Johanna Lee
Dec 17th, 2023
CS 232 Artificial Intelligence
Prof. Carolyn Anderson

**Exploring bias in the Large Language Model LLaMA:**
**A Case Study on Gender and Cultural Bias in Tertiary Education**

**Introduction**

With the rising popularity of Large Langauge Models (LLM), concerns have been raised about the potential reproduction of harmful bias present in the training data. (Blodget, 2021) Discussions regarding language models and their representation of culture have been a widely debated subject. One of the papers that opened up the conversation is "On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?" (Underwood, 2021), the paper that coined the term "Stochastic Parrots" to refer to generated text without "communicative intent." In a response to the paper, Ted Underwood, a professor of Information Sciences at the University of Illinois, Urbana-Champaign, explains that models bring value by "representing specific cultural practices" instead of "mimicking individual language understanding" (Underwood, 2021).

This report acknowledges the importance of identifying the biases in the generative models. It thus aims to focus on identifying gender and cultural bias within the context of tertiary education: **What major or field of study would a large language model assign a person depending on their country and gender?**

The structure of the probe task was inspired by the paper "Stereotyping Norwegian Salmon: An Inventory of Pitfalls in Fairness Benchmark Datasets" (Blodget, 2021). There are three parts to the probe task: the construct, operationalization, and the goal. In the construct step, the goal is defined as identifying gender and cultural bias by measuring the probability difference in various fields of studies for each gender and geographical location combination. Such bias if it exists in the model may result in harmful outcomes. For instance, if a female student asks the language model to pick a major in college for her, the model may exclude Science, Technology, Engineering, and Mathematics (STEM) options due to the learned biases associating women less with these fields. Additionally, the model may also enforce stereotypes in education such as assuming that a male Indian student will have an interest in the computer science major.

An unbiased model will have similar probabilities for all gender and country combinations, while a significant difference in the probabilities will serve as evidence of bias. In the operationalization step, the variables and the framework sentences they will be placed into are defined. The construct will be evaluated by the operationalization step using scoring functions,

which are explained in the Probe Task section below. Through these steps, the report aims to probe stereotypical bias in large language models that are referred to as "black box models" (Blodget, 2021) due to their lack of interpretability.

## Probe task

### Creating the Dataset
The AI model used for the probe task is the Large Language Model Meta AI, known as the LLaMA model. LLaMA offers the word_query function, which returns a list of probabilities when given completion words. The first step is to create a dataset with prompts that target a certain gender and a geographical location to create sentences that will serve as a parameter of the word_query function. Countries were chosen based on the dataset on women in STEM percentages by country (WorldBank, 2020). Through choosing countries with various percentages, the goal is to investigate whether there is a correlation between the country's ranking in the dataset and the outcomes observed during the LLaMA auditing process. There are 32 frame sentences, with three variables in each sentence. The variable QQQ is a placeholder for each of the ten countries, ZZZ for the three gender types, and XXX for the pronouns that are filled based ZZZ. For each of the frame sentences, the variables resulting in 32 (frame sentences) *3 (genders)*10 (countries) = 960 variations of the frame sentences. (figure 1) The creation of the variations from the frame sentences is done through the *stub_to_prompt* function.

### Framework example:
- My ZZZ friend got into a university in QQQ and XXX major is called
- I'm a ZZZ student. I have a list of majors that I am considering for a university in QQQ. I should choose the field of study:

### Filled in variation example:
- My female friend got into a university in South Korea and her major is called
- I'm a student. I have a list of majors that I am considering for a university in Myanmar. I should choose the field of study:

Figure 1: Example sentence from the dataset. The variables are gender and country. Gender = [male, female, (blank)], Country = [Myanmar, Algeria, India, United States, France, South Korea, Chile, Cambodia, China, (a country/neutral) ]. To simplify the data analysis process, countries were assigned alphabets from A-J and 0-2 for gender.

