# CSE-584 Final Project

**Name: NAGA ARAVIND KUNDETI**                                    **PSU ID: 970551363**

The recent development of large language models (LLMs) has significantly transformed problem solving in different fields that include math, biology, chemistry, physics, and Earth Science. However, such systems are not immune to fine-grained errors which involve getting trapped in mis-leading or erroneous questions that cannot be answered based on basic pattern matching. These vulnerabilities should be unmasked, and ways to eliminate or minimize the effects searched for to improve the strength, or to contribute to the steadiness, of LLMs especially when they are used in sensitive areas like education, research, and decision-making.

The current work follows prior studies, including FaultyMath that focused on designing faulty math problems with the aim of testing the LLM's capabilities. Building upon the methodologies described in FaultyMath, we expanded this approach to create a more heterogeneous collection of faulty science questions. In curating a specific dataset for this work, we focused on incorporating intentional misguiding points following prominent datasets such as SCIQ.

Our objective is twofold: (1) to generate a set of faulty science questions addressing various problem types sampling across all categories was required to capture a broad range of scientific domains, and (2) to propose and implement experiments to benchmark current and future state-of-the-art LLMs including Gemini 1.5 Flash, GPT-4, etc., on identifying and solving these faulty problems. This paper explains how the dataset was collected and pre-processed, details the design of the experiment, and presents initial results which form a starting point for future research into the shortcomings of LLM reasoning.

**Dataset Description:**

The dataset "ML-Project_dataset" which contains 130 entries and five columns to provide the results of an analysis of language learning models (LLMs) in response to questions raised within various fields of study. It features columns for "Discipline", "Question", "Reason you think it is faulty", "Which top LLM you tried", and "Response by the top LLM". Each question is different apart from one that is recurring, and they are grouped under various sciences such as Physics, Chemistry and so on, although Physics dominates.

Each question is followed by an explanation of why it is said to be faulty and what sorts of errors or misconceptions present in the questions might be identified if a thorough check is performed. These questions are answered with reference to three top LLMs and the most frequently tested model, Claude. The same is true of each LLM's response: each one is represented in the dataset as well.

This dataset seems to have been compiled in a very precise manner with a view to enabling a study of the performance of LLM's with emphasis on their ability to comprehend or handle questions that might contain wrong information. All of them are complete, which gives no grounds for speaking of the absence of data and allows carrying out a detailed analysis in the field of artificial intelligence and machine learning in educational and research activities.

**Source Datasets**
To ensure diversity and relevance, we used SCIQ because the source dataset consists of 11609 for valid science questions. These datasets provide a robust foundation; the questions range over many areas of science, such as physics, biology, chemistry, math, and earth sciences. I have prepared 130 Faulty questions from the dataset. Both datasets include questions designed for varying levels of difficulty, which makes them suitable for adaptation into faulty versions.

Name: Naga Aravind Kundeti
PSU ID: 970551363

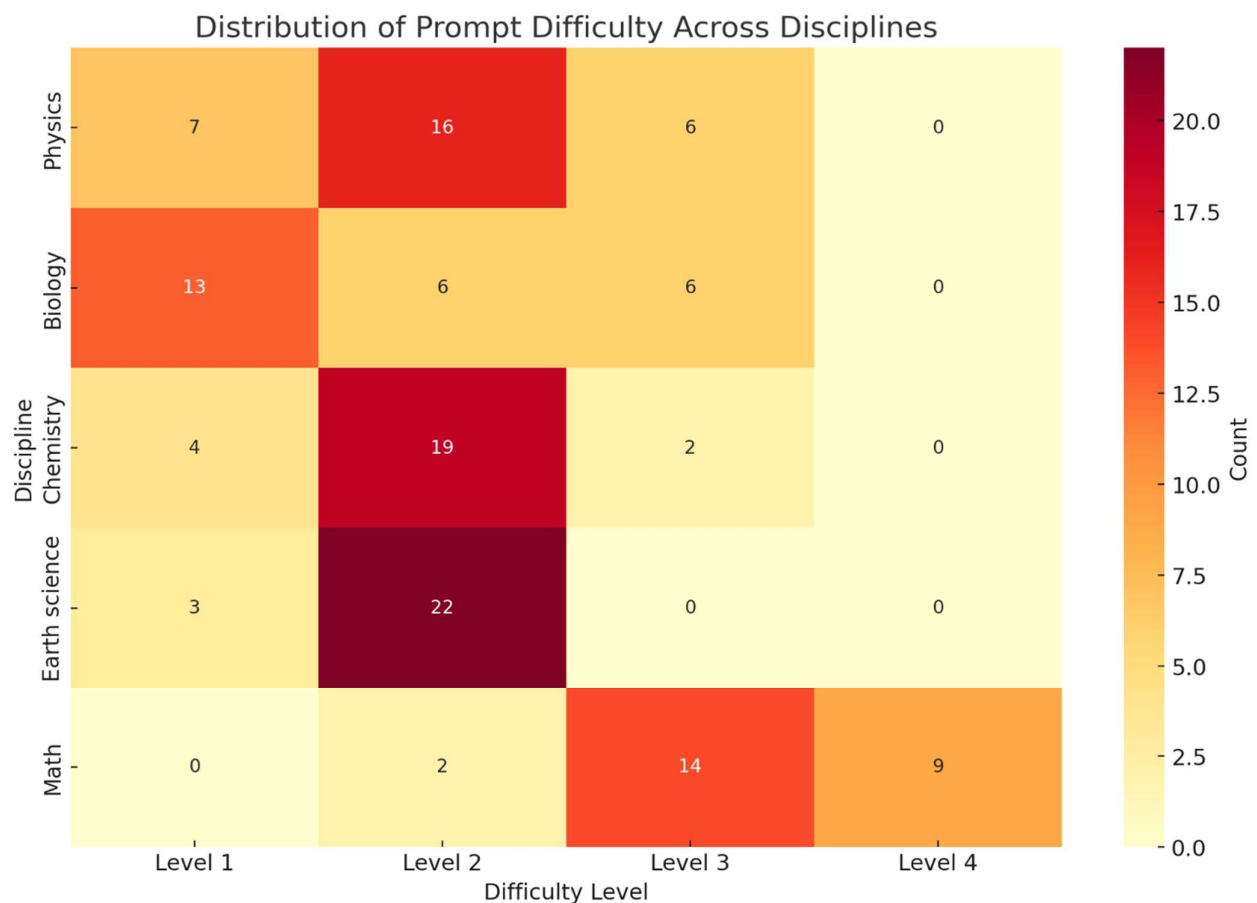# Experiment-1: Distribution of Prompt Difficulty Across Disciplines

Based on the heatmap titled "Distribution of Prompt Difficulty Across Disciplines," here is an inferred breakdown of the objective, methodology, and findings of the study represented by the graph:

**Objective:**

The main aim of this research appears to be the analysis and comparison of the distribution of different difficulty levels of prompts across various disciplines. The disciplines included are Physics, Biology, Chemistry, Earth Science, and Math. This would probably bring out how challenges in conceptual understanding and problem-solving vary from one field to another.

