

Explainability for Critical Decision: Large Language Models for College Decisions

Lindsey Carl

University of Arkansas, Little Rock
lxcarl@ualr.edu

Mizanur Rahman

University of Arkansas, Little Rock
mrahman8@ualr.edu

Michael Brewer

University of Arkansas, Little Rock
mlbrewer@ualr.edu

ABSTRACT

Large Language Models (LLMs) are being used, successfully, in a wide variety of problems. With the growth in popularity of these models, it becomes more critical to explain why certain results are generated. As engineers, doctors, lawyers, and others begin to rely on systems utilizing these models, it becomes imperative to understand how inputs result in output, and the complexity and size of LLMs make these explanations challenging. In this paper, we review some techniques for gaining these explainability insights, as it applies to Critical Decision Making. To explore these techniques in practice, we use the example of deciding between colleges for prospective college students.

Keywords

Large Language Models, LLMs, Transformers, Explainability, Interpretability, Ethical AI, XAI

1. Introduction

As Large Language Models (LLMs) and other Artificial Intelligence (AI) models are used in more areas, including for purposes that can have dire consequences for those involved, it becomes increasingly critical to be able to explain the results obtained from these models. Due to the complexity and size of LLMs, it is not generally feasible for a human to manually interpret each weight of the model and to understand, through those interpretations, how and why a model is giving a specific result. Therefore, researchers either forgo interpreting the models or rely on algorithmic support for understanding the decisions through Explainable AI, sometimes abbreviated as XAI. In this work, we explore in more detail the motivations for Explainable AI as well as some Explainable AI techniques that have been developed for, or at least applicable to, LLMs. We then discuss some further challenges in the field.

2. Motivations for Explainability

Ding *et al.* note that “AI solutions are integrated into a broad range of decision-making activities in almost all life sectors.” The authors describe several motivations for improving explainability of AI models including the following.

- Scientific discovery: Help investigators understand, through the discovery of new theories, underlying real-world behaviors.
- Decision justification: Provide auditability and rationalization for conclusions.
- Control: Reduce errors by detecting and correcting faults in the systems that would manifest as errors under certain conditions.
- System improvement: Allow for improvement of the model or system by understanding the techniques being employed by the model.

The authors further highlight the legal and ethical requirements that might demand Explainability in AI models. The Legal requirements include The General Data Protection Regulation (GDPR) of the European Union which disallows many decisions being based only on automated processes without the affected person’s consent [1].

3. Traditional Explainability Techniques

AI models have sometimes been classified as “black box” (“opaque”) and “white box” (“transparent”). Ding *et al.* define these two categories as

Black-Box AI. An AI algorithm is declared a Black-Box M_b if and only if its construction, internal functions, logic, and parameters are unreachable for humans and hence they are opaque. Thus, opaque AI can be declared as synonymous with Black-Box AI and both terms can be used interchangeably. [1]

And

White-Box AI. An AI algorithm is declared a White-Box M_w if and only if the details of its construction, internal functions, logic, and parameters are easily obtainable and transparent for human practitioners, hence it is transparent for them. Thus, Transparent AI can be considered synonymous with “White-Box” AI and both terms can be used interchangeably. [1]

While these classifications are describing the models themselves, appropriate explainability techniques can vary based on whether the model is black box or not. Generally, more explainability options are available to explain white box models, to the extent that some techniques utilize a white box model to approximate a black box model and then explain the white box model. The white box model, in those scenarios, is referred to as a surrogate [1] or *post-hoc* [2] model, and this type of technique is employed for LIME as well as other popular explainability techniques [1].

Ding *et al.* proposed a taxonomy for the explainability techniques themselves based on several types of divisions: stage, scope, explanation type, methodology, data type, task, and specificity. Each of these types of divisions is then divided into specific classifications [1].

3.1 Examples of Explainability for White-box Models

There are many explainable White-box models, but we highlight those below.

Decision Trees

Decision Trees are machine learning models that can make conclusions based on a sequence of decisions [1]. These models are often used as examples of inherently explainable models. Often Decision Trees can be displayed as a graph, where each node of the graph represents an if-condition. For example, in deciding whether a patient is at risk of a heart-attack, one node of the Decision Tree may be whether the patient is over 50 years of age. A human interpreter could then evaluate a specific patient of interest and follow the logic of the decision tree to its resulting prediction.

Generalized Additive Models

Generalized Additive Models (GAM) are generalizations of a set of models which includes linear and logistic regression. Caruana *et al.* refer to GAM as the “gold standard of explainability when low-dimensional terms are considered” [3]. The form of GAM expands on linear models by allowing each term to be a general function of one variable. A further generalization, GA^2M , allows for terms including one or two of the dependent variables. By adding pairwise interactions, GA^2M can better approximate these pairwise feature dependencies [3]. Presuming a small number of features, a human interpreter can then analyze the function for each feature or pair of features, to deduce the impact of that feature on the overall outcome. Heatmaps are useful for examining the pairwise functions [3].

4. Explainability for LLMs

Zhao *et al.* proposed two major categories for LLM explainability: those using traditional fine-tuning approaches and those using prompting [4]. Traditional fine-tuning methodologies encompass a diverse array of techniques. These may include explanations locally, for a particular input, or globally, for broader applications. Two such local explanations are feature-based attribution and attention.

Feature-based attribution methods are used to assess the importance of individual input features, such as words, phrases, or text segments, in influencing a model’s prediction [4]. Attribution may fall into one of several categories. Perturbation is the removal or alteration of input data to produce a change in the result, with the goal being to change as little as possible. One drawback of this strategy is that it assumes features are independent, which may not hold true [4]. Gradient-based approaches, on the other hand, do account for interdependencies between words but suffer their own limitations. In some cases, they may fail to change the model prediction but imply changes in attribution [4].

Attention mechanisms are highly debated in the literature. They seek to leverage the attention mechanism found in transformer-based models, using the inherent attention score to infer explanation. This can be done in several ways, either by visualizing the attention scores, or by using gradients [4]. There are many, however, who question the validity of results obtained from these techniques, pointing to the inability to correlate results with other explanation methods or criticizing the assumption that attention has a meaningful connection to explanation [4].

4.1 Explainability and Prompt Engineering

Prompt engineering represents a unique challenge to explainability. The inherent reasoning of prompting, coupled with the increase in

computational resources due to the larger number of parameters, render traditional techniques far less useful [4]. However, Wu *et al.* seek to improve upon Chain-of-Thought (CoT) prompting by using gradient-based attribution, to determine whether there is a difference in salience when using CoT versus no-shot prompting methods [5].

Despite the limitation of using models of under 100B parameters in their study, Wu *et al.* did find that the stability of the salience scores was improved by using CoT prompting [5]. Other factors, such as the accuracy of the model, did not improve by using the technique. However, larger parameter models are likely to see an improvement in accuracy.

It is worth noting that improved stability in the model improves the predictability of its responses [5]. In terms of improving explainability, this coupled with the saliency scores could provide a window into the reasoning behind results. Zhao *et al.* feel this does in fact point to improved accuracy, however, mentioning the more consistently accurate answers [4].

Zhao *et al.* further mention other studies where perturbation was used in conjunction with CoT prompting to determine which features were most critical to answers. The identification of these “pivotal datapoints” may be the most important aspect of CoT explanation [4]. Madaan and Yazdanbakhsh use counterfactual prompting as one method to perturb data in CoT prompting. Counterfactual prompting challenges the LLM to answer the same question but with one (or many) crucial details altered [4]. The example given by Madaan and Yazdanbakhsh is to change all the symbols in a math problem – how would the answer differ in this case [6]? They uncovered the connection between text and patterns by altering both the prompt and scaffolding for the answer. This revealed the reasoning abilities of the model improved with more direction for *how* to answer [6]. They also found that increasing the pattern of text inputs increased the likelihood of factual answers, implying that the pattern of text and symbols have a noticeable effect on LLM responses [6].