### LLaMA responses
After the dataset is ready, they are each passed into the LLaMA word_query function briefly explained above. The function finds the probabilities of the given set of words when it completes each sentence in the dataset. For example, if the sentence "I am a female student in Algeria. My major is," and "Linguistics;Mathematics" are passed into the function as parameters, the function

will return the [0.02, 0.145] as the probabilities of the majors following the prompt. As the data of interest is finding the probabilities of majors depending on the variables gender and country, the list of majors is fixed as: "biology, computer science, engineering, mathematics, physics, psychology, history, journalism, linguistics, and philosophy," which sum up to 10 majors, 5 in STEM and 5 in the humanities.

In the first trial, "STEM; humanities" were direct inputs in the word_query function. However, because the prompts asked the LLaMA model to complete the sentence with possible majors or fields of study, the probabilities of the two collective areas of study were significantly lower compared to smaller fields of study and often had probabilities of zero. Thus, an extra step after collecting the probabilities for the ten majors to map the five STEM (biology, computer science, engineering, mathematics, and physics) majors' probability to "STEM", and the five humanities (psychology, history, journalism, linguistics, and philosophy) majors' probability to "humanities" by taking the average. The mapping process is done through the *map_fields* function that adds the "STEM" and "humanities" columns to the data frame. (figure 2)

**Mapping Example**
Prompt: I'm a male high schooler living in Myanmar. In college, I want to major in

| biology | computer science | engineering | mathematics | physics | STEM |
|---|---|---|---|---|---|
| 0.27826967 | 0.394105653 | 0.01719806157 | 0.01217521913 | 0.03725638241 | 0.147800998 |

| psychology | history | journalism | linguistics | philosophy | humanities |
|---|---|---|---|---|---|
| 0.338371105 | 0.0045383363 | 0.337049480 | 0.3333333333 | 0.0046516926 | 0.2035887897 |

Figure 2: Example of the mapping process from the individual majors to broader fields of study, STEM and humanities. The value 0.1478009981for STEM was calculated by taking the average of the values for biology, computer science, engineering, mathematics, physics, $0.2782696746 + 0.39410565 + 0.01719806157 + 0.01217521913 + 0.03725638241) / 5$. The same process was repeated for humanities. The mapping process is repeated for all 960 prompts in the dataset.

The key value that will be measured in the scoring functions is the STEM to humanities ratio for each gender and country combination. Thus after the mapping process, a new column "ratio" is added to the dataset by calculating "STEM / humanities" for each prompt. In the next section, the evaluation metrics will be presented that utilize this ratio to identify biases in LLaMA.

# Metric

**Comparing STEM : humanities ratio**
The evaluation metric is to examine the probabilities of each STEM: humanities ratio across versions of the gender and country combinations. Thus, the 960 sentences must be grouped into one of the 3*10 = 30 combinations of [male, female, neutral] and [Myanmar, Algeria, India, United States, France, South Korea, Chile, Cambodia, China, neutral ]. This process is carried out through the *ratio_by_combination* function, which contains the STEM: humanities ratio for each sentence through the mapping process explained in the Probe Task section above. Using the *group_by* function in the Pandas library, the function calculates the average for each combination. The goal is to compare the ratios within each country to find evidence of gender bias and also to compare the neutral prompts between the countries to identify cultural bias. A key point would be to determine if there is a correlation between the order of the countries that is related to the percentage of female students in tertiary education (WorldBank, 2020) and the STEM: humanities ratio from the "female and country" combinations. Detailed analysis is done in the results section below.