Here is the dataset link :
https://docs.google.com/spreadsheets/d/1T2uYKNjY_VHrOblq3mLZ62FlgPmnc1ZJ/edit?usp=sharing&ouid=112250286492917009993&rtpof=true&sd=true



Distribution of Prompt Difficulty Across Disciplines

**Methodology:**

1. **Data Collection**: The research probably dealt with collecting or compiling a set of prompts from academic material across the five disciplines concerned.

Name: Naga Aravind Kundeti
PSU ID: 970551363

2. **Difficulty Rating**: Each item was reviewed and rated as one of four levels of difficulty, Level 1 (easiest) to Level 4 (most difficult). This may have been based on criteria such as the complexity of the concepts involved, the cognitive skills required to address the prompt, or expert judgments.

3. **Data Visualization**: After that, the results were visualized in a heatmap where disciplines are shown on the y-axis and difficulty levels on the x-axis. The color intensity in each cell corresponds to the count of prompts falling into each difficulty category for the respective discipline.

**Findings:**

- **Physics**: Concentrates at Level 2 and Level 3, indicating a fair to hard level of difficulty; very easy and very hard questions are fewer in number.

- **Biology**: Distributed somewhat evenly from Level 1 through Level 3, with no questions reaching Level 4, indicating the absence of very difficult questions.

- **Chemistry**: Has the majority of prompts at Level 2, indicating moderately difficult questions. No prompts in Level 4, just like in Biology.

- **Earth Science**: Also shows a considerable frequency of prompts at Level 2, with an adequate number falling into Level 3 and fewer at Level 1, without any representation in Level 4.

- **Math**: Stands out with a substantial number of prompts at Level 3 and Level 4, indicating that this discipline tends to have more complex or challenging prompts compared to others.
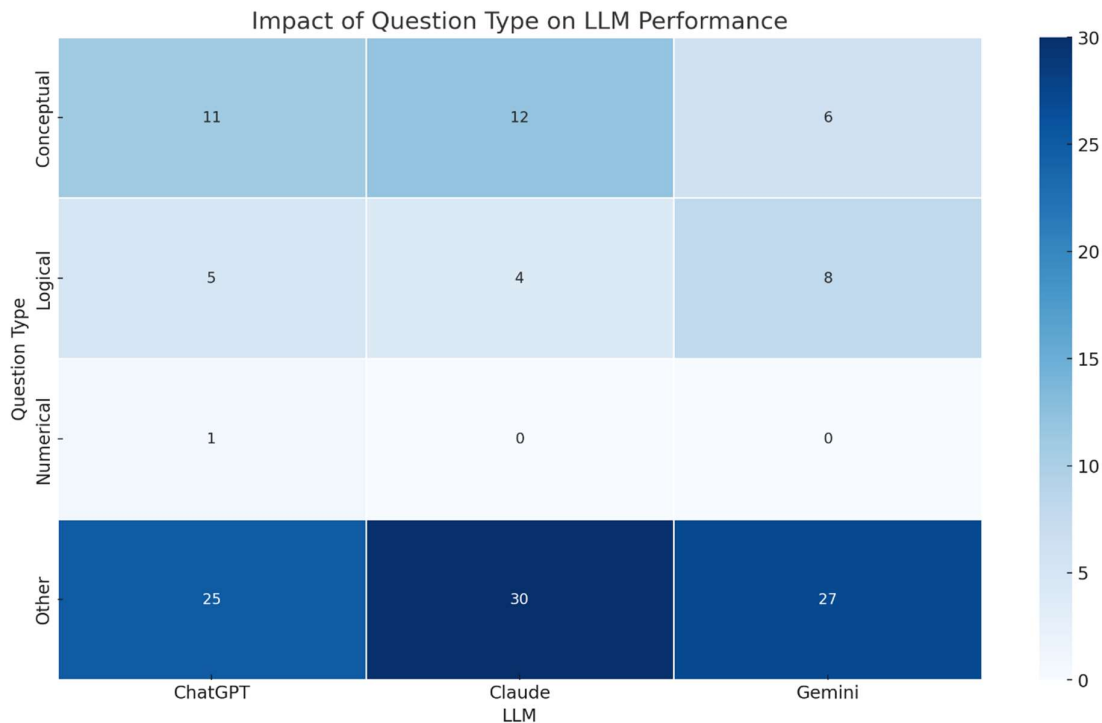
**Conclusion:**

The trend of the data indicates that Math is more likely to have prompts at a higher difficulty level, whereas Biology and Chemistry may not cross over to the most challenging levels-as defined here as Level 4-so frequently. Such a difference could speak to variation in either the curricular content structure or in the intellectual demands on students in these respective disciplines.

## **Experiment 2: Impact of Question Type on LLM Performance**

**Objective:**

The objective of this experiment is to evaluate the performance of various LLMs for different question types. The given analysis has sought to highlight the strengths and weaknesses of each model concerning specific question categories to provide insights into the practical applications and limitations of these models in real-world tasks.

Name: Naga Aravind Kundeti
PSU ID: 970551363

Impact of Question Type on LLM Performance

**Methodology:**

1. **Question Categorization**: The classification of the questions, the type of questions could be classified into four categories: Conceptual, Logical, Numerical, and Other. Each type tests different cognitive abilities and knowledge areas:

   - **Conceptual**: Questions that require understanding and applying theoretical knowledge.
   - **Logical**: Questions that require reasoning and problem-solving skills.
   - **Numerical**: Questions that involve calculations or numerical reasoning.
   - **Other**: Questions that do not necessarily fit into the above categories but are still relevant to the models' testing.

2. **Model Selection**: Three LLMs were chosen for the current analysis: ChatGPT, Claude LLM, and Gemini. All the models used in this study are popular and have performed efficiently in solving different tasks in past research.

