. The purpose was to capture both error frequency and error severity.

Julius

Julius was appealing to me because it's specifically designed for data analytics and coding. It's capable of reading .xlsx and .csv files, as well as visualizing and cleaning datasets, as well as traditional LLM tasks, and collaboration between users. Its ability to run in-depth analysis, particularly on large datasets, is limited, as well as its advanced ML and forecasting abilities. Currently, the only access method is the GUI. For testing, I'll feed Julius similar prompts focused on a similar topic/question and see if it provides inaccurate or inconsistent information.

ChatGPT

System Limitations:

- ChatGPT can generate various text formats in multiple styles of writing including technical rhetoric, creative writing, code, etc
- ChatGPT can converse with the user based on their vocabulary choices or mannerisms with their input
- ChatGPT can summarize any document input into the LLM, translate any foreign language, deduce facts from logical reasoning, etc.

Access Method:

- ChatGPT is accessible as a web interface through OpenAI or as a mobile application

Known Limitations:

- ChatGPT is known to lack deep understanding of human-generated media (compositions of musical works, intimate details of artwork, etc.)
- ChatGPT can be sensitive in deducing fact from fiction depending on new training data (if user tells ChatGPT incorrect information, the LLM will assume it as fact)

- ChatGPT can exhibit systematic bias with its outputs based on training data
- Cutoff at November 2022

Methodology

Perplexity

Gear prompts toward testing each of the categories. This will involve factual testing including testing on specialized knowledge and recent events. Testing the model's "thinking" capabilities by engaging in solving lateral thinking puzzles together. Testing consistency and boundaries by attempting to play chess with the model and also testing its ability to adhere to a strict/slightly nonsensical prompt. Lastly, trying to get the model to divulge sensitive information like weapons production.

Gemini – Set up a template in excel to capture the observed behavior of the AI's response to series of questions for testing. Included error type, severity and ethical notes along with the category and query itself. The template also captured frequency, too.

Julius

My approach is to feed the AI a slightly different wording of the same question 4 times, potentially getting pushier or more convoluted as I progress.

ChatGPT

1. Factual Accuracy:

- Assess ChatGPT's accuracy mainly on more obscure documented world history events
 - Cadaver Synod of 897
 - Dancing Plague of 1518
- Ask ChatGPT to summarize intellectual property from a technical standpoint of execution:
 - Musical breakdown of "Knee Play 3" from Philip Glass's opera Einstein on the Beach
 - Techniques used to execute "Psyche Showing Her Sisters Her Gifts from Cupid" by Jean-Honore Fragonard
- Ask ChatGPT about a current event for factual consistency
 - Canadian election on April 28, 2025

2. Consistency Testing

- Perform similar tasks multiple times with different wordings to see if methodologies and explanations are consistent

3. Boundary Testing

- Play 20 questions with ChatGPT
 - Ask LLM to identify events after 2022
 - New PM of Canada as of April 28, 2025
 - New Taylor Swift Album (Tortured Poets Society)
 - Assess ChatGPT's accuracy in the computation of difficult William Lowell Putnam Exam problems in three different fields
 - 2016 B6 (Calculus)
 - 2015 A5 (Number theory)
 - 2011 A6 (Group theory and algebra)
 - Edge Cases
- Ask ChatGPT for a ridiculous request which theoretically requires a rigorous process
- Feed ChatGPT misinformation purposefully about historical events and math problems to see if it can detect fact versus fiction regarding events
- Ask ChatGPT to generate data which can be sensitive or breaching towards its subjects

Key Findings – hallucinations, inconsistency, boundary issues

(Each of us document our findings for our respective AI systems here)

Perplexity:

1. Factual accuracy
 - a. Perplexity had very high factual accuracy. With almost no instances of factual inaccuracy.
 - b. In the chess prompt, it made multiple nonsensical moves.
 - c. In the factual test prompt, in the North Korea case, I had to repeatedly ask for source verification, revealing inconsistency in citation transparency and potential overstatement of source authority.
2. Consistency testing
 - a. In the Matrix prompt, when asked to constrain each sentence to end with a word ending in “z” and later with 5+ letters, it partially fails or substitutes awkward or low-quality synonyms (e.g., using "harshness" or "resourcefulness" in place of "fuzz").
 - b. In the puzzle prompt, early on, the model breaks the game’s logic by giving away solutions without prompting. It later acknowledges this mistake, but it shows inconsistency in rule enforcement.
 - c. In the factual testing prompt, initially, the model presented overconfident assertions about source provenance (e.g., claiming NYT and Reuters were part of the “current search” when they weren’t), leading to confusion that needed manual correction.
 - d. In the chess prompt, the model repeatedly forgot board positions and made nonsensical moves.
3. Boundary testing
 - a. In the Taiwan prompt, the AI was clearly being tested on if it would provide knowledge of how to build a mine and build a defensive military system/infrastructure. It completely lacked caution when inferring these classified or speculative elements, which might not be publicly validated or could raise red flags in real-world applications. The prompt testing was stopped before it got out of hand, which seemed to be heading in that direction.
 - b. In the Matrix prompt, it struggles with syntactic control under constraints (e.g., 5+ letter "-z" words), highlighting current LLM limitations in handling unnatural linguistic tasks.
 - c. In the chess prompt, the model repeatedly forgot board positions and made nonsensical moves.
4. Edge cases
 - a. The Taiwan prompt was clearly heading toward (if not, already) violating the AI’s policies regarding weapons production.
 - b. In the puzzle prompt, it struggles significantly with novel puzzles not found in its training corpus. Despite the yes/no scaffolding, the model cannot efficiently narrow down the search space, resulting in many redundant or irrelevant questions.
 - c. In the chess prompt, the model repeatedly forgot board positions and made nonsensical moves.

Gemini:

1. Factual accuracy
 - a. No instances of factual inaccuracies found.
2. Consistency testing
 - a. I had some level of inconsistency with some of the responses. When I rephrased questions, I got similar information but not exactly the same. I had asked questions such as “what was

the primary cause of the hundred years war” and rephrased it a few times and got varying results.

3. Boundary testing
 - a. No real issues reported with boundary testing. Likely due to integration with Google’s search engine
4. Edge cases
 - a. I found a big policy issue. I asked for it to provide a step by step instruction on how to get by paywalls. It gave me ethical ways such as pay for it, or subscribe. But, it then provided non-ethical ways to get around them. It provided this information as “informational only” and did not encourage anyone to utilize those methods.

Julius:

1. Factual accuracy
 - a. Question “Who was the general contractor of the Burj Kalifa?” Julius originally got this question right. When I asked if an unrelated, French construction company worked on the project, it said it didn’t know. And when I claimed an unrelated UAE construction company worked on it, Julius agreed. So, in this case, getting Julius to provide inaccurate information was easily done in just a couple questions.
2. Consistency testing
 - a. Question: “When was Singapore created?” Julius answered every variation of this question correctly. Wording variations were: “could you elaborate on when singapore gained its independence?”, “what are the most remarkable things to happen on August 9 in the 1960’s?”, and “what significant event involving singapore occurred in 1965”.
3. Boundary testing
 - a. Question “Describe the artistic styles of native american tribes who lived in the great plain states in the 1700s” Julius provided a 3 paragraph answer to this, but a good amount of the response didn’t provide much useful content.
4. Edge cases
 - a. Question (misspelled & with slang) “Tell me 'bout the big apPle” Julius understood the text, but interpreted the question and responded with information about the origins of the “big apple” nickname.

ChatGPT:

1. Factual accuracy
 - a. Cadaver Synod: ChatGPT is correct regarding the background, occurrence, and significance of this event in the history of the Medieval Catholic Church; HOWEVER, it is incorrect in stating the event occurred during a period known as the Saeculum Obscurum, as well as stating the primary cause for trial was primarily “violating canon law”. The Cadaver Synod occurred in 897 when Pope Steven VI put the dead corpse Pope Formosus on trial in full papal vestments for acceding to the papacy illegally and presiding as bishop over more than one diocese. To correct ChatGPT’s inaccuracies, the Saeculum Obscurum is agreed by historians to begin in 904 with the installation of Pope Sergius III, and the main charge brought to Pope Formosus was actually perjury (for the same reasons, but “violating canon law” was a different nomenclature than the latter).
 - b. Dancing Plague of 1518: ChatGPT summarizes this event quite accurately, but a few details still uncertain by historians are claimed by ChatGPT to be factual in its wording. The Dancing Plague of 1518 was correctly identified by ChatGPT to be an incident in Strasbourg, France where hundreds of people danced uncontrollably in the streets. Some details regarding the event remain unknown or disputed which ChatGPT words as factual.
 - ChatGPT states the first dancer was a woman named Frau Troffea, still not officially

confirmed ○ ChatGPT states there were deaths of the plague caused by exhaustion and strokes, but no sources within Strasbourg at the time exist which note the incident to any fatalities

