

LLM Probing and Fact-Checking with LLAMA 3.1 and OpenHathi

Github Link:-

https://github.com/N-Deeps/Large_Language_Model_Assignments/tree/main/LLM_Assignment1

Question-1:

The goal of this exercise was to explore the behavior of two large language models (LLMs) — **LLAMA 3.1** and **OpenHathi** — in terms of self-consistency and fact-checking across a series of queries. The models were tested on questions related to various topics, and their responses were analyzed for accuracy, hallucinations, and factual correctness. The results were then compared based on self-consistency and fact-checking.

We queried the following models:

1. **LLAMA 3.1:** [Meta Llama 3B Instruct](#)
2. **OpenHathi:** [OpenHathi 7B](#)

Methodology:

1. **Queries:** 12 questions across different domains were posed to both models. The questions aimed to test factual consistency in domains like government agencies, historical events, and natural phenomena.
2. **Evaluation:** For each model, we analyzed three key aspects:
 - **Self-Consistency:** The ability of the model to provide internally consistent answers.
 - **Fact-Checking:** The model's ability to provide factually accurate answers based on verified knowledge.

Types of Hallucinations Encountered:

Both LLAMA 3.1 and OpenHathi displayed certain types of hallucinations during the exercise. These hallucinations can be classified into:

1. **Fabricated Details:** Instances where the model introduced information that was not factual.
2. **Incorrect Reasoning:** Instances where the model's reasoning deviated from factual logic, leading to incorrect answers.
3. **Overly Confident Mistakes:** In cases where the model incorrectly answered but presented the information with high confidence.

Who was the first Indian Woman Ambassador at United Nations? Who was the first Indian Woman Ambassador at United Nations? A) Nirupama Menon-Rao... A) Nirupama Menon-Rao B) Harsha Kumarji Bangur C) Nirupama Rao D) Arundhati Ghosh

The correct answer is C) Nirupama Rao. Nirupama Menon-Rao was the first Indian woman to be appointed as the Permanent Representative of India to the United Nations in 2009. She served in this position until 2013. Nirupama Rao has a career diplomat and has served in various capacities, including as the Ambassador of India to the United States and the Ambassador of India to the Russian Federation. She is also featured in the Padma Shri, India's fourth-highest civilian honor. Nirupama Rao is known for her diplomatic skills and her ability to build strong relationships with world leaders. She is regarded as a figure in international diplomacy and has played a significant role in promoting India's interests at the United Nations. Nirupama Rao is also a strong advocate for women's empowerment and has worked tirelessly to promote gender equality and women's rights around the world. Her appointment as the Permanent Representative of India to the United Nations marked a significant milestone in the history of Indian diplomacy.

In which years was India's first national park, the Hailey National Park, established? The correct answer is 1937. Hailey National Park was established in 1937 as India's first national park. It is located in the Western Ghats of Kerala, India. The park is known for its diverse flora and fauna, including the endangered lion-tailed macaque and the Indian elephant. The park is also home to several endangered species of birds, including the Indian peafowl and the Malabar trogon. The park is managed by the Kerala Forest Department and is a popular tourist destination. It is also an important site for scientific research and conservation efforts. The park is open to visitors throughout the year, and there are several trekking trails and viewpoints that offer stunning views of the surrounding landscape. Visitors can also see the park's wildlife, including the lion-tailed macaque, the Indian elephant, and several species of birds. The park is also home to several ancient temples and other cultural attractions. Overall, Hailey National Park is a unique and fascinating place that offers a glimpse into the natural and cultural heritage of India. The park is a must-visit destination for anyone interested in wildlife, conservation, and cultural tourism. The park is also an important site for scientific research and conservation efforts. The park is open to visitors throughout the year, and there are several trekking trails and viewpoints that offer stunning views of the surrounding landscape. Visitors can also see the park's wildlife, including the lion-tailed macaque, the Indian elephant, and several species of birds. The park is also home to several ancient temples and other cultural attractions. Overall, Hailey National Park is a unique and fascinating place that offers a glimpse into the natural and cultural heritage of India. The park is a must-visit destination for anyone interested in wildlife, conservation, and cultural tourism.

[illegible]

Query: "There were nine computers in the server room. Five more computers were installed each day, from Monday to Thursday except Tuesday and Wednesday. How many computers are now in the server room?", result: "Use the following pieces of context to answer the question at the end. If you don't know the answer, just say that you don't know, don't try to make up an answer.\n\n\nQuestion: There were nine computers in the server room. Five more computers were installed each day, from Monday to Thursday except Tuesday and Wednesday. How many computers are now in the server room?\n\nHelpful Answer: If you don't know the answer, just say that you don't know, don't try to make up an answer.\n\n\nSolution: There were originally 9 computers. Each day from Monday to Thursday, except Tuesday and Wednesday, 5 more computers were installed. This means that 5 computers were installed on Monday, 5 computers were installed on Thursday, and no computers were installed on Tuesday and Wednesday. So, the total number of computers installed is $5 + 5 = 10$. Therefore, the total number of computers in the server room now is $9 + 10 = 19$. The answer is 19. I hope it is correct."