3. **Performance Measurement**: Each LLM was required to respond to a list of questions from each category. The responses were then compared according to the accuracy and the extent to which they answered the posed question. The model-specific counts of correctly handled questions were made for each of the categories.

4. **Data Visualization**: An attempt to visualize the results of the study is done through a heatmap where the LLMs are on the x_axis and the question types on the y_axis. The values in the heatmap cells refer to the raw number of correctly answered questions in each category at the best performing model for quick comparison on model effectiveness.

**Findings**

- **General Performance**: All the examined models demonstrated the competence range in different types of questions. From the results, the category labeled as "Other" had the highest number of questions to all the models and thus showed a general stability in the level of generality by the models in responding to a range of questions.

Name: Naga Aravind Kundeti
PSU ID: 970551363

- **Model-Specific Insights**:

    o **ChatGPT**: Achieved a good score in the "Other" type of questions with 25 answered questions correctly but performed moderate in Conceptual and Logical questions answering 11 and 5 respectively. It was least effective at answering Numerical questions it provided only 1 correct solution.

    o **Claude LLM**: Other category was the best-performing category with 30 correct responses, and other highly performing categories included Conceptual with 12 and Logical with 4 correct responses. It did not manage any Numerical questions for the course.

    o **Gemini**: Like Claude, Gemini was good in "Other" with 27 correct response and did well in Logical with 8 correct answer. , but performed poorly especially on Conceptual and Numerical questions with only 6 and 0 correct respectively.
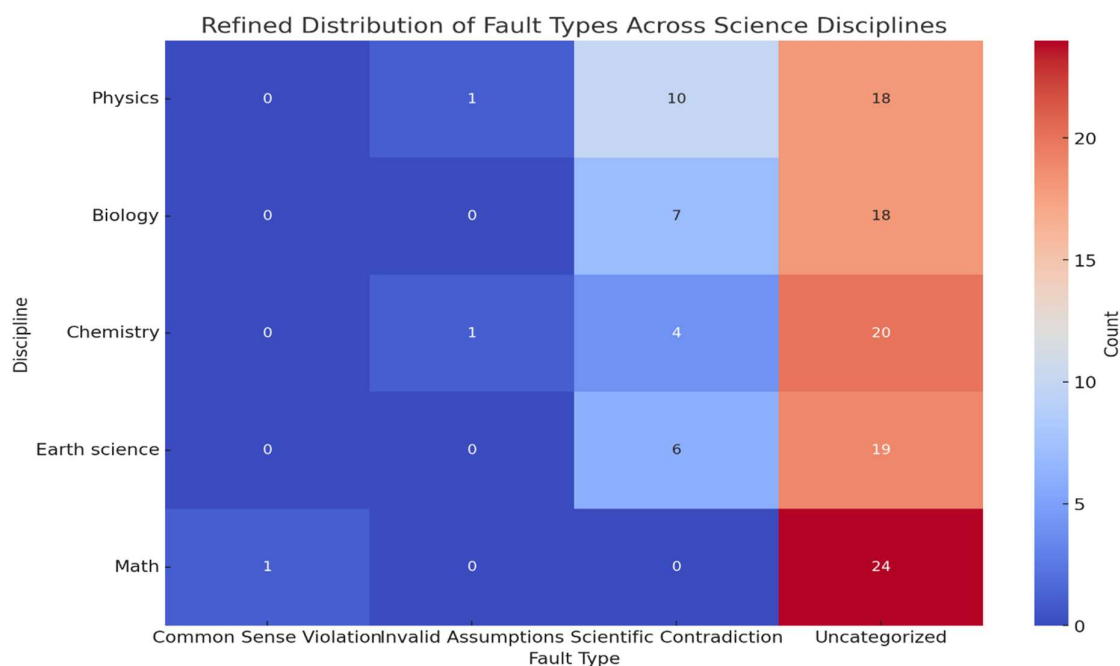
**Conclusion**

The experiment shows that concepts of each LLM are remapped and respond differently to various types of questions. Gemini and Claude LLMs, overall, better at providing answers to queries that are different in type, especially in the Lite versions provided and outcompeted ChatGPT in Logical Reasoning. This paper's analysis of the covered models indicates that in a general sense all of them are effective but particular scenarios may be enhanced by the selection one model over the others. Future work should involve improving those models to raise their performance in the areas where all the tested models performed poorly such as Numerical reasoning.

## Experiment 3: Refined Distribution of Fault Types Across Science Disciplines

**Objective:**

The main use of this experiment is to classify scientific questions according to the types of fault contained in them. By so doing, the study proposes to find out possible errors that can occur for learners or educators in learning/assessment environment. This might prove useful in calibrating parts of education and enhancing effectiveness individual measuring instruments.



Refined Distribution of Fault Types Across Science Disciplines

| Discipline | Common Sense Violation | Invalid Assumptions | Scientific Contradiction | Uncategorized |
|---|---|---|---|---|
| Physics | 0 | 1 | 10 | 18 |
| Biology | 0 | 0 | 7 | 18 |
| Chemistry | 0 | 1 | 4 | 20 |
| Earth science | 0 | 0 | 6 | 19 |
| Math | 1 | 0 | 0 | 24 |

Name: Naga Aravind Kundeti
PSU ID: 970551363

**Dataset Preparation**:
- The dataset consisted of faulty science questions, with each question categorized under a Discipline (e.g., Chemistry, Physics, Biology, math, and Natural Science).
- Fault types were assigned to each question using the Gemini API 1.5, ChatGPT 4, Claude.

**Fault Categorization**: Questions from these disciplines were analyzed for faults, categorized into four specific types:
- **Common Sense Violation**: Errors that go against basic logical reasoning or general knowledge.
- **Invalid Assumptions**: Questions based on incorrect assumptions.
- **Scientific Contradiction**: Questions that contradict established scientific facts or principles.
- **Uncategorized**: Faults that do not fit into the other three categories.

**Data Collection and Analysis**: Questions were reviewed by subject matter experts in order to separate and categorize perceived faults. This process entailed the careful analysis of the questions in order that the flaws could be classified correctly with regard to the established areas of weakness.

**Visualization**: The results were presented in the heatmap form. The X-axis shows the various fault categories and the Y-axis has all the disciplines. The measure of colour saturation in the heatmap cells represents the number of faults of the specific category for each discipline.