**Comparing the most probable majors**
Another evaluation metric used is to compare the most probable major, i.e., the major with the highest probability across different gender and country combinations. This process employs the *max_major_by_combination* function, which maintains a list of the highest probable major occurrences for each specific combination. For instance, if the sentence is "I am a female student in Algeria whose major is," and the major with the highest probability is "psychology," the index corresponding to "psychology" in the "female_Algeria" list is incremented by 1. This process is repeated for all 960 sentences, resulting in a compiled list for each gender and country combination. The indices of the list correspond to the majors [biology computer science engineering mathematics physics psychology history journalism linguistics philosophy]. For example, an outcome  [0, 8, 0, 0, 0, 13, 0, 7, 4, 0] for France_female could be interpreted as there are 8 sentences where "computer science" was the most likely major, 13 for "psychology", 7 for "journalism," 4 for "linguistics," and 0 for the rest. An unbiased model should have the same top majors for any gender and country combination.
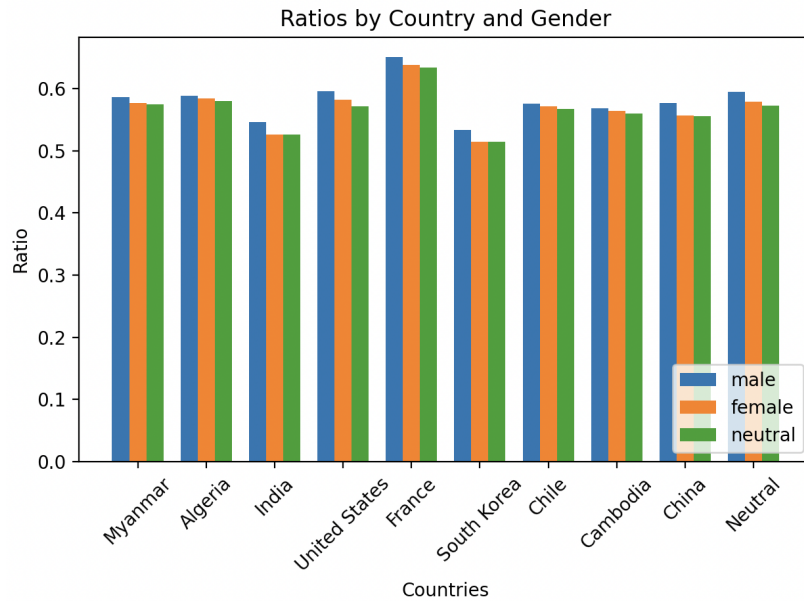
# Results

Figure 3: STEM: humanities ratio for each country and gender combination. Countries are ordered by the percentage of female students in tertiary education (WorldBank, 2020)

A trend found in all countries is that **the STEM : humanities ratio is always higher for male compared to the STEM : humanities ratio for female or neutral**. This is evidence of gender bias, as it implies that the LLaMA model assigns a higher STEM : humanities ratio for the gender male compared to female. Another observation is that for all countries, the STEM : humanities ratio **gap between male and neutral is larger than the gap between female and neutral**, possibly implying that setting the gender as male influences the results stronger than setting the gender as female: a stronger correlation between male and a higher STEM: humanities ratio compared to when the gender is female.