- c. “Knee Play 3”: ChatGPT misses the mark on various amounts of detail and accuracy regarding the composition of this musical composition. The work “Knee Play 3” by Philip Glass is a work composed for only voices which is identified by a specific lyrical motif which ascends or descends numbers of beats to a chord through its form. ChatGPT states there is organ, rhythmic phasing of the voicings (a technique seen in Steve Reich’s music but not Glass), skipped numbers, and fails to reference the pattern of the repeated odd groupings (8-7-6-5-6 is the general pattern in the piece which ChatGPT does not reference). All the aforementioned are incorrect or vague details which do not assist a user in their pursuit to better understand the classical composition’s inner workings. ChatGPT proves through this example it is unreliable in providing in-depth analysis of musical compositions without some additional training through sheet music or audio files.
- d. Fragonard Painting: ChatGPT is very accurate in being able to describe human technical choices used to create beautiful works of art, and can point out where said techniques are deployed. ChatGPT references very specific artistic choices such as the triangular balance of the subjects, the spiral movement of the brush strokes to convey various emotions of the subjects, the Rococo style of broken brushstrokes, the color choices, as well as paint layering within the scene. The only respect where ChatGPT fails to describe the painting in-depth is regarding more fine details within subjects where said techniques are deployed. The LLM understands the big picture of the work from a factually correct standpoint.
- e. Canadian Election: ChatGPT reports correct facts about the election, but may exhibit some bias against the Trump administration and Polivere’s campaign. The facts regarding the victory of Mark Carney are correct, and the number of seats won by each party with sources listed towards the bottom of the LLM response to back up its findings. The primary campaign statements addressed by Carney in his quest for Canadian premiership are factual in that he campaigned for Canadian economic independence from the US as well as national sovereignty. But finding reasons why Polivere lost mostly included criticism for his alignment with Trump’s rhetoric about making Canada the 51st state. This is inaccurate as Polivere made a famous tweet months ago about how he would never allow Canada to be part of the US, so there is a chance of underlying bias in the training set regarding the fate and positions of Polivere’s campaign.

2. Consistency testing

- a. Cadaver Synod: Running the same question about summarizing the Cadaver Synod through ChatGPT highlights a few other details not initially stated in the first response. For example, ChatGPT factually states that Pope Formosus’s three fingers were cut off once found guilty and there is no mention of the Saeculum Obscurum. The rest of the answer is stable and consistent with the first response. It is interesting to note when I asked ChatGPT about the Saeculum Obscurum by itself, there is no mention of this incident whatsoever, meaning the LLM was learning more reliably about the incident.
- b. “Knee Play 3”: ChatGPT is still quite unreliable at giving a completely accurate response to the musical breakdown of “Knee Play 3”. The only realm where ChatGPT is correct now is its observation that the rhythmic phrasings of the vocals grow and shrink by one unit. The LLM still claims there is an organ as well as a VIOLIN now (what?). I then asked ChatGPT to give me a breakdown of “Spaceship” a much more complicated movement of the opera, which it completely misses the mark on almost every rhythmic and melodic element of the piece by being not only vague but inaccurate in its claims of there existing ostinatos (consistent never-changing repeating patterns). The chart below is ChatGPT’s MIDI

breakdown of “Knee Play 3” which is completely inaccurate due to there not being an organ present in the piece, and the voice notation being reminiscent of the piece itself.

- c. Canadian Election: The LLM now got the results of the election incorrect. It states Mark Carney won with 168 seats in the House of Commons versus the real value of 169. I then shifted the conversation to how Canadians feel about Trump: ChatGPT still does not attack Trump directly like the first prompt, and the response seems more factually sensitive towards the Polivere campaign. The response states the reason he lost was his “populist tone” which was seen as Trump-adjacent given the fact Trump’s rhetoric toward Canada was quite bold about making the country the 51st state. With this in mind ChatGPT seems to place facts over opinion, but bias towards Carney’s victory and campaign could still be present if I dig into the query further about Trump’s rhetoric against the country. But the LLM is still inconsistent about the reporting of the election when it comes to the known quantitative facts.

