# Problem Statement

- Existing telemedicine have limited patient engagement and real-time artificial intelligence for documentation and clinical decision support. This leads to delayed diagnosis and significant gaps in providing efficient healthcare

# Literature Review

**1. The Impact of multimodal LLM on healthcare's future**

The current LLM in medicine and healthcare is unimodal and can accept and process text queries. Considering the multimodal nature of  Medicine and healthcare, multimodal LLMS that can handle texts and analyze images, videos, and sounds are ideal and would represent a significant achievement in the field of medicine (Mesko, 2023). The researcher further explains the hypothetical benefits of multimodal LLMS accepting patient text queries, analyzing their images of lesions and x-rays. Moreover, such an LLMS can potentially analyze vocal, lung, and heart sounds for abnormalities and generate comprehensive information and recommendations about the patient

**2. Applying Object Detection and Large Language Model to Establish a Smart Telemedicine Diagnosis System with Chatbot: A Case Study of Pressure Injuries Diagnosis System.**
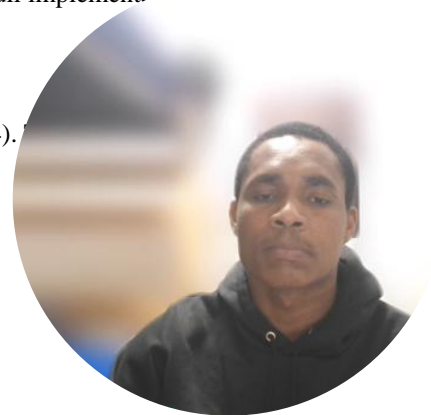
Assessing patients and pressure wounds without physical examination is a major challenge confronting Telemedicine (Chen et al., 2024). The researchers leveraged an auxiliary diagnosis system using YOLOv7 and LLM for enhanced detection and classification of pressure wounds and diagnosis. This approach generated an average F1 score of 0.9238, providing remote real-time diagnosis and staging of wounds.

**3. Transformer Models in Healthcare: A Survey and Thematic Analysis of Potentials, Shortcomings, and Risks**

Large language models have significant benefits in the healthcare sector, such as optimized processes, improved operational efficiency, and efficient clinical documentation (Denecke et al., 2024). Despite these benefits, the study uncovered potential challenges, such as privacy challenges, model development bias, and auditability issues hindering the full implementation of LLM in healthcare.

**4. Transforming virtual healthcare: The potential of ChatGPT-4omni in telemedicine**

The inception of ChatGPT-4o is a significant achievement in virtual healthcare and telemedicine, allowing users to upload medical images (Mohamad-Hani et al., 2024). argue that such milestones would improve remote patient care and access to personalized treatment plans.

# Purpose of the Code

- The purpose of the code is to develop AI-powered Telemedicine leveraging LLM to;