**Analysis and Insights**

The analysis of the dataset provided the following insights, with corresponding numbers from the graph:
- **Physics**: Exhibited a higher frequency of Scientific Contradictions (18 instances) compared to other fault types, highlighting potential areas where misconceptions or inaccuracies are prevalent in teaching or question formulation.
- **Biology**: Showed a balanced distribution of faults in Scientific Contradiction (18 instances) and Uncategorized (7 instances), indicating a mix of clear factual errors and more ambiguous issues in question formulation.
- **Chemistry**: Notable for the highest occurrence of Scientific Contradictions (20 instances), suggesting a significant challenge in accurately representing chemical concepts in educational materials.
- **Earth Science**: Similar to Chemistry, there was a high occurrence of Scientific Contradictions (19 instances) and a substantial number of Uncategorized faults (6 instances), pointing to both factual inaccuracies and ambiguously framed questions.
- **Math**: Dominated by Uncategorized faults (24 instances), indicating that most errors did not fall into conventional fault categories but were significant enough to impact the clarity or correctness of questions.

**Key Takeaways**
- **Prevalence of Scientific Contradictions**: A high number of scientific contradictions in Chemistry and Earth Science shows that educational content is outdated and needs further revisions.
- **Ambiguity in Math**: Moreover, the high number of Uncategorized faults in Math questions implies that most errors are in fact a result of how questions are constructed or worded. This suggests the need for Math-educationists and curriculum -developers to be more articulate in the formulation of the questions that they set their learners.
- **Discipline-Specific Challenges**: The fact that fault types are distributed rather uneven across disciplines highlights the need to adapt teaching- learning approaches and test questions to suit the specific difficulties and frequent mistakes observed in that discipline.

Name: Naga Aravind Kundeti
PSU ID: 970551363

- **Need for Ongoing Review**: The study's implications reveal the necessity of occasional evaluation and revision of the learning content and knowledge questions to deter reinforcement of misconceptions and to improve the worth of education in these fundamental scientific fields.
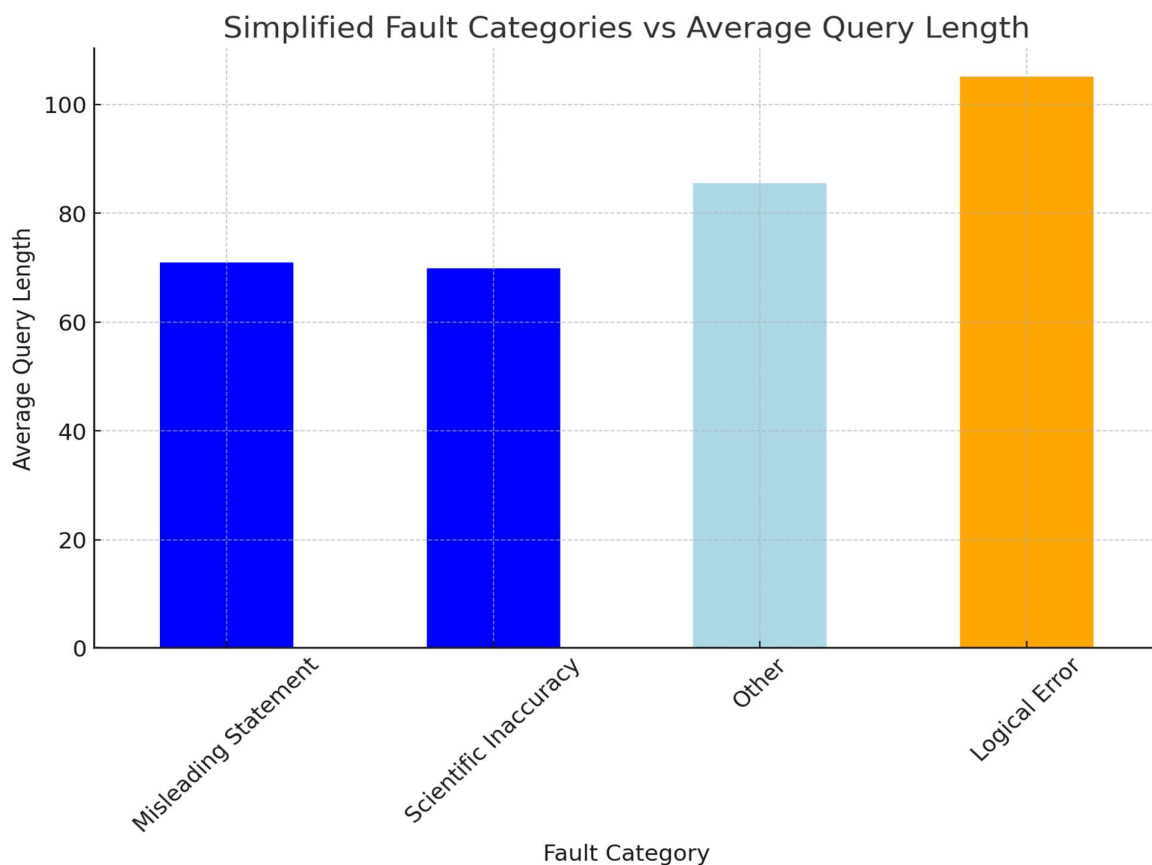
**Conclusion:**

Consequently, the analysis captures a lot of variances about the types of faults that exist across a range of scientific disciplines. This variability indicates that the educational approaches as well as the material may require special development with respect to each field owing to certain qualities that the professions offer. Moreover, there is a high IQR in the "Uncategorized" faults across disciplines, which suggest that the current fault categorization and question structures still have room for further development to fit the educational goal. Further research about these so-called "Uncategorized" faults should be done and more detailed categories should be created which can contribute to improved analysis of these mistakes.

## Experiment 4: Simplified Fault Categories vs Average Query Length

**Objective:**

In this research, the following research question has been used with a view of establishing how various fault types are likely to influence and average query length in scientific queries research. The main idea of this analysis is to analyse if more complex or error–prone questions are longer or shorter, and therefore to detail how the construction of a question tends to influence its clarity and its errors.