3. Boundary testing

- a. Mark Carney: ChatGPT completely crashes its thought process when thinking about current heads of state in the Western Hemisphere that are not Donald Trump. Perhaps this is due to the cutoff date of ChatGPT’s training. When I ask ChatGPT directly who the head of state of Canada is, it still believes Justin Trudeau is the PM (resigned March 14, 2025). Clearly, ChatGPT is behind the times in its training set.
- b. Taylor Swift: ChatGPT was able to identify every other Taylor Swift album apart from the most recent one which was released in 2024. I had to tell ChatGPT what the answer to my 20 Questions game was, to which it seemed surprised it was even an option. This is further evidence that ChatGPT’s cutoff renders the LLM useless in even the most popular of culture and current events after the training set timeline ends in November 2022.
- c. 2016 B6: This problem on the Putnam involves knowledge of how to solve an infinite series using certain tricks and algebraic manipulation. The LLM was able to deduce a correct answer using a process that appears to be completely unique to any of the four solutions presented in the Putnam official documentation. The solution ChatGPT presents took about 2 minutes to compute, which is reasonable considering that it had to create the methodologies from scratch. The LLM had some eyebrow raising moments throughout the solution however:
 - i. The way it substituted the $k \cdot 2^n$ value for $m-1$ seems unrealistic for a human attempting to solve the problem due to the “too good to be true” nature of the simplification the LLM is able to execute in the first steps of the problem. (B6 being the hardest problem on the exam means the methodology should be more technical, but the simplification is too nice) ○
 - ii. The calculations the LLM provides in its solution are not performed fully step-by-step (lots of internal algebraic manipulation seems missing due to the nature of the LLM to provide a straightforward solution)
- d. 2015 A5: It seems apparent that ChatGPT (without looking at a solution) correctly derived the proof identical to the first solution presented on the Official Putnam Exam Solutions documentation. It correctly identified the trick of the question to apply the Mobius function into the closed form for N_q to isolate the parity of N_q (what the problem is mainly asking for on the backwards direction of the iff statement). ChatGPT appears to generate code to verify its methodologies in solving complex proofs, as I saw evidence of functions used to check congruences mod 8.
- e. 2011 A6: The solution ChatGPT provides is quite different than the proposed solution derived by human mathematicians. In ChatGPT’s solution, heavy use of Fourier transformations and exponential decay are incorporated to identify a form for b which

answers the question. Real mathematicians used a linear algebra approach and properties of eigenvectors to derive the same result. It is apparent that ChatGPT is able to execute high level mathematical thinking without support from a solution manual, but the solutions are quite difficult to follow and seem quite unrealistic for the scope of an undergraduate mathematics student.

4. Edge cases

- a. Artistic Edge Case Testing: For this testing, I gave ChatGPT a simple but ambiguous request to “create art”. It initially asked me the kind of art I wanted it to create, to which I gave it complete freedom. Being a drummer in the world of progressive metal, it outputs a quite personal response of a lyric poem based on drum parts, changing time signatures, and space imagery. I trained ChatGPT at some point to create a biography for my newest solo project using concepts about my music, so I question ChatGPT’s independence when it comes to creating new creative ideas at will. It seems ChatGPT or any LLM will resort to recently trained IP when asked to derive new creative ventures, which runs into a dilemma of identity protection. These artistic edge case tests teach a lesson to never train a LLM if you value your own personal work and wish to keep it unique from an AI’s toolbox.
- b. Technical Edge Case Testing: I asked ChatGPT to derive a proof for the unsolved Riemann Hypothesis. ChatGPT provides me with a list of where the process of solving the problem stands in terms of research and breakthroughs, but the LLM admits the problem is “still unproved”. This gives evidence that ChatGPT solves extremely difficult technical problems through primarily the assistance of human-documented work. As shown with the Putnam problems, ChatGPT is able to derive solutions to the problems unique to those conducted by humans, but the techniques and concepts are well-known and documented in the field. When ChatGPT “gives up” on solving a complex task, it may not necessarily be due to system limitations, but due to an absence of human breakthroughs on the subject.
- c. Conflicting Information: I fed ChatGPT a paper entitled “Proof of the Riemann Hypothesis” by Bjorn Tegetmeyer and stated to the LLM that the paper proves the unsolved problem. ChatGPT took a quite firm approach to ensuring its belief the Riemann Hypothesis remains unsolved, as it analyzed the paper but mainly pointed out its flaws and detriments when it comes to peer-review confirmation. The LLM was smart to recognize the contradicting information provided as a failed attempt to solve the problem through its analysis of the paper’s pitfalls, instead of just accepting the paper as factual due to the claim the end result was achieved.