- Generate answers to patient queries

- Enhance patient engagement,

- Improve healthcare accessibility and
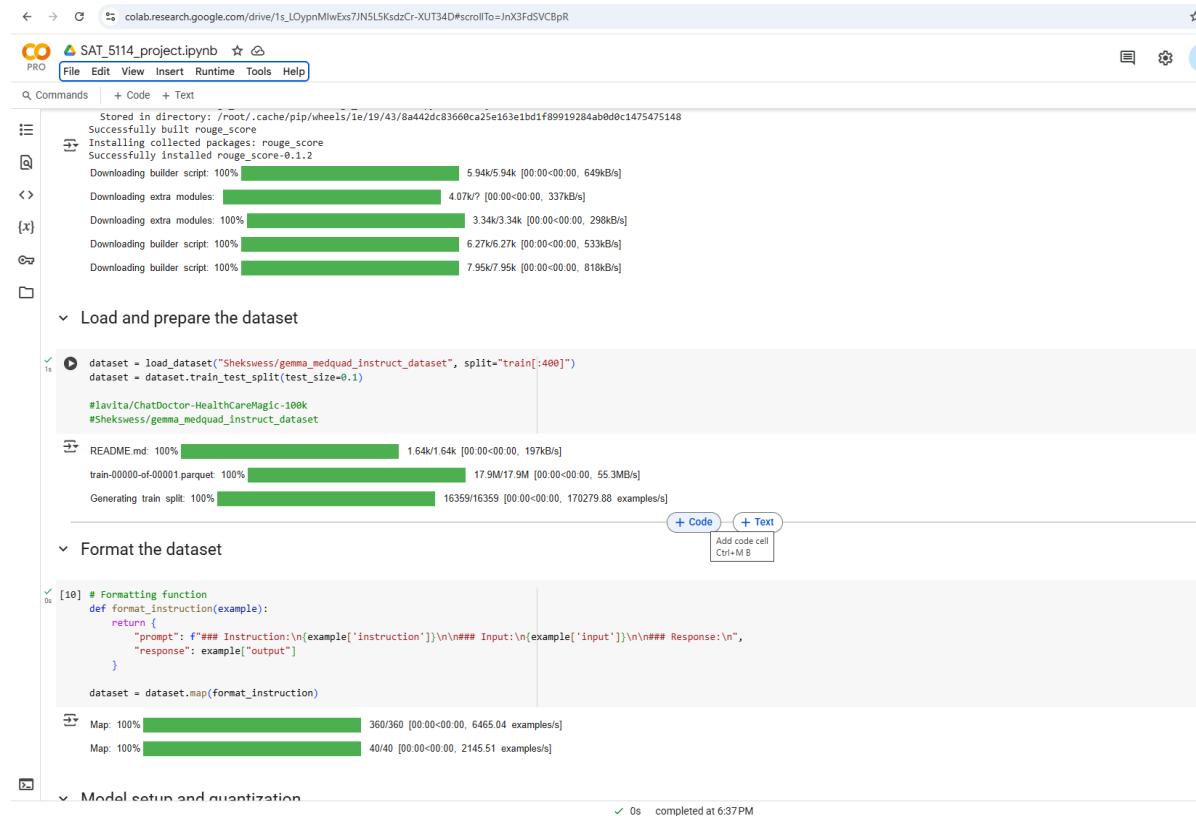
- Optimize workflows

# Code Demo (Dataset and Model)

The dataset is from Hugging Face
(https://huggingface.co/datasets/Shekswess/gemma_medquad_instruct_dataset)
Model: Llama 3 8b from Unsloth

# Code Demo- Cont'd (LoRA and Performance Metrics)

Bleu : 0.45

BertF1 score = 0.90

RougeL = 0.66

# Code Demo-Cont'd (Text Generation in GUI)

Test the model by inputting a query to generate a response
GUI using Gradio

# Conclusion

- This model performed well regarding generating responses to queries; however, it was trained on a significantly small portion of the dataset to conserve computer power which affected its performance

- It can be upgraded to accept queries in different formats, like audio and video queries

# References

1.  Meskó B

    The Impact of Multimodal Large Language Models on Health Care's Future

    J Med Internet Res 2023;25:e52865

    URL: https://www.jmir.org/2023/1/e52865

    DOI: 10.2196/52865


    2. Applying Object Detection and Large Language Model to Establish a Smart Telemedicine Diagnosis System with Chatbot: A Case Study of Pressure Injuries Diagnosis System

       Chun-Chia Chen, Chia-Jung Wei, Tsung-Yu Tseng, Ming-Chuan Chiu, and Chi-Chang Chang

       Telemedicine and e-Health 2024 30:6, e1705-e1712. https://www.liebertpub.com/action/showCitFormats?doi=10.1089%2Ftmj.2023.0715

3.  Denecke, K., May, R. & Rivera-Romero, O. Transformer Models in Healthcare: A Survey and Thematic Analysis of Potentials, Shortcomings and Risks. *J Med Syst* 48,
    (2024). https://doi.org/10.1007/s10916-024-02043-5


4.  Mohamad-Hani, T., Jamal, A., Khalid, A., Fadi, A., Ibraheem, A., Malki, K. H., . . . Al-Eyadhy, A. (2024). Transforming virtual healthcare: The potentials of
    telemedicine. *Cureus, 16*(5) doi: https://doi.org/10.7759/cureus.61377