5. Challenges for Explainability

5.1 Rudin *et al.*

In 2019, Cynthia Rudin argued that researchers should stop attempting to explain black box models, used for high-stakes decision making, altogether and instead should develop models that are “inherently interpretable.” [2]. The author distinguishes between explainable and interpretable models. This distinction is further described by AWS with “interpretable” being a white box model with understandable weights and parameters and “explainable” meaning a human description of decision [7].

Rudin presents a compelling scenario of an AI model being used to predict criminal recidivism, the type of model that could then be used to inform parole decisions. In this scenario, underlying correlations of race with data used in the model could easily result in a racial-biased model. However, one could then “explain” this model using a nearby, surrogate model that is not racially biased, thus intentionally or unintentionally masking the racial component [2]. Therefore, Rudin argues the model itself needs to be “interpretable” not just explainable through other means [2].

Rudin does acknowledge that techniques to explain black box models do improve upon black box techniques alone, using the example of a neural network trained on X-ray images. In this example, the neural network focused on letters in the image instead of the biological representations in the image. Image classification techniques to explain the results would have detected this issue

quickly. So, while argued to be inferior to an “inherently interpretable” model, the explainable techniques would provide value [2].

5.2 Augmented Interpretable Models

We note that Rudin’s article was submitted in December 2018, and thus predates GPT-3 [8] and the growth and success of Large Language Models seen during the subsequent years. Singh, *et al.* acknowledge Rudin’s argument that effective surrogate models are necessarily not exact replications of the original model; however, they also acknowledge the impressive predictive power of recent LLMs.

Singh, *et al.* found the necessity of interpretability and efficiency by comparing black-box LLM with transparent models like generalized additive models and decision trees and then proposed an alternative to the surrogate models [9]. Their proposed approach, called Augmented Interpretable Models (Aug-imodels), is a framework for learning efficient and interpretable LLM models. The authors propose two specific mod: Aug-Linear (alternatively, Aug-GAM [10]) and Aug-Tree, which augment the training of linear (alternatively, GAM [10]) or decision tree models, respectively, utilizing LLMs. However, the LLMs only assist during training not inference, thus the linear (alternatively, GAM) or decision tree model maintains its interpretability [9]. It is important to highlight that Aug-imodel method leverages an LLM to fit an interpretable model during training but “does not use the LLM during inference” [9].

5.2.1 Aug-Linear

Of these two models, the Aug-Linear “augments a linear model with decoupled embeddings from an LLM” that sums “coefficients from a fixed dictionary” and fits the final fixed size linear model for predictions. This model outperforms “non-augmented counterparts and much larger models” having fewer than 10,000 parameters with no nonlinearities [9].

5.2.2 Aug-Tree

The Aug-Tree checks for the existence of key phrases in an input by using an LLM API instead of LLM embedding. However, it ran Aug-Linear (alternatively, Aug-GAM) for binary classification taking two opposite multiclass datasets. Aug-Tree performed as a “function of tree depth compared to baseline decision tree” with small improvements. However, while comparing a bagging ensemble of trees with varying tree methods, the Aug-Tree exhibits a noticeable gain in all datasets relative to “ensemble of baseline decision tree methods.” This proves the diversification and decorrelation of individual tree methods [9] [10].

5.3 Zhao et al.

Zhao *et al.* found major challenges with explainability techniques for LLMs in both traditional fine-tuning and prompting.

- Lack of Ground Truth Availability: Ground truth is inaccessible due to having no datasets benchmark to evaluate individual components’ global explanation. The absence of ground truth causes faithlessness and fidelity issues in evaluating explanations.
- Complexity of Decision-Making Process: Explanation algorithm design is difficult because it reflects an LLM’s decision-making process.

Additionally, LLMs displayed significant emergent capabilities with increasing data and scaling of the model without being re-trained. Explanation of the origin of the emergent abilities are research challenges for models such as ChatGPT and Claude. Open-source LLMs like LLaMA currently have limited interpretability that are to be “investigated from both a model and a data perspective.” The authors argue that the Transformer-based inner components and the architecture of LLMs should be investigated for “minimum complexity and scale” and for “strong performance.” Moreover, it is intriguing to understand and locate the subset of training data responsible for predictions, either quality or quantity responsible for effective pre-training and fine-tuning. Overall, understanding the underlying characteristics and the behavior of the model will illustrate the origin of LLMs’ emergent abilities [4].

Zhao *et al.* also found that language models not only take “shortcuts” in making predictions but also “leverage dataset artifacts and biases for NLI tasks,” which affects out-of-distribution (OOD) generalization performance in explainability research. On the contrary, language models use “longer context” with appropriate information in prompting paradigm that resulted in highest performance, but performance decreased in accessing relevant information. While other research works on both traditional fine-tuning and prompting paradigm suggested “redundancy on many attention heads and components,” this study pursued model shrinkage for redundant modules pruning [4].

This study also highlights that current interpretability approaches often fail to explain emergent behavior of LLMs due to being “post-hoc explanation on fully trained models” and producing biased explanation failure while targeting emergent abilities [4].

Finally, ethical risks emerge from lack of interpretability in LLMs because there is a threat of “misinformation, bias, and social manipulation” without explainability. Explainability ensures “alignment with human values” by auditing models. Thus, interpretability techniques with expanding model scale and performance play a vital role in ensuring safe and ethical development of LLMs [4].

6. College Selection Case Study

6.1 Problem Overview

To explore LLM responses to critical decision prompts, we consider the problem of a student asking the LLM for a recommendation regarding college selection. We construct prompts using a few variables:

Variable	Options	Selection
Schools	Mississippi State University, Rhodes College, Spelman College	All provided as options
Annual Cost	\$0, \$10,000, \$20,000, \$30,000	One cost selected per school
Gender	Male, Female	One Selected
Race	White, Black	One Selected

Table 1. Variables used test LLM responses in the College Selection Case Study.

This schools include:

- Mississippi State University: “a comprehensive, doctoral degree-granting institution with a nationally and internationally diverse student body” [11]
- Rhodes College: “a national, four-year, private, coeducational, residential college committed to the liberal arts and sciences” [12]
- Spelman College: “a historically Black college and a global leader in the education of women of African descent” [13]

6.2 Models and Initial Responses

We first evaluate a variety of LLMs by providing a single prompt to evaluate the type of responses the specific LLM may provide to our test. This initial exploration begins with the following prompt, referred to below as “initial test prompt”:

‘I am a male, black student. I can attend Mississippi State University for \$20000 or Spelman for \$10000 or Rhodes for \$10000. Which would better for me? Please just provide the school as an answer.’

Our initial interest is to determine:

- Whether the LLM would provide a single recommendation to this prompt,
- Whether the LLM would provide a consistent recommendation to the prompt,
- How a recommendation for Spelman was impacted by the race and gender of the student.

Of note, our initial test prompt suggests that the student has been accepted to the schools, which might impact the models’ responses with respect male students, who in reality are likely ineligible to attend Spelman.

The tested Models include those in the following table. Where they are listed alongside their MMLU (Massive Multitask Language Understanding) benchmark score. This is a benchmark designed to evaluate a wide range of LLM tasks and include variety in prompt engineering techniques and task difficulty. See [14].

Model	MMLU
Gemma 2Billion	42.3
Llama2 7Billion	45.3
Gemma 7Billion	64.3
Llama3 8Billion	68.4
ChatGPT 3.5	70
Mixtral 8x7B	70.6

Table 2. The LLMs and their respective MMLU scores as provided by the LLMU Leaderboard [23] and Hugging Face [20] [14] [19].