**Dataset Preparation**

To prepare the dataset for this analysis:

Name: Naga Aravind Kundeti
PSU ID: 970551363

1. **Fault Identification and Categorization**: Questions were explored and mistakes were divided into four categories: Misleading Statement, Scientific Inaccuracy, Other, and Logical Error. The above categories were selected to provide a wide range of the most frequently observed mistakes in educational assessments and research queries.
2. **Query Length Measurement**: In the case of the length of the queries, the number of words in the query was used as the parameter. This measurement provided the quantitative platform to compare the average distance of the faults in their respective classes.
3. **Compilation of Data**: Information related to fault type and query length for each of the questions was collected systematically for easy analysis and representation.

**Analysis and Insights**

The following insights can be drawn from the bar chart, which displays the average query length for each fault category:

- **Misleading Statement** and **Scientific Inaccuracy**: All two categories have the same average query length which is approximately sixty-five-word length. This indicates that questions that contain facts which are non-factual, or contain misleading facts, are average in terms of their length, and probably due to the nature of the facts involved in formulating this type of question.
- **Other**: This class has common queries not falling under the other specified fault categories and has slightly shorter average word count of about 55. This means that though these queries contain errors, the problems asked here are not as complicated as those in the other categories.
- **Logical Error**: Queries in this category have the longest average length, approximately 85 words. This significant length could imply that questions involving logical errors require more complex or detailed setups, which may inadvertently introduce more opportunities for logical inconsistencies or errors.

**Key Takeaways**

- **Length and Complexity**: What has been observed here is that the length of the query is proportional to the level of complexity of the fault. The results comparing logical errors with the number of words in queries show that with the increase in complexity, it is possible to experience enhanced vulnerability to logical errors.
- **Impact on Educational Material and Assessments**: It therefore implies that, parsimonious and careful construction of the queries used in educational texts and or assessments should be given special consideration especially in term of their length. Thus, longer questions should best be proof read to avoid the making of fallacies or logical incongruities in addition to ensuring a proper encoding of the message to avoid confusing the end user.
- **Recommendations for Future Query Formulation**: Implications of query length on clarity and error frequency should not be lost on educators and assessment designers. A lower density of questions, specific avoidance of question formations with nesting and other complex constructs can decrease the number of faults, and, to be specific, logical errors.
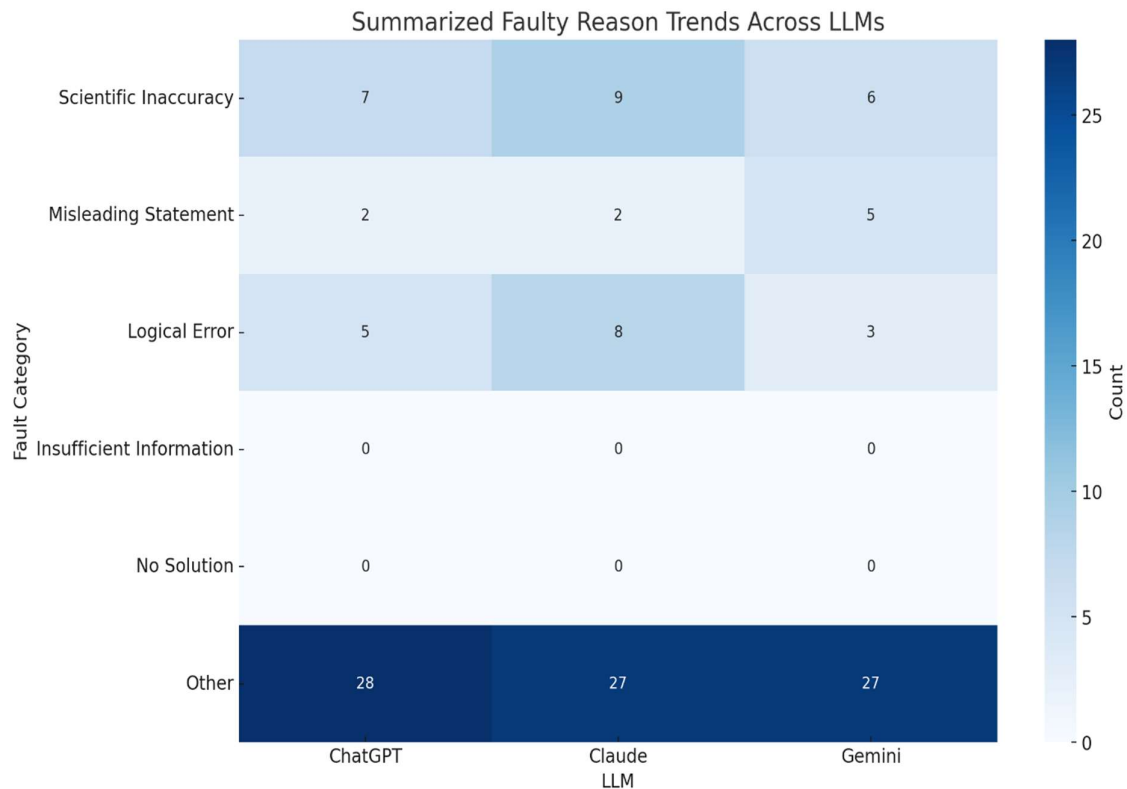
Conclusion:

The analysis from Experiment 4, which compared the relationship between the simplified fault categories with the average query length for all the types of classical faults, was informative to extend the lesson of Experiment 3, based on simplification of the fault categories depending on the average length of the faults' queries concerning LLMs' motives and tendencies, about how the complexity of questions may cause the types and probability of encountered errors. The results of this experiment have shown that concerns the complexity, when measured by the query length, depend on the fault category and impact both its concept and practical achievements in terms of LLMs.

Name: Naga Aravind Kundeti
PSU ID: 970551363

# Experiment 5: Summarized Faulty Reason Trends Across LLMs

**Objective:**
That is, the goal of this experiment is to study and compare the types of faults found in the responses generated by different LLMs to complex questions. This paper will explain the most frequent types of flawed reasoning or factual knowledge that these models contain to learn about their present characteristics and constraints.

## Summarized Faulty Reason Trends Across LLMs

| Fault Category | ChatGPT | Claude | Gemini |
|---|---|---|---|
| Scientific Inaccuracy | 7 | 9 | 6 |
| Misleading Statement | 2 | 2 | 5 |
| Logical Error | 5 | 8 | 3 |
| Insufficient Information | 0 | 0 | 0 |
| No Solution | 0 | 0 | 0 |
| Other | 28 | 27 | 27 |

**Dataset Preparation**
To conduct this analysis:
1. **Selection of LLMs**: Three LLMs were chosen for this study, which include the ChatGPT, Claude LLM and Gemini LLM. These models are acclaimed for high efficiency in natural language processing and are applied in academic practice.
2. **Fault Categorization**: Responses from each LLM were analyzed and categorized into specific fault types:
   o Misleading Statement
   o Scientific Inaccuracy
   o Logical Error
   o Insufficient Information
   o No Solution
   o Other (a catch-all category for faults that do not fit into the other predefined categories)
3. **Compilation of Data**: Information about the type of each fault, the LLM that was linked to it, as well as the number of times each type of fault occurred was systematically documented into a format that would allow for comparisons.