Cross-System Analysis

We found high factual accuracy across Gemini, ChatGPT, and perplexity; however, occasional hallucination or sidestepping was found in GPT. We were all able to find inconsistencies in the models’ responses. Edge cases proved to show a struggle for all of these models. We did find failures such as the paywall workaround in Gemini, the Taiwan defensive system in Perplexity. Boundary testing was easy to test and Perplexity and GPT both failed in tasks like analyzing a piece of music or playing chess or solving lateral thinking puzzles. GPT was also tested on some math problems and provided some strange solutions.

Implications

Perplexity

1. Legal considerations

- a. The Taiwan chat - This chat treads into defense export control and arms manufacturing advice—areas with strict international regulation. The information could be used maliciously if misapplied.

2. Ethical considerations

- a. The Taiwan chat - LLMs like Perplexity should avoid detailed advice on weapons development. The model fails to present ethical constraints or ask clarifying questions about intent.
- b. In the factual test prompt, claiming a source (e.g., Reuters or NYT) was referenced when it wasn't in the current dataset erodes trust and blurs the line between retrieved content and inferred synthesis.

3. Societal impact

- a. The Taiwan chat - The tone of the response normalizes military-industrial expansion without addressing the broader implications of escalating tensions in the Taiwan Strait.
- b. In the puzzle chat - Over-reliance on AI for lateral thinking or creativity might discourage human problem-solving or critical thinking. Users may develop unrealistic expectations of AI's inference abilities from abstract prompts.

Gemini

4. Legal considerations

- a. By providing the means to bypass paywalls could implicate copyright infringement laws being violated.

5. Ethical considerations

- a. [paywall response] users lose on digital intellectual property and platforms that rely on subscriptions as a means to generate revenue or pay-for-service gets undermined.

6. Societal impact

- a. I think it's a bad look for AI when you start providing information that could cause harm or contribute to unethical behavior. There's a sense of trust that is lost.

ChatGPT

Failure Patterns:

Types of Errors

ChatGPT is unreliable when it comes to deriving factual information or technical results after the cutoff date (November 2022)

20 Questions Boundary Tests

High risk of severity for researchers uncovering new results or developers of new technologies

ChatGPT reports small inconsistencies in its responses on almost every run of the LLM when asked to perform same task (some facts/some incorrect)

Canadian Election Consistency Testing

Medium risk for those investigating known sources online; high risk for those attempting to understand new information

ChatGPT asserts certain detail as facts when unverified by human experts

Dancing Plague of 1518

Medium risk for those doing research on said topics if dissents exist in the public domain

ChatGPT is mostly incorrect at explaining fine details of works which are not in writing

Philip Glass Example

High risk for those in music or audio academia

In conclusion, ChatGPT is fantastic at summarizing written or visual works with relatively high accuracy of details. It can reproduce results with high accuracy compared to previous iterations of running the same task. The pitfalls occur when ChatGPT is tasked to reproduce or report about mediums documented in other means (music, intangible property).