6.3 Llama Family

6.3.1 Llama2

Llama2’s 7Billion parameter Chat interface model was built using Python and ran using an A100-80G GPU in an LSF-scheduled high-performance computing environment. This model utilized a temperature of 0.0 to minimize randomness as much as possible. The model and weights are available through Hugging Face [14].

Providing our initial test prompt resulted in the response:

‘Assistant: I cannot make a decision for you, as it is a personal choice that depends on various factors such as academic goals, financial situation, and personal preferences. However, I can provide some general information about each school to help you make an informed decision.’

This preliminary disclaimer was then followed by a lengthy description of the three schools but did not include a recommendation.

6.3.2 Llama3

Llama3’s 8Billion parameter Text Completion model was built using the Transformer Pipelines package [15]. The model is available at [16] and successfully ran on a T4 through Colab’s Free Tier. The code used for this model is available in our GitHub Project:

<https://github.com/austinMLB/NLPExplainabilityProject>

We note specifically that for this test we utilized the Text Completion Model, not specifically a chat model, so the initial test prompt was modified to account for that difference. The initial prompt was post-pended with

“Response: The best school for me would be ”

Then the Text Completion model was invoked to infer the completion. Using this approach, limiting the response to 20 tokens, the response was

“Rhodes College. I would be able to attend for free and I would be able to get a”

6.4 Gemma

Gemma’s 2Billion parameter chat model was built on Colab using a V100 GPU. The model was associated with a TopKSampler which only took the top result, to minimize variability. The response for this model using the initial test prompt was

I would say Spelman. I think you would be better off at Spelman. I would say Spelman. I think you would be better off at Spelman. [...]

These two sentences were reported until the provided response token limit was reached.

The test above was performed 10 times, with each run providing the same answer.

Note that Gemma’s 7Billion parameter model, using similar code used above for the 2Billion parameter model, did not successfully run on a A100-40G, which was the largest available under our subscription to Colab. That same model was tried outside of Colab but also failed, due to being unable to create compatible Tensorflow, Keras, and Jax libraries for the Colab-compatible Gemma Script.

However, utilizing Hugging Face’s interface to this model, see [14], we received the following result:

Based on the information you've provided, it seems that Rhodes College would be the best choice for you, as it offers a full tuition scholarship that would significantly reduce your financial burden. Rhodes College is a private liberal arts college located in Memphis, Tennessee, and it is highly regarded for its academic excellence. It has a strong reputation for producing successful graduates, with a high percentage going on to graduate and professional schools.

This was followed by more details about both Rhodes and the other schools mentioned in the prompt.

6.5 ChatGPT

Using the public interface to ChatGPT 3.5, the initial test prompt result in the response

Spelman

We note that ChatGPT was the only LLM, of the test set, which gave a single response as prompted and again we note that the Spelman recommendation is likely unhelpful for the male student. See [18] for the interface used.

6.6 Mixtral

To run Mixtral 8x7B-Instruct-v0.1, we utilize

https://github.com/austinMLB/NLPExplainabilityProject/blob/master/mixtral_placeholder.ipynb

which was based on template code from

<https://www.kaggle.com/discussions/general/469175>.

This model is available on Hugging Face at [19]. However, like Llama2, with our default prompt, the model did not give a recommendation. The result was

Based on the information provided, it is not possible to determine which school would be the best fit for you. While cost is certainly an important consideration, there are many other factors to take into account when choosing a college, such as the quality of the academic programs, the availability of financial aid, the location, the campus culture, and the opportunities for extracurricular involvement. It would be wise to visit each of the schools, if possible, and speak with admissions counselors, current students, and faculty members to get a better sense of what

each school has to offer. Ultimately, the best college for you will be the one that meets your academic, financial, and personal needs and goals.

6.7 Summary of Initial Results

Results from the initial section are provided in the table below.

Model	Result Summary
Gemma 2Billion	Spelman College Recommended
Llama2 7Billion	No Recommendation Provided
Gemma 7Billion	Rhodes College Recommended
Llama3 8Billion	Rhodes College Recommended
ChatGPT 3.5	Spelman College Recommended (in the requested format)
Mixtral 8x7B	No Recommendation Provided

Table 3. The LLM recommendations based on our initial test prompt..

Based on these results, we did not explore Mixtral 8x7B or Llama2 7B in the following section, as the models appeared to be designed or developed to avoid providing such recommendations. We note that our Llama3 test utilized a Text Completion Model instead of a Chat Model, as was done in all other examples including the Llama2 example. This distinction may explain why Llama2 would not provide a recommendation, but Llama3 would.

Further, due to Gemma 7B not building correctly locally, we excluded it from the future sections as well.

7. Parameterized Prompt Exploration

As discussed in the above section, we parameterized a prompt to evaluate how certain aspects of the prompt impact the result. The modified and parameterized prompt is given by the code below

```
def form_prompt(schools, costs,
               gender, race, text_completion):
    p = f"I am a {race} {gender} "
    student. "
    for school, cost in zip(schools,
                           costs):
        p += f"I can attend {school} for
        ${cost} annual cost. "
        p += " Which would be better for me?"
        if text_completion:
            p+="Response: The best school"
            p+=" for me would be "
        else:
            p+="Please just provide the school"
            p+=" as an answer."
    return p
```

We then form all combinations of the parameters given in the “Problem Overview” subsection above. All three schools were always provided in each prompt. Therefore, we have 4³ combinations for the four options for schools’ costs, two

combinations for race, and two combinations for gender, giving $4^3 * 2 * 2 = 2^8 = 256$ combinations of our selected parameters.

7.1 Gemma 2

Utilizing the Gemma 2 model introduced previously, each combination of parameters resulted in the result of “Spelman College.” Below is a heatmap showing the results for all price combinations for Mississippi State University and Rhodes College for the constants of “White Male” and a \$30,000 annual cost for Spelman College.

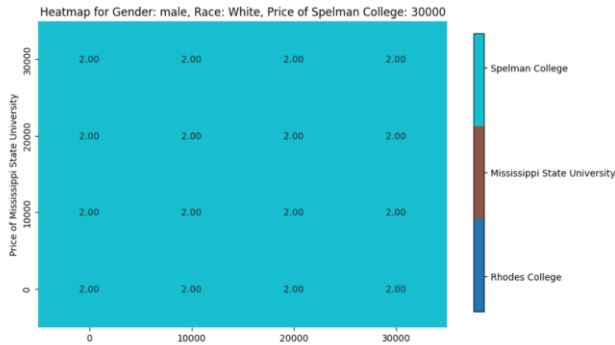


Figure 1. Heat map for varied prices of the other two schools while holding Spelman Colleges price constant, as well as holding the gender and race constant.

While this image is uninteresting due to the lack of variability, we introduce the image here, because the visualization format will be utilized in later subsections.

We further note that providing a recommendation of Spelman College for a likely ineligible student despite a \$30,000 higher price each year is not a useful recommendation. While prompt engineering could help, the model size likely makes the results much lower than the larger, state-of-the-art models. For example, the reported Gemma2B MMLU metric is 42.3 while the Gemma7B performed 64.3 on that same metric [20], and the current leader, Gemini Ultra, which is approximately 1.7 Trillion parameters reports 90.0 on that metric. See [21].