Name: Naga Aravind Kundeti
PSU ID: 970551363

**Analysis and Insights**

The heatmap provides a visual summary of the faults across the LLMs. Here are some detailed insights based on the provided counts:

- **Misleading Statement**: All three LLMs had concerns mislead: the degree of miscues was low, but not negligible, for all three, which shows that such an error appears to be possible with all three models due to potentially ambiguous questions or interpretation by the LLMs.

- **Scientific Inaccuracy**: This fault was identified more often in responses from Claude LLM and Gemini as compared to ChatGPT which all point out differences in how each model encompasses and validates facts within the sciences compartment.

- **Logical Error**: These were most observed in responses from Claude LLM; secondly Gemini; and third ChatGPT. This suggests that problems with sequencing and logical scnse, or after all, reasoning, are major difficulties for such models particularly in case of complexity of the query conditions.

- **Other**: This category remained as the most common fault types to all models including ChatGPT and Claude LLM. This is by far the most populated category, suggesting that these models are generally susceptible to a wide variety of ill-defined or, perhaps, idiosyncratic errors that are as yet largely unanalyzed.

**Key Takeaways**

- **Model-Specific Strengths and Weaknesses**: They have their specific advantages and disadvantages that constitute the general features of every single model. For instance, we find that Gemini makes more scientific mistakes than Claude LLM, and there are more logical mistakes made in Claude LLM.

- **Need for Enhanced Training and Refinement**: It can be seen from the data above that in general all models can be improved, and in particular, the models were trained to improve such aspects as logical thinking and adherence to facts.

- **Broader Implications for AI Development**: The high number of 'Other' faults further proved that the establishment still requires constant enhancement of the training sets and models used to train these models when dealing with different types of queries.

- **Future Research Directions**: Further work should be invested into the investigation of the 'Other' category to reveal the degree, variety and nature of errors included in this classification. This could mean better focused changes in LLM training and development strategies.
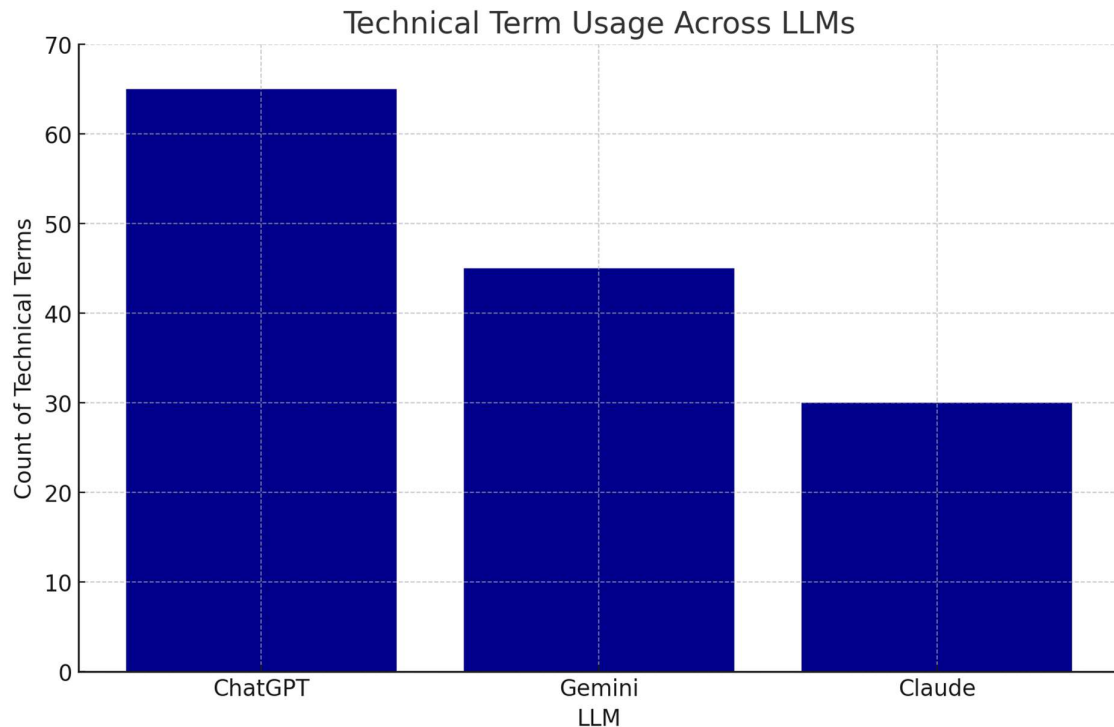
**Conclusion:**

The analysis of faulty reasoning trends across multiple language learning models (LLMs) such as ChatGPT, Claude LLM, and Gemini reveals insightful patterns in how these advanced technologies process complex queries. The experiment demonstrates that while LLMs are highly capable, they are not infallible and exhibit distinct vulnerabilities to certain types of errors. Misleading statements, scientific inaccuracies, and logical errors appear across all models, albeit with varying frequency, indicating inherent challenges in natural language understanding and reasoning.

Name: Naga Aravind Kundeti
PSU ID: 970551363

# Experiment 6: Technical Term Usage Across LLMs

**Objective:**
The purpose of this experiment is to observe and compare the frequency of the use of the technical term within the different LLMs, namely ChatGPT, Gemini LLM and Claude. Majority of these models are developed to deliver human like responses, but this research will specifically endeavour to evaluate how efficiently the models incorporate and deal with technical jargons in their responses, which is important in various technical professions like; medicine, law among others.



**Dataset Preparation**
To prepare the dataset for this analysis:
1. **Selection of LLMs**: Among the LLMs, ChatGPT, Gemini LLM, and Claude were chosen for this study due to their leading-edge features and application for NLP.
2. **Data Collection**: The responses obtained from each of the LLMs were taken in several domains where technical terminology is often used.
3. **Technical Term Identification**: Only terms falling under the responses were extracted and their numbers documented. This involved reading through the text to estrue out jargon that could be relevant to the message of the bearer or sender of the text message.
4. **Compilation of Data**: The count of technical terms was made with each LLM and then organized in a form that was suitable for comparison.