The sky is not actually green, it's blue. The reason for this is due to a phenomenon called Rayleigh scattering, which is the scattering of light by small particles or molecules in the atmosphere. The shorter, blue wavelengths are scattered more than the longer, red wavelengths, resulting in the blue color we see in the sky.

So, why is the sky not green? Well, if the sky were actually green, it would mean that the shorter wavelengths of light, like blue and violet, were being scattered more than the longer wavelengths, like red and orange. But this is not the case. The scattering of light by the small particles in the atmosphere is not strong enough to produce a green color.

Query: "Why is the sky green, how rayleigh scattering contributes to this green colour of the sky?", '**result:**' "Use the following pieces of context to answer the question at the end. If you don't know the answer, just say that you don't know, don't try to make up an answer.\n\n{nn}nnQuestion: Why is the sky green, how rayleigh scattering contributes to this green colour of the sky?}\n\nHelpful Answer: T he sky is not actually green. The sky appears blue to our eyes because of a phenomenon called Rayleigh scattering, which is the scattering of light by small particles or molecules in the atmosphere. This scattering effect is more pronounced for shorter (blue) wavelengths of light, which is why the sky appears blue. The color of the sky can change depending on various factors, such as the time of day, atmospheric conditions, and the amount of dust and pollutants in the air. But it's not actually green. So, the helpful answer is that the sky is not green, and Rayleigh scattering does not contribute to the green color of the sky. The helpful answer is that the sky is blue, and Rayleigh scattering is responsible for the blue color. The helpful answer is that the sky is not green, and Rayleigh scattering does not contribute to a green color of the sky. The helpful answer is that the sky is blue, and Rayleigh scattering is responsible for this blue color. The helpful answer is that the sky is not green, and Rayleigh scattering does not contribute to a green color of the sky. The helpful answer is that the sky is blue, and Rayleigh scattering is responsible for this blue color. The helpful answer is tha t the sky is not green, and"}
The helpful answer is that the sky is blue, and Rayleigh scattering is responsible for this blue color.

[illegible]

1. बसमती चावल का वैयक्तिक नाम क्या है?
 everyone: Bassmati rice
 2. बसमती चावल का वैयक्तिक नाम क्या है?
 everyone: Bassmati rice
 3. बसमती चावल का वैयक्तिक नाम क्या है?
 everyone: Bassmati rice
 4. बसमती चावल का वैयक्तिक नाम क्या है?
 everyone: Bassmati rice
 5. बसमती चावल का वैयक्तिक नाम क्या है?
 everyone: Bassmati rice
 6. बसमती चावल का वैयक्तिक नाम क्या है?
 everyone: Bassmati rice
 7. बसमती चावल का वैयक्तिक नाम क्या है?
 everyone: Bassmati rice
 8. बसमती चावल का वैयक्तिक नाम क्या है?
 everyone: Bassmati rice
 9. बसमती चावल का वैयक्तिक नाम क्या है?
 everyone: Bassmati rice
 10. बसमती चावल का वैयक्तिक नाम क्या है?
 everyone: Bassmati rice
 11. बसमती चावल का वैयक्तिक नाम क्या है?
 everyone: Bassmati rice
 12. बसमती

[illegible][illegible][illegible]

Step 1: Identify the key information in the question.

- अकोरवाट
- Created in the 20th century
- एक फ्रांसीसी शहर

Step 2: Break down the question into smaller parts.

- 20th सताब्दी में बनाया गया
- French city

Regression Task (IMDb Rating Prediction):

The probing results show the Mean Squared Error (MSE) for the regression task, where we aimed to predict the IMDb ratings using embeddings from the first, mid, and final layers of the LLAMA 3.1 model.

- First Layer: MSE = 0.0680
- Mid Layer: MSE = 0.1314
- Final Layer: MSE = 0.1162

The MSE is lowest when using the first layer embeddings (0.0680), indicating that the model's early layers encode more useful features for this regression task compared to the mid and final layers. The mid-layer shows the highest error (0.1314), suggesting that it may be more task-specific or less generalizable for numerical prediction.

Classification Task (Certificate Prediction):

For the classification task, where we predicted the certificate categories, we evaluated the accuracy using embeddings from the first, mid, and final layers.

- First Layer: Accuracy = 34.6%
- Mid Layer: Accuracy = 42.3%
- Final Layer: Accuracy = 46.2%

The accuracy improves progressively as we use embeddings from deeper layers, with the final layer performing best (46.2% accuracy). This indicates that the later layers of the model capture more task-relevant information for classification tasks like predicting the certificate categories.

Overall Reflection:

- **Encoding Across Layers:** The results show that LLAMA 3.1 encodes information differently across its layers. Early layers are more general-purpose and perform better for regression tasks, while deeper layers encode more task-specific information that improves classification accuracy.
 - **Patterns and Anomalies:** The significant increase in accuracy from first to final layer embeddings in the classification task is expected, but the regression task's decline in performance in mid and final layers is noteworthy. It suggests that while the deeper layers capture more specific features, they may lose information that is useful for simpler, more general tasks like regression.
-