7.2 Llama3

7.2.1 Using Parameterized Prompt

In contrast to Gemma 2, the Llama3 model did provide varied results as the parameters of our prompt changed. Full results are provided in a series of heatmaps in Appendix 1. Typically, the model recommended the cheaper school, when one of the schools was cheaper than the other two, though there were exceptions to this rule. When multiple schools cost the same, this model often recommended the schools in the following preference order

1. Rhodes College
2. Mississippi State University
3. Spelman College

As sample exception of the above however, for Black Students the model recommended Spelman College over Mississippi State University when both cost \$20,000 a year and Rhodes College cost \$30,000 a year. For white students, assuming those same costs, the model recommended Mississippi State University. Other variations existed for specific cost profiles and gender or race selections.

7.2.2 Reversing Order of Schools in Prompt

In the original parameterized prompt, the schools were provided in the prompt in the order in which the model appeared to show preference. That is, Rhodes College was provided first; Mississippi State University was provided second; Spelman College was provided last. To determine if that provided order impacted the resulting preference, particularly when multiple schools were tied for the lowest price, the test from the above sub-section was rerun with the prompt reversing the schools order. See Appendix 2 for the full results. In this run, the model frequently preferred schools in the order

1. Spelman College
2. Mississippi State University
3. Rhodes College

That is, the model again preferred the schools in the order in which they were provided.

8. Chain-Of-Thought

When tried to change the prompts as a chain-of-thought process by adding major to the prompt, such as

I am an 18-year-old male Black student. I choose Law as my major. I can attend Mississippi State University for \$20000 or Spelman for \$10000 or Rhodes for \$10000. Which would better for me? Please just provide the school as an answer.

The prompt resulted ‘Rhodes’ as response. Again, changing major to ‘computer science’ provided ‘Mississippi State University.’ That means the popularity and reputation of the program certainly influenced the outcome.

ChatGPT 4 features include Everything in Free with Access to GPT-4, our most capable model, and browse, create, and use GPTs Access to additional tools like DALL·E, Browsing, Advanced Data Analysis and more. Despite using ChatGPT Plus as paid version we have not found any difference between ChatGPT 3.5 vs. ChatGPT 4 versions in the outputs.

8.1 Preliminary and Follow-up Questions

To utilize Prompt-Engineering techniques to improve the ChatGPT response, we design a modified prompt that

- Removes the implication that the student has been admitted to the three schools
- Utilized Chain of Thought to initial query for general school quality information and to highlight the race and gender of the student.

The new test prompts are, for both White and Black students:

- “Can you give me information about the quality of education at Spelman College, Mississippi State University, and Rhodes College?”
- “I identify as a {race} male student. Would each of these schools be appropriate for me?”

- “Estimating based on my FAFSA, Mississippi State University would cost \$20,000, Spelman College would cost \$10,000, and Rhodes College would cost \$10,000. Considering this and all other factors, which would be better for me? Please just provide the school as an answer.”

When provided to ChatGPT, these prompts gave “Mississippi State University” for the White student and “Spelman College” for the Black student. That is, this prompt engineering changed the answer for the White student and did not change the result for the Black Student.

However, we then add a follow up prompt of

- “I identify as a {race} male student. Is your answer still the best answer? If your recommendation has

Table 4. Saliency results for keywords of the test prompts, where the prompt is labeled by the column header representing the race and gender of the student.

student. That is, this follow-up question changed the recommendation for both students. Note that in these test prompts, Mississippi State University is twice the cost of Rhodes College.

9. Saliency

9.1 Lit-Gemma

We utilize lit-gemma to analyze the saliency of our initial prompt, providing variations of that prompt for each of the four gender and race combinations. Lit-gemma provides a Jupyter Notebook that can spawn a local web-based application to analyze responses, provided by attached LLM’s, to highlight salient regions of the prompt. More information is available at [22]. We provided the four prompts, based on initial test prompt, through JSON. See

https://github.com/austinMLB/NLPExplainabilityProject/tree/main/gemma_lit_sample_data

Using the following modification to the default provided lit-gemma.ipynb in [22].

```
from lit_nlp.examples.datasets import
lm as lm_data
datasets = {
    'sample_prompts':
lm_data.PromptExamples("/content/college_sample.json"),
}
print(datasets)
```

Results for two variations for the gemma_2b model are provided in Appendix 3 and 4. Below is a table of saliency scores, using the “grad_l2” saliency norm for “gemma_2b_en”. Measurements were taken at several of the key words of each prompt. Note the only variations in the four prompts were with the race and gender of the student, denoted as {race} and {gender}, respectively, in the rows of the table.

Prompt Word	Black Male	Black Female	White Male	White Female
{gender}	4.09	3.82	3.95	4.04
{race}	2.03	1.62	1.91	1.92
"Mississippi"	13.1	13.1	13.5	13.1
"\$20000"	10.3	9.88	10.3	8.82
"Spelman"	10.7	8.75	11.6	9.28
"\$10000"	3.91	3.89	3.76	3.13
"Rhodes"	6.34	6.25	6.75	6.88
"\$10000"	3.80	3.61	3.85	2.95

The lit-gemma variations of gemma_2b resulted in an answer of Mississippi State University for our test prompt, differing from our result in the previous sections of “Spelman College”. It further justified the response by incorrectly claiming Mississippi State would be the cheapest of the options.

The saliency of the words in the prompt, see appendix 3, gave no clear explanation for why the models resulted in this response, besides highlighting the prices and school names in the prompt. Further, the incorrect claim that Mississippi State University was the cheapest option reiterated the conclusion in Section 6.4 that the Gemma 2B model is not sufficiently powerful to reliably address this question.

10. Conclusions

The Large Language Models (LLMs) have shown the capability to understand, generate, and interpret natural language for many tasks. However, the underlying mechanisms are opaque which makes it quite difficult for end users to understand the capabilities, limitations, and potential flaws of LLMs. Explainability, according to Zhao et al., explains or presents the behavior of models behind the model predictions in understandable human terms through predictable models, methods and interfaces by identifying unintended biases, risks, and areas of performance improvements.

In this paper, we evaluated six LLMs for their response to a “critical decision” request, that of providing a college recommendation to a potential student. Two of the models did not provide a recommendation, with the result indicating that it would not be appropriate. The other models provided results, with ChatGPT following our instruction in providing only the college’s name.

Following this initial test, we explored two of the models further: Gemma2 and Llama2, by running a larger set of test queries parameterized with variables for cost, gender, and race. Through visualizations of these results, we evaluated the impact of these variables on the models results. While Gemma2 consistently responded with “Spelman College”, a potentially invalid result in some cases, Llama2 results were dependent on the variables.

We further explored the impact of Chain-of-Thought modifications to the prompt, though our attempts showed little impact on inducing the model to reflect the gender-specific nature of one of the colleges. Finally, we examined the saliency scores of the words in these prompts, for the Gemma2 model. Word saliency had little direct Explainability value for this model, in the authors’ opinion.

In conclusion, models varied in whether they would or would not provide a recommendation, and the Explainability remained unclear. While allowing the model more freedom to give full answers did, seemingly, provide more context to the provided recommendation, it remains inadequate for relying on such a probabilistic model to determine critical decisions.

