# 24KIDS442 XAI LLM LAB

## Stefi Shobika Sukumar

Understanding how large language models (LLMs) process and reason about text is a major challenge in modern AI research — especially when these models are deployed in black-box settings, where internal decision processes are hidden from users.

As part of the XAI LLM Lab Internship, this project focuses on building a lightweight, explainable interface for small, quantized LLMs (≤7B) that can run on CPU devices. The goal is to help users explore how the model processes input prompts, by combining visual neuron activation maps with natural language rationale explanations.
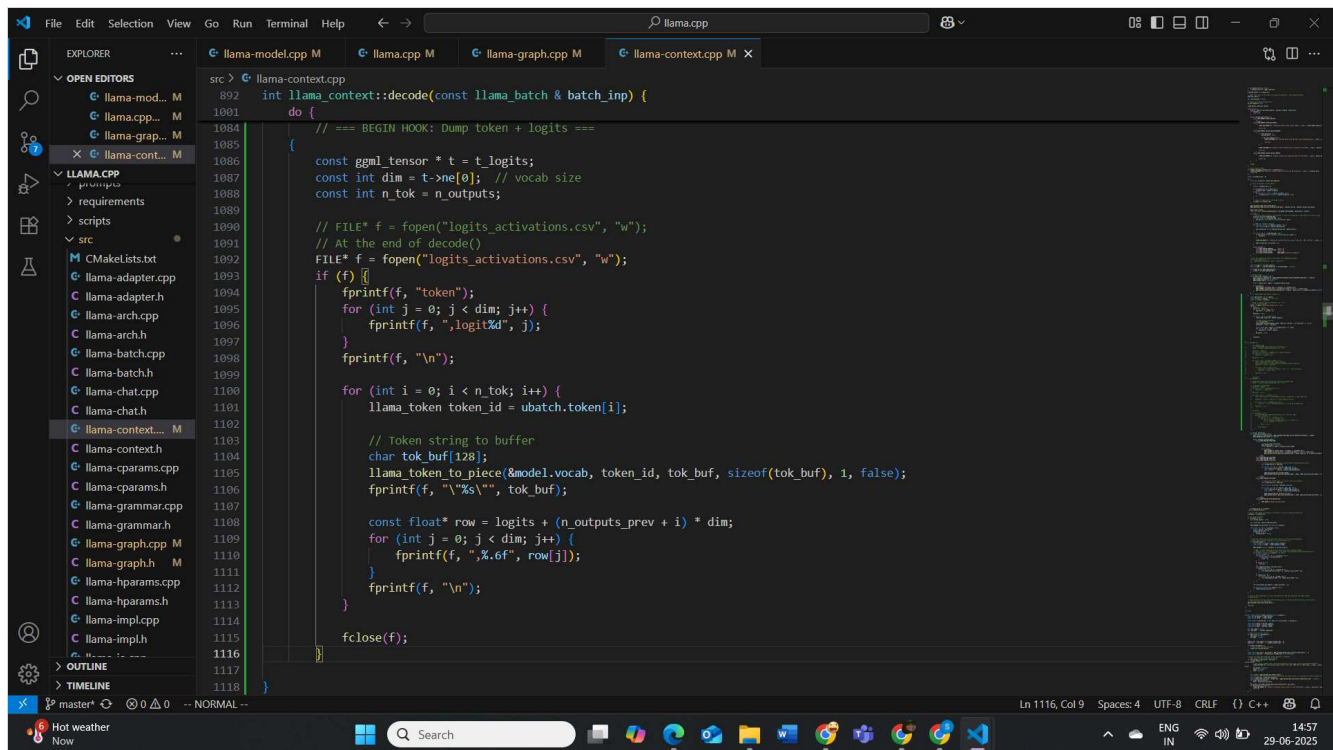
## System Architecture and Methodology

This project is implemented as a two-stage pipeline:

1. **Offline Backend (Computation):** Prompts are executed using quantized GGUF models (TinyLlama Q3_K_M and Phi Q4_K_M) in `llama.cpp`. A custom-modified decoder caches layer-wise neuron activations into `logits_activations.csv`.