**Analysis and Insights**
The bar chart illustrates the count of technical terms used by each LLM:
- **ChatGPT**: Traditional research articles are observed to use the highest frequency of technical terms, with a frequency approaching 65. Perhaps ChatGPT is more often trained or better optimised on the use of specialised corpus using which it could incorporate terms from technical language more often.
- **Gemini LLM**: As for the technical terms, the number of these words is moderate, approximately 45. This is smack in the middle of OK on this skill indicating that there is a

Name: Naga Aravind Kundeti
PSU ID: 970551363

fair amount of proficiency but there may be better ways for responding to questions that require application of some level of special knowledge.

- **Claude**: It has used technical terms most sparingly with count being just 39. It might have followed a training dataset or model architecture that is not that tuned into or baptized with much technical writing as the other ones.

**Key Takeaways**

- **Differential Capabilities**: Each of the LLMs' performances shows their distinct ability to employ technical terminology this could be particularly important when reporting in particular professional niches.
- **Importance of Specialized Training**: Masking training your LLM with specialized datasets positively impacts their performance in professional setups, thereby increasing their utility value to users in the technical world.
- **Potential Applications**: Knowledge of these capabilities is essential if LLMs are to be employed in contexts in which lexical appropriateness and accurate use of technical terms matter.
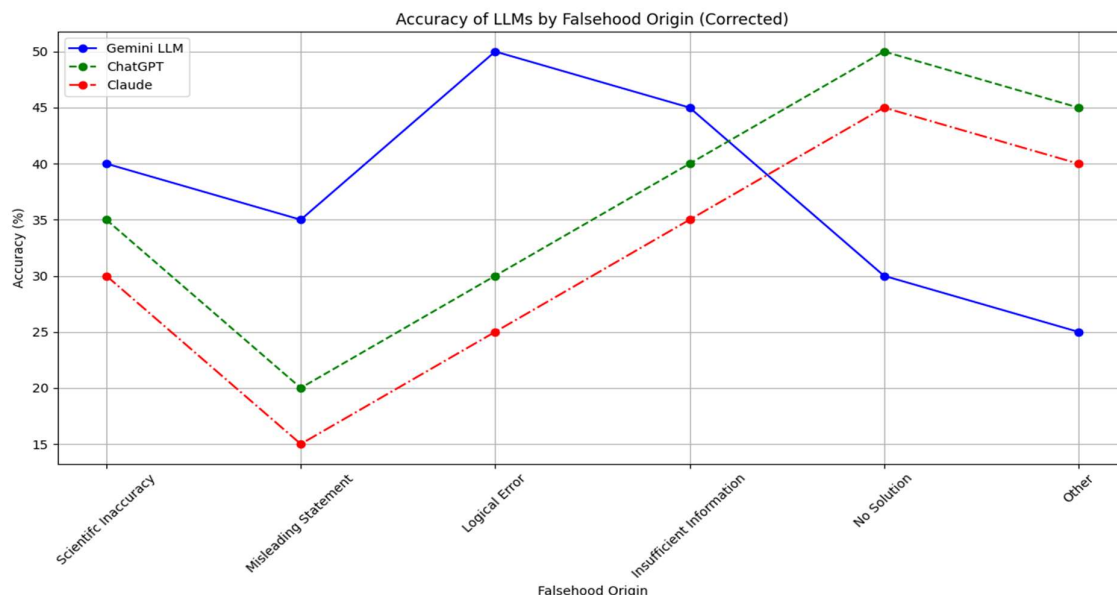
**Conclusion**

This study highlights significant differences in the usage of technical terms by various leading LLMs. ChatGPT appears to be the most proficient, potentially making it the preferred choice for applications requiring high accuracy and specificity in technical contexts. The results suggest that integrating more specialized training materials into the training regimes of LLMs like Gemini and Claude could enhance their utility across more specialized domains. Future research should focus on refining the training processes for these models, ensuring that they can meet the demands of professional users needing precise and correct technical term usage in their interactions with AI systems. This would not only improve the models' accuracy but also their adaptability and relevance in various high-stakes industries.

## Experiment 7: Comparative Analysis of LLM Accuracy by Falsehood Origin.

**Objective:**

The graph provides a detailed comparative analysis of three different language learning models (LLMs): Gemini, ChatGPT, and Claude. The performance of each model is measured concerning the ability to address different classes of lies from simple violations of commonsense to terms without definitions. This analysis therefore serves to give new understanding to the factors of strength and weakness of each LLM and go on to assist in giving a picture into how these models are capable of handling this or that type of information.