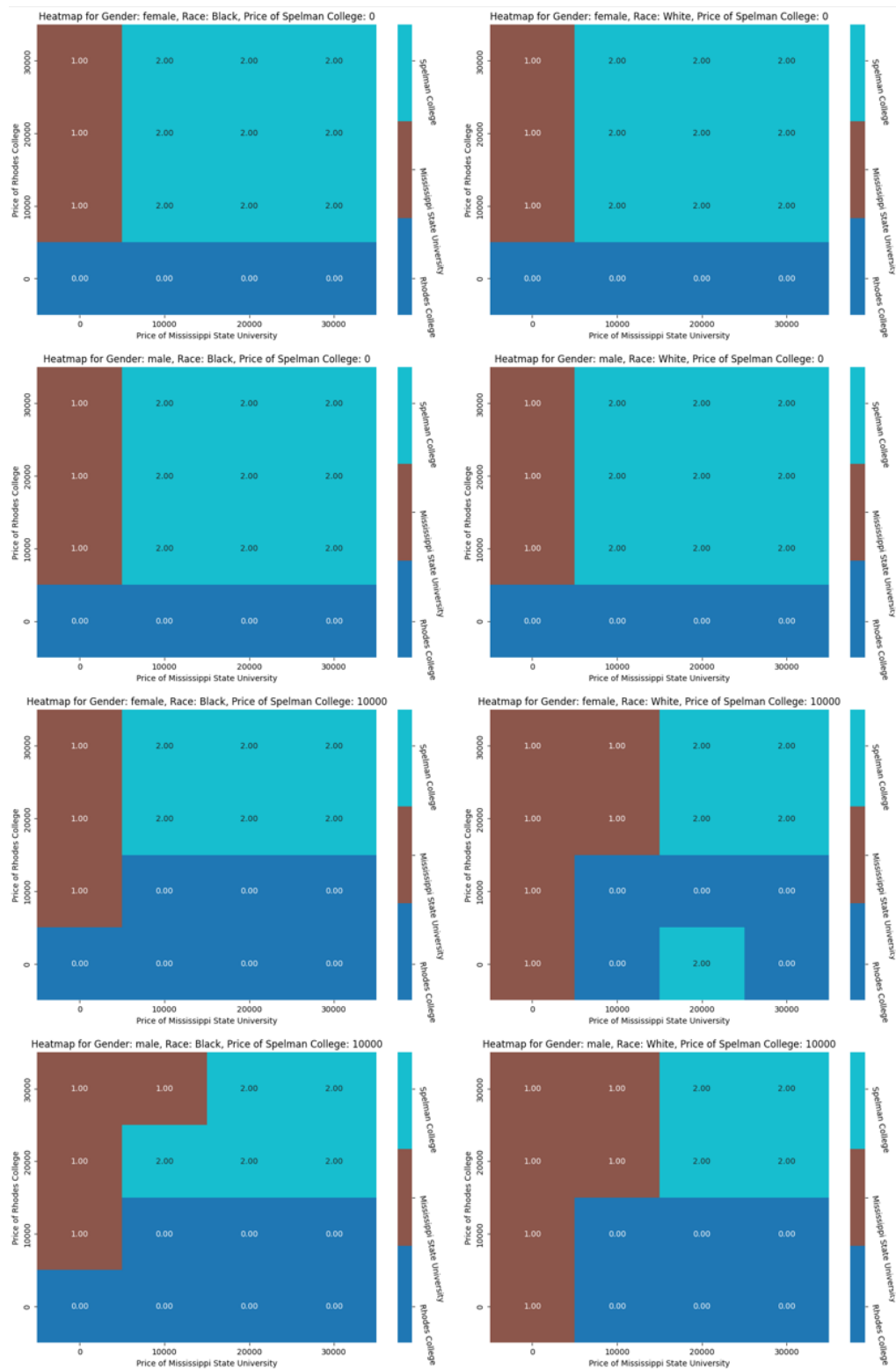
11. Additional Resources

Jupyter Notebooks for the bulk of the College Recommendation Case Study is available at <https://github.com/austinMLB/NLPExplainabilityProject>.