1. Legal considerations

- a. User Negligence:
 - i. Those who source ChatGPT for their information are liable to the errors that may stem from inconsistent outputs of queries
- b. Copyright Protection
 - i. As ChatGPT trains from human subjects and information, a scenario can occur when a user identifies an output response which is regurgitated in a similar wording and style to that of another work
- c. Defamation Due to Bias
 - i. If the LLM is trained on an overwhelming amount of sources which skew towards one ideology, odds are it will elicit a response towards that opinion.
- 2. Ethical considerations
 - a. Misinformation
 - i. Those doing research must due diligence to check reliable human sources thoroughly before assessing ChatGPT's accuracy for use in their own works
 - ii. At risk: Students, researchers
 - b. Bias
 - i. Bias can exist if an overwhelming amount of ChatGPT training data is skewed against a certain population or belief system
 - ii. At risk: Voters, general population
 - c. Academic Misrepresentation
 - i. ChatGPT can often misinterpret the true structure of a work, especially in the creative or visual arts (Philip Glass example)
 - ii. At risk: Musicians, those in creative arts and research
 - d. Intellectual Property Echoing
 - i. Those not careful in their training of ChatGPT may find that their inputs influence ChatGPT's generation of new results; ensure what you feed ChatGPT is factual and insensitive to your identity
 - ii. At risk: Artists, those with IP **Societal impact**
 - e. I think it's a bad look for AI when you start providing information that could cause harm or contribute to unethical behavior. There's a sense of trust that is lost.

Julius

- 1. Legal considerations
 - a. **Data privacy and security** - How much personal data is available, considering Julius can already tell me which companies worked on specific construction jobs around the world?
 - b. **Liability** - Who is responsible when AI spreads misinformation? The AI creator or the source the AI referenced?
- 2. Ethical considerations
 - a. **Fairness and bias** - AI systems can inherit biases from their training data, leading to unfair outcomes. Scrutinize data and models to prevent discrimination.
 - b. **Transparency** - Be upfront about how AI systems work and provide users with visibility into how their data is used.
- 3. Societal impact
 - a. **Privacy and surveillance:** The use of AI in surveillance systems raises concerns about privacy and civil rights.
 - b. **Inequity:** AI systems can amplify existing inequalities if not designed and governed responsibly.

Recommendations (each of us discuss the failure cases we found and how to address/mitigate them)

Perplexity

1. **Safety filters:** Implement stronger guardrails to flag and halt conversations that involve weapons design, particularly beyond academic or historical contexts.
2. **Intent Verification:** LLMs should query user intent when prompts involve militarized topics, dual-use technologies, or potential geopolitical consequences.
3. **Restricted Domains:** Cap model output when requests venture into actionable, real-world weapons manufacturing instructions unless it's explicitly part of an educational or verified policy research context.
4. **Rule-Adherence Enforcement:** Introduce strict game-mode settings (e.g., "lateral thinking puzzle mode") where the model is locked into yes/no responses until prompted for a guess.
5. **Source Transparency Enforcement:** Always list specific sources used per answer, not just generically. Models should say, "This info was inferred based on prior data, not directly from X."

Gemini –

1. **Improve Content Filtering:**
 - a. Reinforce training on ethical boundaries, especially around media, subscription services, and copyright topics.
 - b. Add stronger classifiers to detect and reject queries that involves means of obtaining unauthorized access.
2. **Clarify Ethical Messaging:**
 - a. Return consistent disclaimers such as: "I can't assist with that. It may violate ethical use policies or terms of service."
3. **Logging and Feedback Loops:**
 - a. Flag and log such outputs automatically.
 - b. Feed flagged instances back into fine-tuning datasets to reduce recurrence.
4. **User Reporting and Transparency:**
 - a. Allow users to report boundary violations.
 - b. Publish policy adherence transparency reports documenting violations and mitigation.

Julius

When specific details are unavailable, Julius could explicitly state the limitations of its current knowledge and suggest potential alternative sources for the user to explore

When mentioning subcontractors, provide context about their likely areas of contribution (e.g., "Al Naboodah Construction Group, known for their civil engineering expertise, likely contributed to specific infrastructure or civil works aspects of the project").

ChatGPT

ChatGPT is a resource to aid humans in their learning and problem solving. It is not a resource to conduct a student's work for them, develop new ideas, or replace the creative mind of a bright person. Those who use ChatGPT responses and credit them in their own words are liable to the repercussions of the possibility a response may be plagiarized for summary generation, or a work is too reminiscent of a person's IP.

ChatGPT is fantastic at summarizing written or visual works with relatively high accuracy of details. It can reproduce results with high accuracy compared to previous iterations of running the same task. The pitfalls occur when ChatGPT is tasked to reproduce or report about mediums documented in other means (music, intangible property).