2. **Token Explanation:** The Zephyr model (Mistral-based) is used to generate natural language rationales for each token. These are saved in `rationales.jsonl`.

3. **Streamlit Dashboard:** The app loads precomputed activations and rationales. When the user enters a prompt, it retrieves the cached data and visualizes the token-wise attention using Plotly heatmaps alongside the generated rationales.

This decoupled design allows for CPU-efficient local analysis and avoids the need to re-run inference live inside the Streamlit app.

**Backend Pipeline (Screenshots)**

*Loading quantized GGUF models with llama.cpp*

```
ophia_ona@Sophia:/mnt/c/Users/Sophia Sona/llama.cpp/build$ ./bin/llama-cli -m ../models/tinyllama-1.1b-chat-v1.0.Q3_K_M.gguf -p "what is a black hole?"
build: 5716 (d27b3ca1) with cc (Ubuntu 13.3.0-6ubuntu2~24.04) 13.3.0 for x86_64-linux-gnu
main: llama backend init
main: load the model and apply lora adapter, if any
llama_model_loader: loaded meta data with 23 key-value pairs and 201 tensors from ../models/tinyllama-1.1b-chat-v1.0.Q3_K_M.gguf (version GGUF V3 (latest))
llama_model_loader: Dumping metadata keys/values. Note: KV overrides do not apply in this output.
llama_model_loader: - kv   0:                       general.architecture str              = llama
llama_model_loader: - kv   1:                               general.name str              = tinyllama_tinyllama-1.1b-chat-v1.0
llama_model_loader: - kv   2:                     llama.context_length u32              = 2048
llama_model_loader: - kv   3:                   llama.embedding_length u32              = 2048
llama_model_loader: - kv   4:                        llama.block_count u32              = 22
llama_model_loader: - kv   5:                  llama.feed_forward_length u32              = 5632
llama_model_loader: - kv   6:                 llama.rope.dimension_count u32              = 64
llama_model_loader: - kv   7:                 llama.attention.head_count u32              = 32
llama_model_loader: - kv   8:              llama.attention.head_count_kv u32              = 4
llama_model_loader: - kv   9:        llama.attention.layer_norm_rms_epsilon f32           = 0.000010
llama_model_loader: - kv  10:                      llama.rope.freq_base f32              = 10000.000000
llama_model_loader: - kv  11:                       general.file_type u32              = 12
llama_model_loader: - kv  12:                     tokenizer.ggml.model str              = llama
llama_model_loader: - kv  13:                    tokenizer.ggml.tokens arr[str,32000]   = ["<unk>", "<s>", "</s>", "<0x00>", "<...
llama_model_loader: - kv  14:                    tokenizer.ggml.scores arr[f32,32000]   = [0.000000, 0.000000, 0.000000, 0.0000...
llama_model_loader: - kv  15:                tokenizer.ggml.token_type arr[i32,32000]   = [2, 3, 3, 6, 6, 6, 6, 6, 6, 6, 6, 6, ...
llama_model_loader: - kv  16:                    tokenizer.ggml.merges arr[str,61249]   = ["_ t", "e r", "i n", "_ a", "e n...
llama_model_loader: - kv  17:                tokenizer.ggml.bos_token_id u32            = 1
llama_model_loader: - kv  18:                tokenizer.ggml.eos_token_id u32            = 2
llama_model_loader: - kv  19:            tokenizer.ggml.unknown_token_id u32            = 0
llama_model_loader: - kv  20:            tokenizer.ggml.padding_token_id u32            = 2
llama_model_loader: - kv  21:                    tokenizer.chat_template str            = {% for message in messages %}\n{% if m...
llama_model_loader: - kv  22:               general.quantization_version u32            = 2
llama_model_loader: - type  f32:    45 tensors
llama_model_loader: - type q3_K:    89 tensors
llama_model_loader: - type q4_K:    62 tensors
llama_model_loader: - type q5_K:     4 tensors
llama_model_loader: - type q6_K:     1 tensors
print_info: file format = GGUF V3 (latest)
```

```
 <|user|>
what is a black hole?
<|assistant|>
A black hole is a region in space where the gravitational pull is so strong that nothing, not even light, can escape its pull. It is a theoretical object that ha
s never been observed directly. The name "black hole" is derived from the idea that the object's matter is black, or very dark, due to the strong gravitational p
ull.

>
```