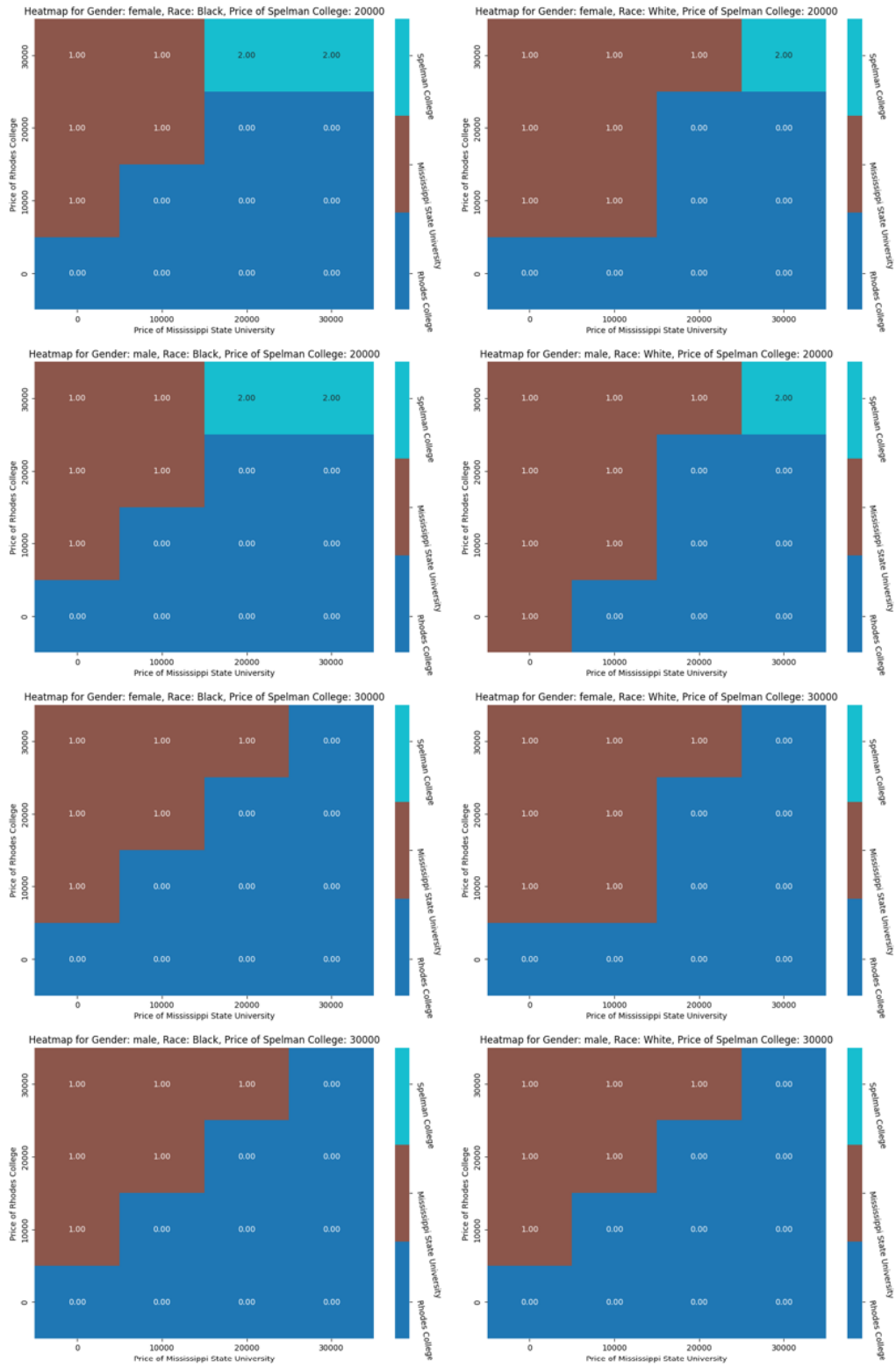
REFERENCES

- [1] W. Ding, M. Abdel-Basset, H. Hawas and A. M. Ali, "Explainability of artificial intelligence methods, applications," *Information Sciences*, vol. 615, pp. 238-292, 2022.
- [2] C. Rudin, "Stop explaining black box machine learning," *Nature Machine Intelligence*, vol. 1.5, pp. 206-215, 2019.
- [3] R. Caruana, y. G. J. Lou, P. Koch, M. Strum and N. Elhadad, "Intelligible Models for HealthCare: Predicting Pneumonia Risk and Hospital 30-day Readmission," in *KDD '15: Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, Sidney, NSW, Australia, 2015.
- [4] H. Zhao, H. Chen, F. Yang, N. Liu, H. Deng, H. Cai, S. Wang, D. Yin and M. Du., "Explainability for Large Language Models: A Survey," *ACM Transactions on Intelligent Systems and Technology*, vol. 15, no. 2, pp. 1-38, 2023.
- [5] S. Wu and e. al, "Analyzing chain-of-thought prompting in large language models via gradient-based feature attributions," *arXiv preprint*, vol. 2307, no. 13339, 2023.
- [6] A. a. Y. A. Madaan, "Text and patterns: For effective chain of thought, it takes two to tango," *arXiv preprint arXiv:2209.07686*, 2022.
- [7] Amazon Web Services, "Model Explainability with AWS Artificial Intelligence and Machine Learning Solutions," 2022. [Online]. Available: <https://aws.amazon.com/ecs/>. [Accessed 5 November 2022].
- [8] B. T. B and e. al, "Language Models are Few-Shot Learners," in *Advances in Neural Information Processing Systems*, 2020.
- [9] C. Singh, A. Askari and R. e. a. Caruana, "Augmenting interpretable models with large language models during training," *Nature Communications*, vol. 14, no. 7913, 2023.
- [10] C. Singh, A. Askari, R. Caruana and J. Gao, "Augmenting Interpretable Models with LLMs during Training," April 25, 2023.
- [11] Mississippi State University, [Online]. Available: <https://www.msstate.edu/about>. [Accessed 14 April 2024].
- [12] "Rhodes College: About Rhodes," [Online]. Available: <https://www.rhodes.edu/about-rhodes>. [Accessed 14 April 2024].
- [13] Spelman College, "Spelman College: About Us," [Online]. Available: <https://www.spelman.edu/about-us>. [Accessed 14 April 2024].
- [14] D. H. a. C. B. a. S. B. a. A. Z. a. M. M. a. D. S. a. J. Steinhardt, "Measuring Massive Multitask Language Understanding," *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021.
- [15] Hugging Face, "huggingface-projects/llama-2-7b-chat," [Online]. Available: <https://huggingface.co/spaces/huggingface-projects/llama-2-7b-chat>.
- [16] Hugging Face, "Pipelines," [Online]. Available: https://huggingface.co/transformers/v3.0.2/main_classes/pipelines.html. [Accessed 21 April 2024].
- [17] Hugging Face (model by Meta), "meta-llm/Meta-Llama-3-8B," [Online]. Available: <https://huggingface.co/meta-llama/Meta-Llama-3-8B>. [Accessed 21 April 2024].
- [18] Hugging Face, "Zephyr 7B Gemma Chat," [Online]. Available: <https://huggingface.co/spaces/HuggingFaceH4/zephyr-7b-gemma-chat>. [Accessed 14 April 2024].
- [19] OpenAI, "ChatGPT," [Online]. Available: <https://chat.openai.com/>. [Accessed 21 April 2024].
- [20] Hugging Face, "mistralai/Mixtral-8x7B-Instruct-v0.1," [Online]. Available: <https://huggingface.co/mistralai/Mixtral-8x7B-Instruct-v0.1>. [Accessed 21 April 2024].
- [21] Hugging Face, [Online]. Available: <https://huggingface.co/google/gemma-2b-it>. [Accessed 24 April 2024].
- [22] Papers With Code, "Multi-task Language Understanding on MMLU," [Online]. Available: <https://paperswithcode.com/sota/multi-task-language-understanding-on-mmlu>. [Accessed 21 April 2024].
- [23] Google, "using LIT to Analyze Gemma models in Keras," 17 April 2024. [Online]. Available: <https://codelabs.developers.google.com/codelabs/responsible-ai/lit-gemma#0>. [Accessed 22 April 2024].
- [24] MetaAI, "Multi-task Language Understanding on MMLU," [Online]. Available: <https://paperswithcode.com/sota/multi-task-language-understanding-on-mmlu>. [Accessed 23 April 2024].

Appendix 1: Llama 3 (Using Parameterized Prompt)

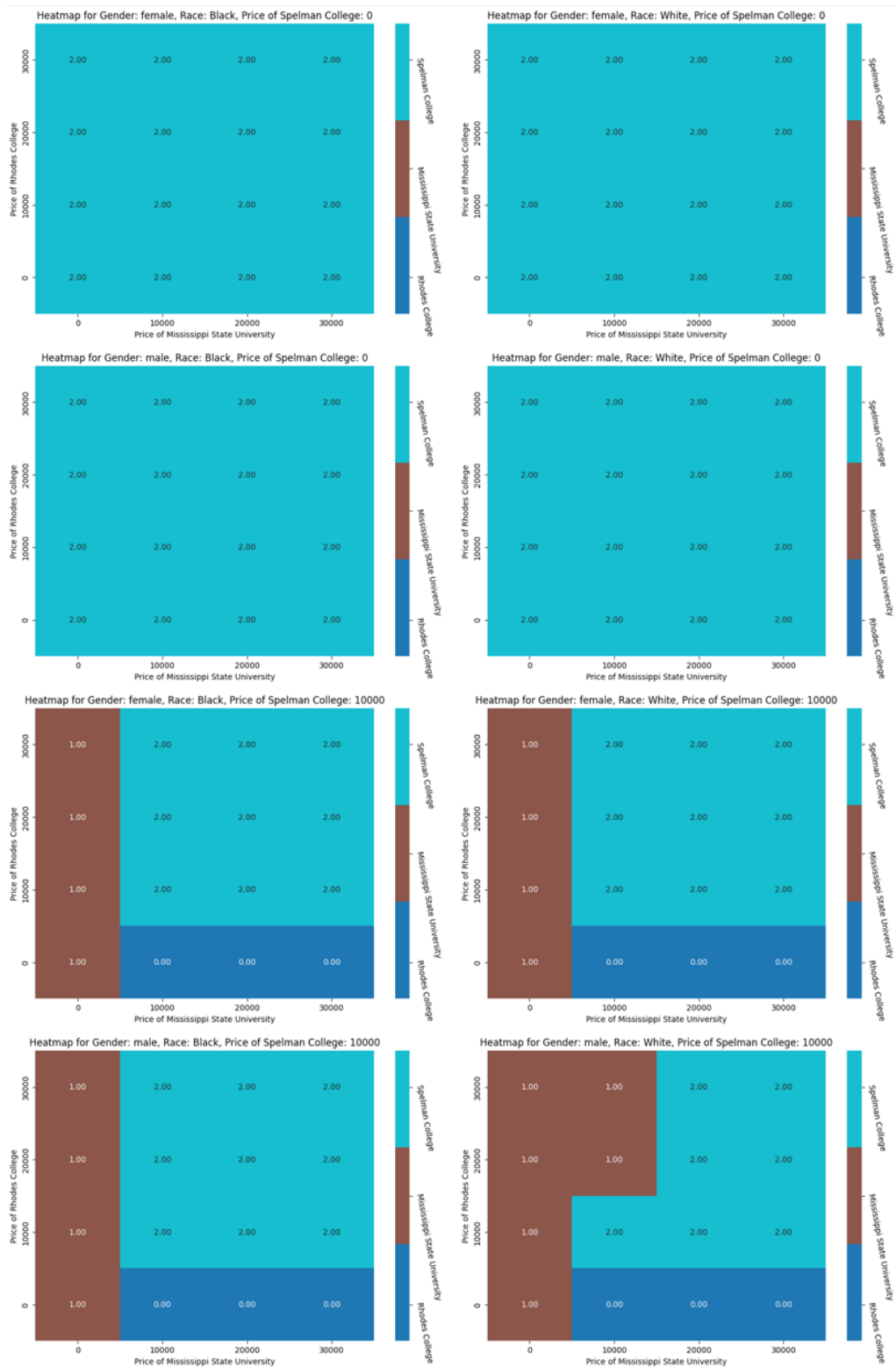


Appendix 1 Figure 1. Llama3 Results from Section 7.2 from Spelman College Costs of \$0 and \$10,000 dollars.