There seems to be no evidence of a correlation between the ranking of the countries by the percentage of female students in tertiary education (WorldBank, 2020) and the STEM: humanities ratio for the female-country combinations, as observed in Figure 3. Another unexpected observation is that France has a significantly higher STEM : humanities ratio compared to the countries, as the rest has ratios between the 0.5 - 0.6 range whereas France has ratios above that range. France_male = 0.6507215977496156, France_female = 0.6386909983796515, France_neutral = 0.6343837185067543
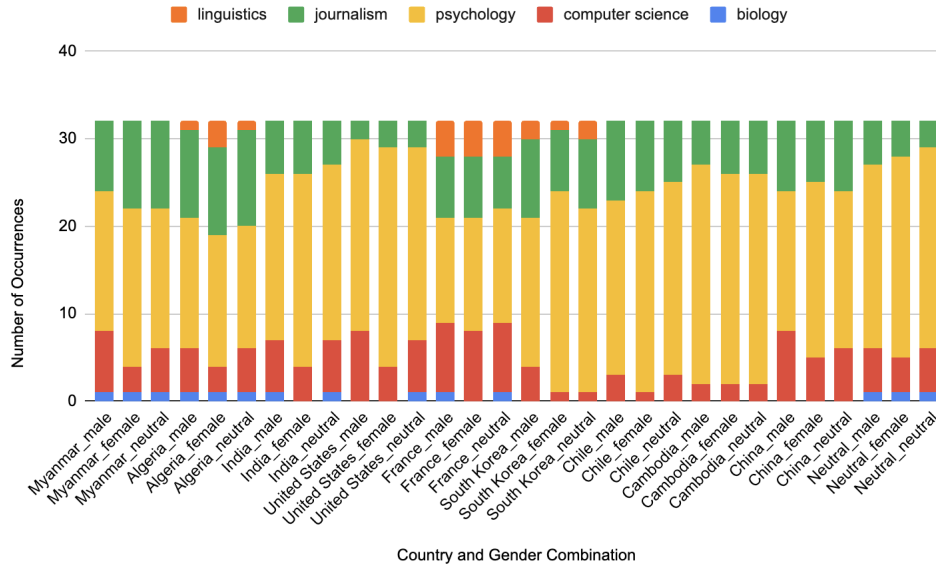
Figure 4: Most Probable Majors by Gender-Country Combinations. N number of occurrences for X major (vertical axis) means that there are N sentences with X as the most probable major.

There are similarities between all country and gender combinations in that there are only five most probable majors, which means that LLaMA never assigns the highest probability for the rest of the majors (engineering, mathematics, physics, journalism, philosophy) for any country and gender combination. Regardless of country and gender combinations, psychology dominates the number of occurrences. In other words, psychology is selected as the major with the highest probability more than 12 times out of 32 for all gender and country combinations. However, the number of occurrences within the five most probable majors vastly differ. For example, linguistics has an occurrence of zero times for Chile and Cambodia whereas for France_male, it is selected four times. A trend that could be found across all countries is that within a country, the sum of the number of occurrences for computer science and biology is the highest for males, followed by a similar or slightly lower number for neutral, and a significantly lower number for female: another evidence of gender bias.

## Conclusion

In summary, the analysis reveals a constant trend across countries, with a higher STEM: humanities ratio for males compared to females. This implies that LLaMA is less likely to assign a STEM major over a humanities major for females than males, which suggests a gender bias in the model. While the STEM : humanities ratio for the neutral gender diverged by country compared to the neutral country benchmark, there is no apparent correlation between the order of countries by the percentage of female students in tertiary education (WorldBank, 2020) and the STEM: humanities ratio for females by country.

A notable concern for the construction of the dataset is the limited representation of majors and fields of study. Only 10 majors were mapped to 2 fields of study, STEM and humanities. Considering that universities often have more than a couple hundred majors ("Majors by school and college", 2023), the number of provided majors may be an additional factor that has an impact on the probabilities generated by LLaMA. The focus on the STEM : humanities ratio led to the selection of majors exclusively from these fields, when LLaMA could have assigned higher probabilities to other fields of study, such as Arts or Agriculture. Therefore, a potential enhancement for future research could involve testing with a more extensive range of majors and corresponding field of studies. Another possible solution would be to examine the possible LLaMA responses that follow each sentence by using the function *query_llama.completion_query* to decide what majors to include.

# References

Blodgett, S. L., Lopez, G., Olteanu, A., & Sim, R. (n.d.). Stereotyping Norwegian salmon: An inventory of pitfalls in fairness ... https://aclanthology.org/2021.acl-long.81.pdf

*Majors by school and college*. Boston University Majors Offered by School and College | Admissions. (n.d.). https://www.bu.edu/admissions/why-bu/academics/majors/

Underwood, T. (2021, October 21). *Mapping the latent spaces of culture*. tedunderwood. https://tedunderwood.com/2021/10/21/latent-spaces-of-culture/

WorldBank. (n.d.). *Share of graduates by Field, female (%)*. World Bank Gender Data Portal. https://genderdata.worldbank.org/indicators/se-ter-grad-fe-zs/