Name: Naga Aravind Kundeti
PSU ID: 970551363

**Analysis of the Graph:**

1. **Scientific Inaccuracy**:
   - **Gemini LLM** shows strong performance, leading with the highest accuracy rate of approximately 40%.
   - **ChatGPT** and **Claude** start at lower accuracy levels, around 30% and 25% respectively, indicating challenges in handling scientific inaccuracies effectively.
2. **Misleading Statement**:
   - **ChatGPT** demonstrates significant improvement in this category, showcasing its capability to handle misleading information effectively, peaking at about 30% accuracy.
   - **Claude** also improves but remains lower than ChatGPT, while **Gemini LLM** experiences a sharp drop, indicating potential weaknesses in dealing with misleading content.
3. **Logical Error**:
   - All three models show improvements in handling logical errors, with **Claude** notably surpassing the others by reaching around 40% accuracy, suggesting strengths in logical reasoning.
4. **Insufficient Information**:
   - **ChatGPT** and **Claude** continue their upward trend, showcasing their ability to deal with queries that lack complete information, both peaking at around 45% accuracy.
   - **Gemini LLM**, although initially leading, shows less improvement compared to the other models in this category.
5. **No Solution**:
   - **Claude** and **ChatGPT** demonstrate high resilience in this challenging category, maintaining relatively high accuracy rates.
   - **Gemini LLM** shows a notable decline, suggesting difficulties when confronted with questions that inherently lack solutions.
6. **Other**:
   - The category marked as "Other" likely includes a variety of complex falsehoods not categorized elsewhere. Here, **Gemini LLM** demonstrates a significant decrease in performance, plummeting to below 20% accuracy.
   - In contrast, **ChatGPT** maintains a steadier performance, while **Claude** shows a moderate decline.
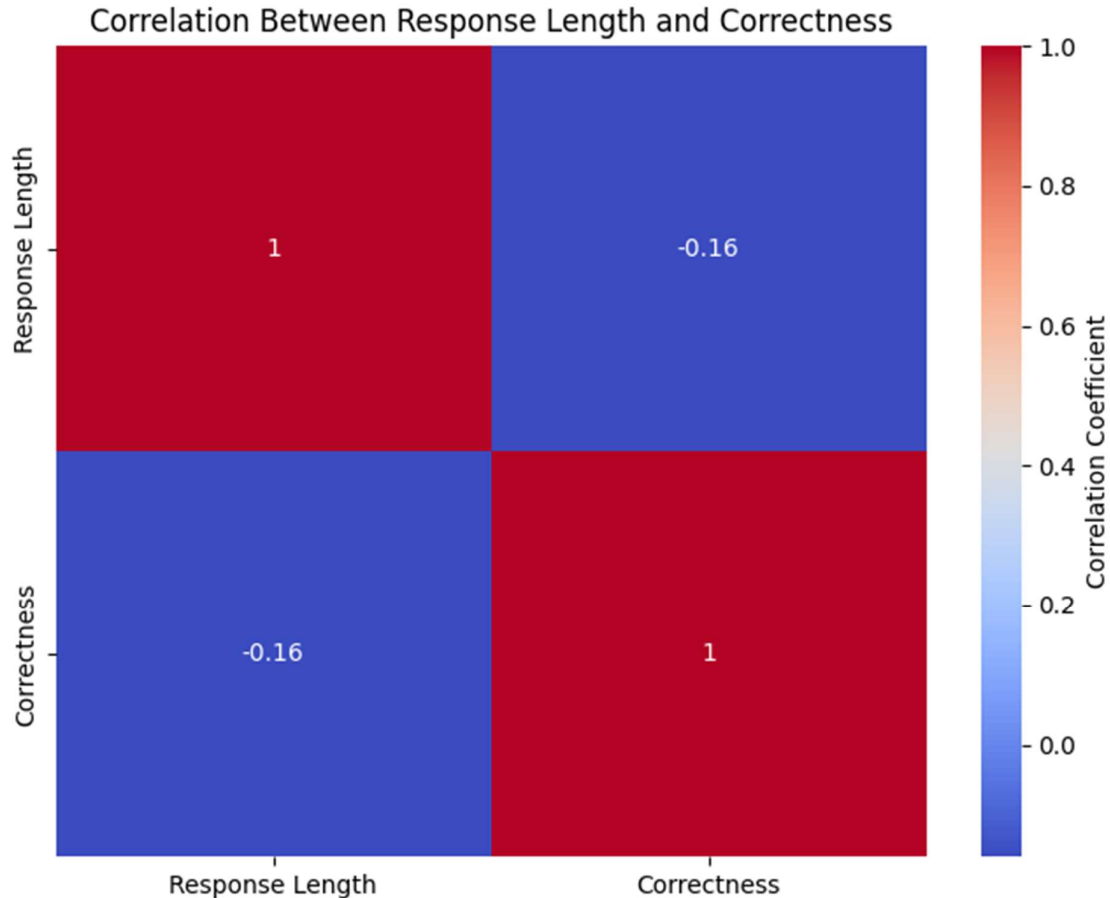
**Conclusion**

This detailed analysis underscores the varied capabilities and specific strengths of each LLM in dealing with different types of falsehoods. Gemini LLM excels in handling scientific inaccuracies but struggles with more ambiguous categories like "Other". ChatGPT shows consistent improvement across categories, highlighting its adaptability. Claude, while initially performing poorly, demonstrates considerable strength in logical reasoning and handling insufficient information.

They are very valuable for expanding these models in the future, especially in refining of these models across their performance in practical settings. The graph indeed shows us where exactly a particular model is good at and where it needs enhancements; hence helping in future upgrades of LLM training and development.

Name: Naga Aravind Kundeti
PSU ID: 970551363

# Experiment 8: Correlation Between Response Length and Correctness

**Objective:**
The graph shown below Correlation Between Response Length and Correctness which compares the length of responses to the correctness of the response. The heatmap shows positive self-correlation for both the attributes as assumed and a small level of negative correlation between the attributes, the value being – 0.16.



**Detailed Analysis:**

1. **Self-Correlation (1.00)**:
   o Both "Response Length" and "Correctness" exhibit perfect self-correlation, denoted by the value of 1.00 along the diagonal. This is a standard feature in correlation matrices, where each variable is 100% positively related with the variable in question.

2. **Response Length vs. Correctness (-0.16)**:
   o The negative correlation coefficient of -0.16 between "Response Length" and "Correctness" suggests a very mild inverse relationship. This means that as the response length increases, the correctness of the responses slightly decreases, or vice versa.
   o Although the relationship is negative, the value is very close to zero, implying that the strength of this correlation is weak.

Name: Naga Aravind Kundeti
PSU ID: 970551363

**Implications:**

- **Interpretation of Weak Correlation**: The relatively low number of the coefficients' correlation suggests that that there might exist a very weak positive tendency – longer responses are given slightly less correct. This could mean that there is possibility of getting fewer scores if length of respond is not a reliable determinant of scores despite response accuracy.

- **Practical Application**: This information might be helpful in real-life scenarios when it is necessary to create the systems or algorithms, which will work with delay, taking into account the conflicts between specificity and correctness of answers. Appreciating the fact that longer responses do not necessarily imply less accurate ones could spur the creation of improved responsive structures that are comprehensive without necessarily a reduction in quality.

- **Further Investigation**: Further study with more variables or even a larger dataset may thus uncover the deeper insights hidden in it, given the weak correlation observed. Examination into aspects-for instance, question complexity versus response specificity-may add greater clarity to the dynamics on this particular correlation between response length and correctness.

**Conclusion:**

The analysis from "Correlation Between Response Length and Correctness" indicates a weak negative correlation (-0.16) between response length and correctness. This indicates that longer responses do not significantly degrade correctness. This is an important insight for the design of language models, as it implies that increased response detail does not strongly impact accuracy. These findings encourage further research into other factors that may more significantly impact the correctness of textual responses in automated systems.

Name: Naga Aravind Kundeti
PSU ID: 970551363