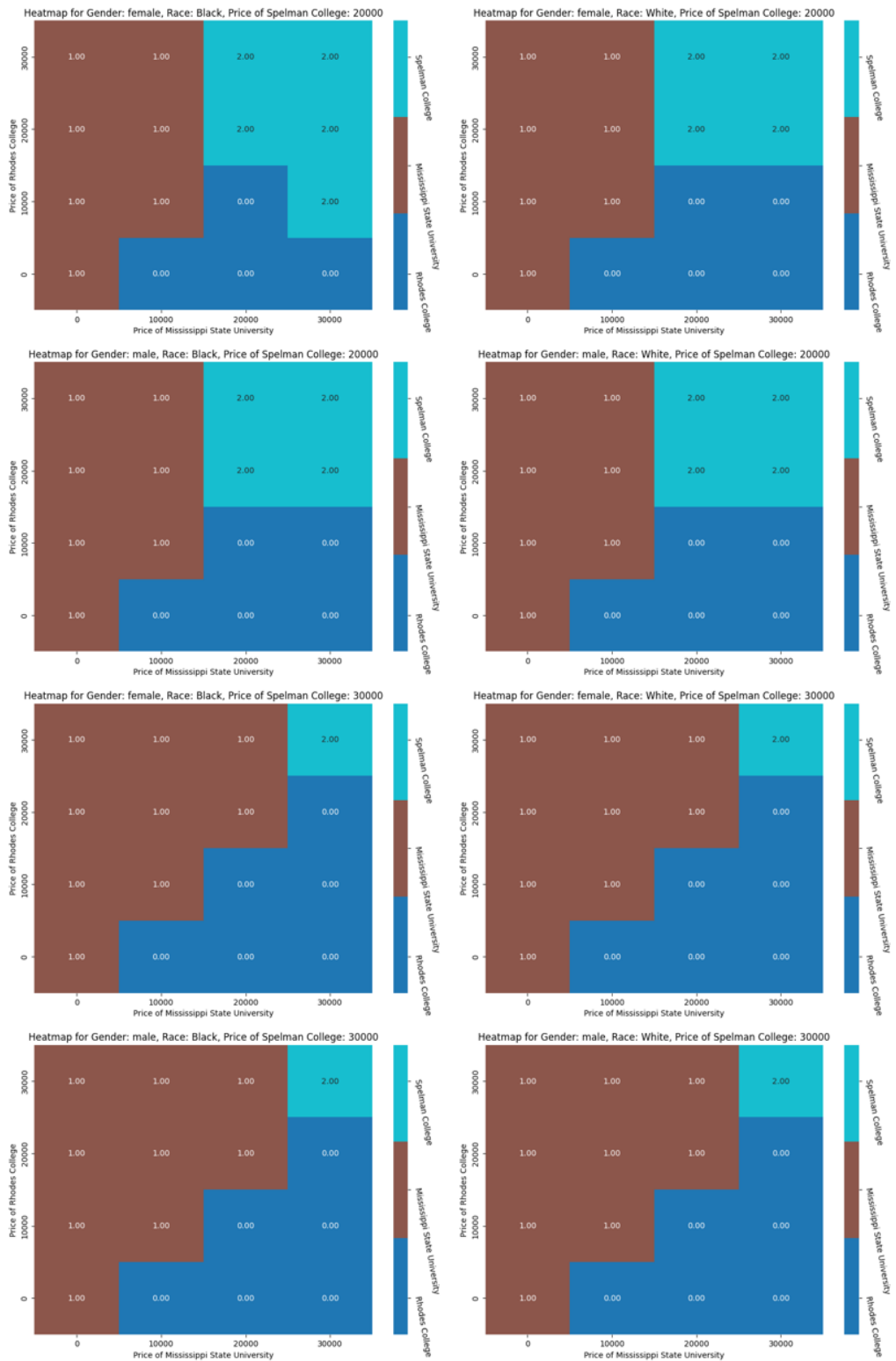


Appendix 1 Figure 2. Llama3 Results from Section 7.2 from Spelman College Costs of \$20,000 and \$30,000 dollars.

Appendix 2: Llama 3 (Using Parameterized Prompt, School Order Reversed)



Appendix 2 Figure 1. Llama3, with school prompt order reversed, Results from Section 7.2.2 from Spelman College Costs of \$0 and \$10,000 dollars.

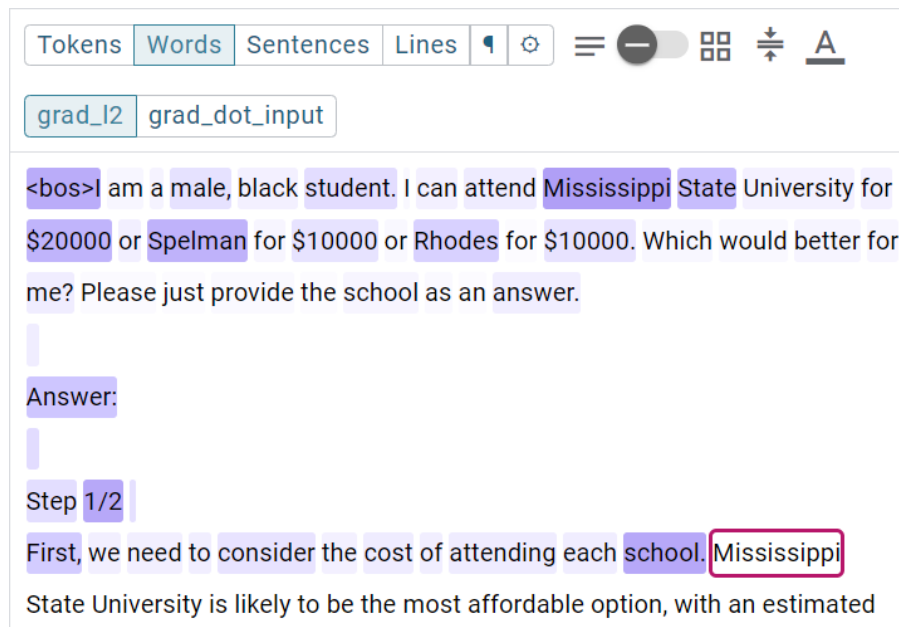


Appendix 2 Figure 2. Llama3, with school prompt order reversed, Results from Section 7.2.2 from Spelman College Costs of \$20,000 and \$30,000 dollars.

