# Local LLM(10 days project)

Industrial Attachment IMH

# Objectives

To study the possible to implement the local LLM in the control environment without internet access

To demonstrate the application implement on the local LLM

Understand the performance of the application with the local LLM

# Local Large Language Model(LLM)

Challenge

Computational resources: LLMs can be very demanding on your computer's resources, especially the CPU and GPU. You may need a high-end computer to run them smoothly.

Storage:  LLMs can also be quite large, so you will need a significant amount of storage space on your computer to download and run them.

Maintenance:  You will be responsible for keeping the LLM software up-to-date and troubleshooting any problems that arise.

# Local Large Language Model(LLM)

Method

GGUF (GPT-Generated Unified Format) is a new framework designed to make it easier to run large language models (LLMs) on your computer. It primarily uses CPU resources, but can also leverage GPUs for specific parts when needed. This makes it well-suited for devices with weaker CPUs, like those found in Apple products. GGUF also compresses the data used by the LLM, making it more efficient to store and run. Overall, GGUF is a more flexible and user-friendly way to run LLMs locally.

# Local Large Language Model(LLM)

Open Source LLM

**Mistral** AI is an artificial intelligence startup that makes large language models (LLMs). Based in Paris, France, and founded by former researchers at Google DeepMind and Meta, Mistral is known for its transparent, portable, customizable and cost-effective models that require fewer computational resources than other popular LLMs.

**Gemma** is a family of open-weights Large Language Model (LLM) by Google DeepMind, based on Gemini research and technology.

**Orca** LLM, a language model developed by Microsoft, seeks to address the limitations of ChatGPT by introducing a logic-based framework that simulates human-like reasoning.

# Local Large Language Model(LLM)

Implementation: Using Langchain - llamacpp

Universal Compatibility: Llama.cpp's design as a CPU-first C++ library means less complexity and seamless integration into other programming environments. This broad compatibility accelerated its adoption across various platforms.

Comprehensive Feature Integration: Acting as a repository for critical low-level features, Llama.cpp mirrors LangChain's approach for high-level capabilities, streamlining the development process albeit with potential future scalability challenges.

Focused Optimization: Llama.cpp focuses on a single model architecture, enabling precise and effective improvements. Its commitment to Llama models through formats like GGML and GGUF has led to substantial efficiency gains.

# Local Large Language Model(LLM)

Implementation

```
model_path = "./models/llm/mistral-7b-instruct-v0.1.Q6_K.gguf"
llm = LlamaCpp(
    model_path = model_path,
    n_gpu_layers=n_gpu_layers,
    n_batch=n_batch,
    n_ctx=8000,
    max_tokens = 1024,
    f16_kv=True,  # MUST set to True, otherwise you will run into problem after a couple of calls
    callback_manager=callback_manager,
    verbose=False,
)
```

## Test Summary of a synthetic medical record with basic prompt



```
prompt = f"""
  Summarise the patient medical
record in the following text.

  Below is the text to be
summarised:   {text}
  Summarise text:

"""
```

# Local Large Language Model(LLM)

Result with CPU only

```
start_time = datetime.now()
llm.invoke(prompt)
time_diff= (datetime.now()  - start_time).total_seconds()
print("Execution time of program is: ", time_diff, "s")
```