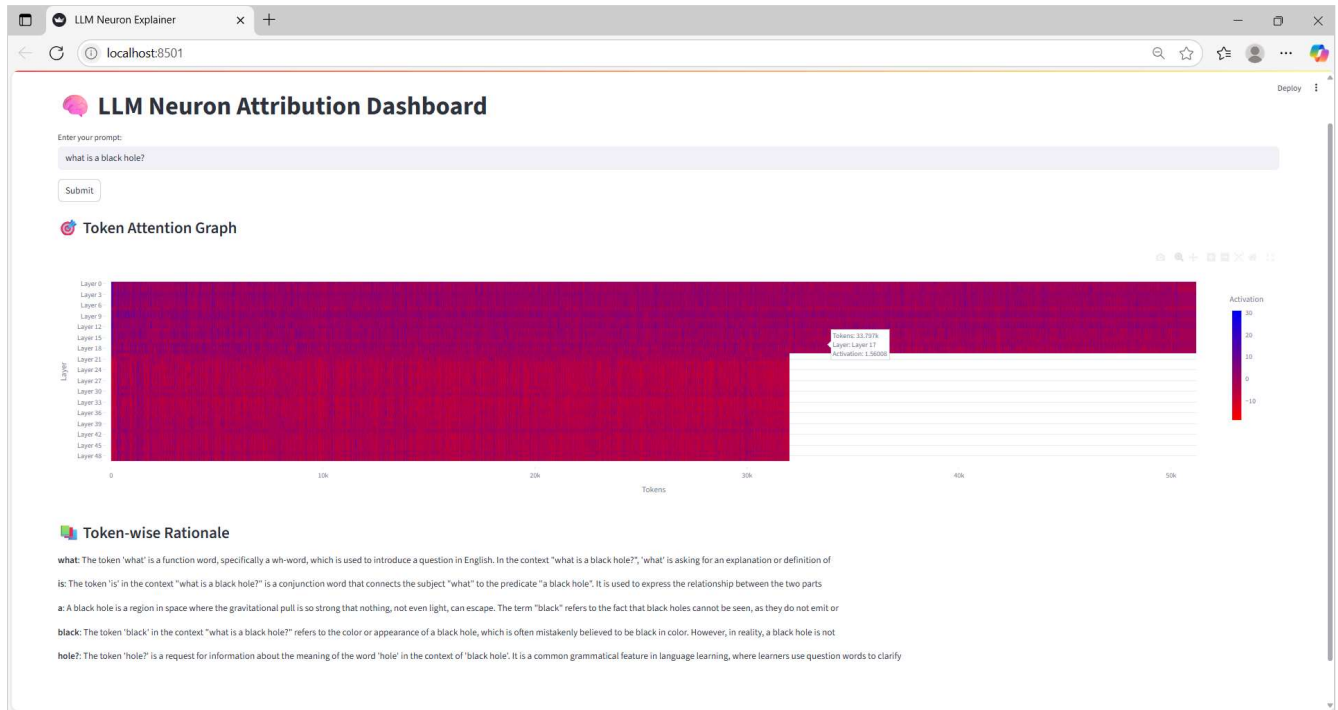*Quantized TinyLlama and Phi models running in CPU with llama.cpp*

*Neuron activations cached in logits_activations.csv*

*Zephyr generates token-level rationales for each prompt*

## Streamlit Visualization Dashboard

- User enters a prompt into the dashboard.
- The app looks up cached activations and explanations.
- Plotly is used to render the token-by-layer attention heatmap.
- Below the heatmap, Zephyr-generated rationales are shown for each token.

This architecture ensures the system remains CPU-friendly, while still offering rich interpretability features for education, debugging, and LLM transparency.