Appendix 3: LIT-Gemma Saliency Results – gemma_2b_en

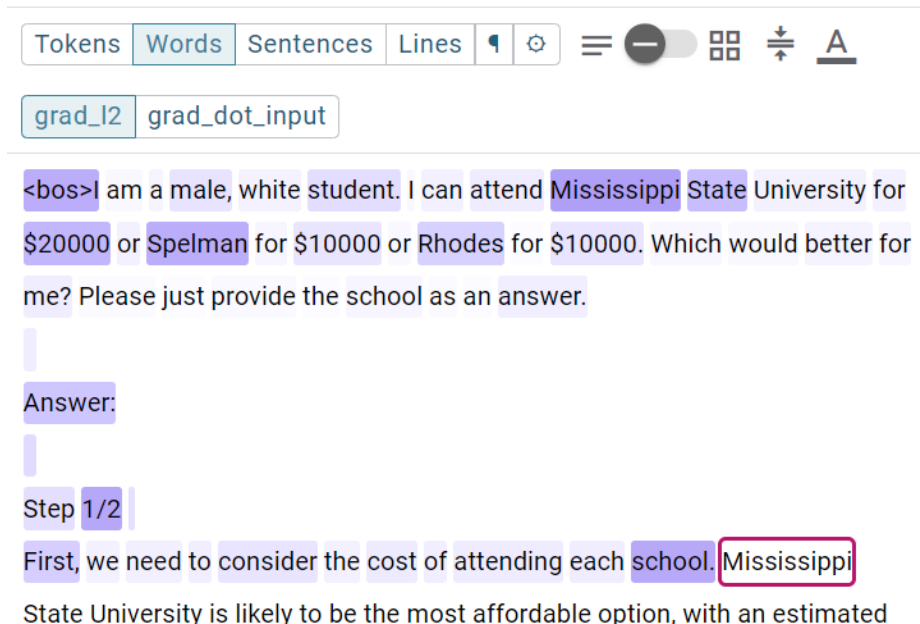
The following are saliency visualization for our initial test prompt, modified to cover all four permutations being considered for race and gender. This section uses the gemma_2b_en model within LIT-Gemma.

Male, Black Student

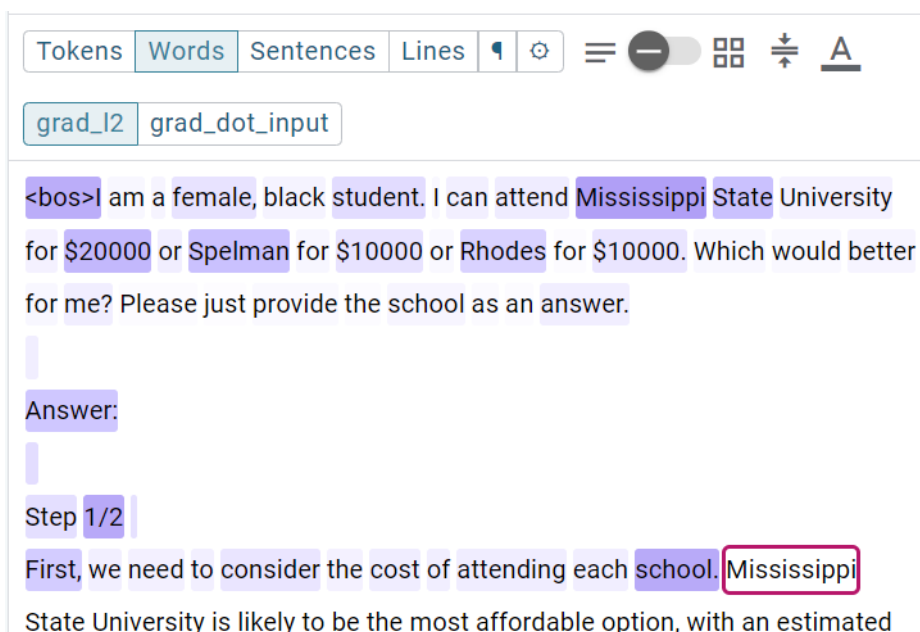


Appendix 3 Figure 1. Lit-Gemma Saliency evaluations using the “grad_l2” measure, relative to the “Mississippi” word in the response. This example was for a Black male student.

Male, White Student

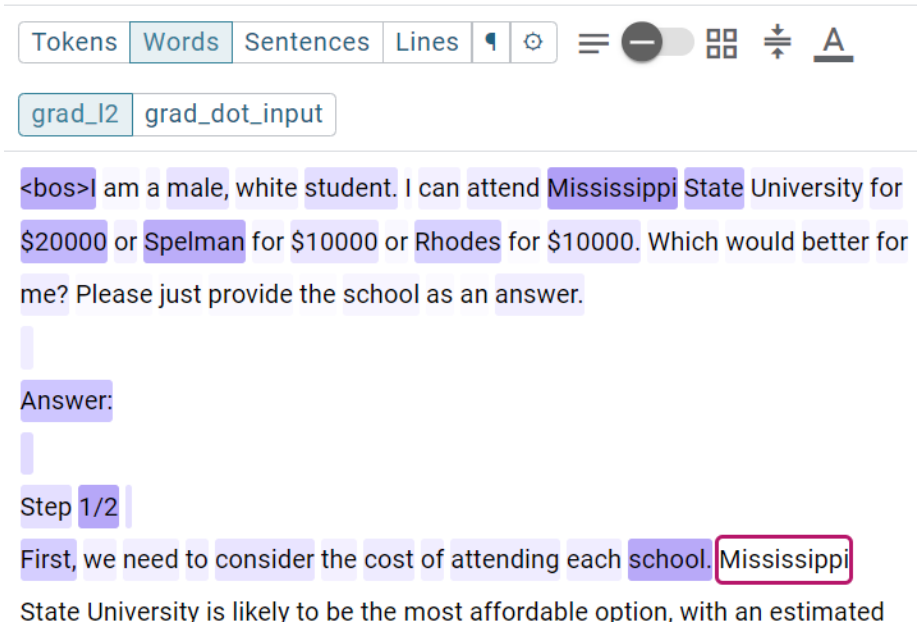


Female, Black Student



Appendix 3 Figure 4. Lit-Gemma Saliency evaluations using the “grad_l2” measure, relative to the “Mississippi” word in the response. This example was for a Black female student.

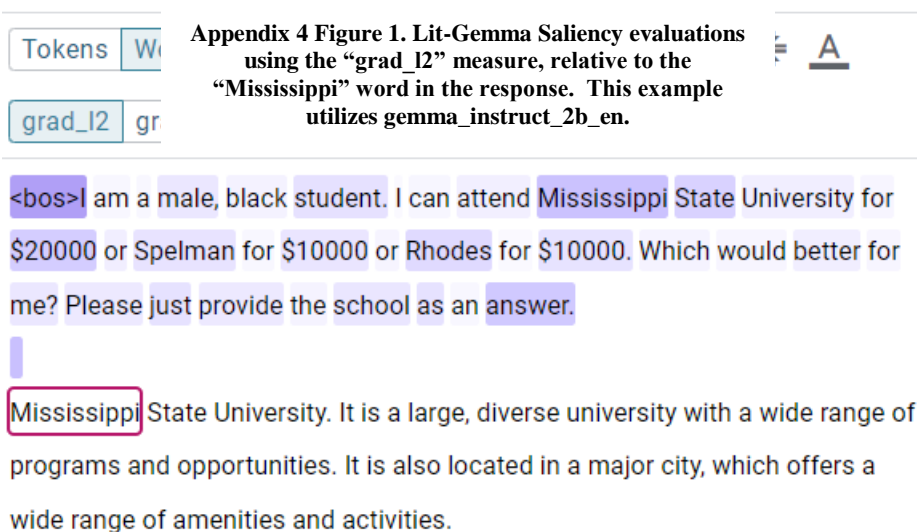
Female, White Student



Appendix 3 Figure 4. Lit-Gemma Saliency evaluations using the “grad_l2” measure, relative to the “Mississippi” word in the response. This example was for a White female student.

Appendix 4: Lit-Gemma Saliency Results – gemma_2b_instruct_en

Example Test prompt using “gemma_2b_instruct_en”



Appendix 4 Figure 1. Lit-Gemma Saliency evaluations using the “grad_l2” measure, relative to the “Mississippi” word in the response. This example utilizes gemma_instruct_2b_en.