Jane Doe is a 44-year-old female with a history of recurrent headaches associated with visual disturbances. She has a past medical history of migraines, hypothyroidism, and seasonal allergies. Her medications include Levothyroxine for hypothyroidism and over-the-counter ibuprofen for headache relief. She denies any recent changes in thyroid function or symptoms of hypothyroidism. Her family history includes a mother with migraines and fathers with hypertension and diabetes. Jane works as an office manager and reports occasional alcohol consumption. During her physical examination, she is alert and oriented and has no focal neurological deficits. She presents with chronic migraines with an aura. Her plan of care includes initiating prophylactic migraine therapy with amitriptyline, monitoring for improvement in headache frequency and severity over the next 4-6 weeks, scheduling a follow-up appointment in 4 weeks to assess response to treatment and consider titration of amitriptyline dosage if needed. She will also continue taking Levothyroxine for her hypothyroidism and be educated on lifestyle modifications including stress management techniques, regular exercise, adequate hydration, and avoidance of known migraine triggers.Execution time of program is:  34.14207 s
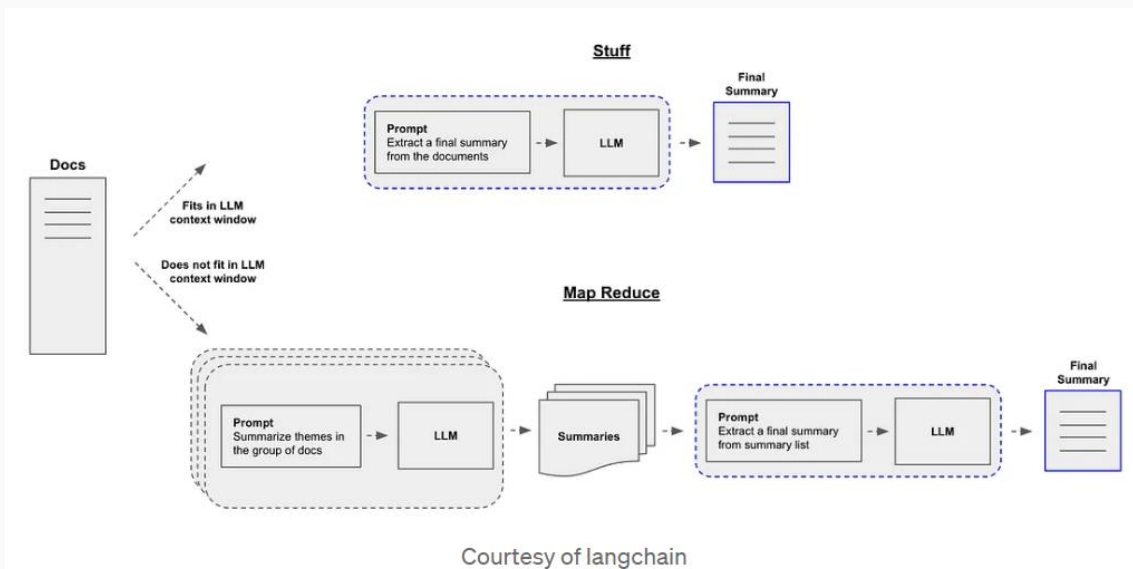
# Local Large Language Model(LLM)

## Result with CPU+GPU

```python
start_time = datetime.now()
llm.invoke(prompt)
time_diff= (datetime.now() - start_time).total_seconds()
print("Execution time of program is: ", time_diff, "s")
```

The patient is a 44-year-old female with a history of recurrent headaches associated with visual disturbances. The patient has been diagnosed with hypothyroidism, seasonal allergies and has a past medical history of migraines. The patient is currently taking Levothyroxine 50 mcg daily for her hypothyroidism and over-the-counter ibuprofen for headache relief as needed. The patient has no known drug allergies. She presents to the clinic with the chief complaint of recurrent headaches associated with visual disturbances that have been progressively worsening in intensity and frequency. The patient is alert and oriented, with no focal neurological deficits or changes in consciousness. Physical examination reveals normal optic discs and no evidence of papilledema. The plan for the patient includes initiating prophylactic migraine therapy with amitriptyline 25 mg orally at bedtime, educating the patient on potential side effects, monitoring the response to treatment, and considering neuro imaging if necessary. Additionally, the patient should continue taking Levothyroxine 50 mcg daily for her hypothyroidism, have annual or clinically indicated thyroid function tests, and make lifestyle modifications including stress management techniques, regular exercise, adequate hydration, and avoiding known migraine triggers.Execution time of program is:  6.342643 s

# Local LLM Application - Text Summarization

## Summarization Technique



Courtesy of langchain

# Local LLM Application - Text Summarization

## Summarization Application (Map Reduce)



1.Load PDF

2. Map Reduce summarization

3.Display Summary

# Text Summarization -Mapreduce

# Local LLM Application - Text Summarization

## Summarization Application (Refine)



1.Load PDF

2. Refine summarization

3.Display Summary

source:https://python.langchain.com/docs/modules/chains/document/refine

# Text Summarization - Refine

# LLM Performance

LLM 7b/13b

Time for summarization time - time taken

Accuracy of summary - BertScore

# LLM Performance

LLM 7b/13b

Time for summarization time - time taken

With the GPU the long text summarization time taken was reasonable

Testing on different text sizes(20K,50K tokens)

# LLM Performance -BertScore

LLM 7b/13b

Accuracy of summary - BertScore

BertScore is a method used to measure the quality of text summarization. This method measures how similar the text summary is to the original text.

Challenge: not able to find good annotated summarized text to compare the generated summary.

# LLM Performance- Custom LLM prompt

Use the LLM prompt to evaluate the summary content(Not implemented yet)

evaluate the text summarization based on the following:

1.Coverage of details from the original text.

2.Factual alignment between the original text and summary.

Give each result in the grading 1-10, where 1 is the lowest and 10 is the highest.

original text:

# LLM Performance- Custom LLM prompt

Use the LLM prompt to evaluate the summary respond(summary from mistral 7b)

Based on the evaluation criteria provided:

1. **Summary text's coverage of details from the original text**: 4/10
   - The summary text covers the general idea of stress's impact on immunity and mentions some key points such as acute versus chronic stress effects, the role of hormones, and implications for vaccination responses. However, it lacks specific details about mechanisms like the HPA axis, specific cytokine interactions, and detailed immune cell responses.
2. **Factual alignment between the original text and summary**: 3/10
   - There are significant discrepancies between the original text and the summary. The summary oversimplifies and generalizes complex interactions and mechanisms described in the original text. It misses crucial details such as specific hormone names, immune cell types, and their functions in stress responses. Moreover, it introduces new information about ChatGPT in medicine, which is unrelated to the original topic.

Overall, while the summary attempts to capture the general concept of stress's effects on immunity, it falls short in terms of detail coverage and factual alignment with the original text.

# LLM Performance- Custom LLM prompt

Use the LLM prompt to evaluate the summary respond(summary from mistral 7b)

Based on the criteria provided:

1. **Summary text's coverage of details from the original text**: 3/10
   - The summary does not adequately cover the details from the original text. It focuses solely on hypertension, which is not the primary focus of the original text on stress and immune response.
2. **Factual alignment between the original text and summary**: 2/10
   - There is minimal alignment between the factual details of the original text and the summary. The summary completely diverges from the original text's focus on stress, immunity, and psychological factors affecting health.

Overall, the summary fails to reflect the main themes, content, and factual accuracy of the original text.

# LLM Performance- Custom LLM prompt

Use the LLM prompt to evaluate the summary respond(summary from Gemma 7b)

Based on the evaluation criteria provided, let's assess the summary text's performance:

1. **Summary text's coverage of details from the original text**: 6/10
   - The summary covers various aspects discussed in the original text, such as stress effects on immunity, cytokine influence, and psychological stress impacts. However, it misses some detailed explanations and examples provided in the original, resulting in a slightly reduced score.
2. **Factual alignment between the original text and summary**: 7/10
   - The summary generally aligns well with the factual content of the original text. It accurately captures the main concepts like stress impact on immune cells, hormonal responses, and immune system modulation. However, some nuances and specific mechanisms described in the original are not fully reflected in the summary, affecting the completeness of alignment.

Overall, while the summary provides a reasonable overview of the original text's main points, it could benefit from including more specific details and examples to improve coverage and alignment.

# LLM Performance- Custom LLM prompt

Use the LLM prompt to evaluate the summary respond(summary from Orca 7b)

Based on the evaluation criteria provided, here are the ratings for the text summarization:

1. **Coverage of details from the original text**: 7/10
   - The summary covers a significant amount of detail from the original text, including discussions on stress effects on immunity, pathways involved, and specific immune responses. However, some details on specific hormones and immune cell types could have been more explicitly included.
2. **Factual alignment between the original text and summary**: 8/10
   - The summary aligns well with the factual content of the original text. It accurately reflects the main points about stress's impact on immune function, including mechanisms and outcomes discussed. However, there were a few instances where specific details or pathways were not entirely mirrored in the summary.

Overall, the summarization effectively captures the essence of the original text on stress and immunity, providing a comprehensive overview while maintaining factual accuracy.