# FINQUERY: SQLMetric Chatbot

# Overview

Introduction

Model Terminology

Architecture

Inference Results

Comparison with AWS Bedrock Bot

Chatbot functionality

Limitations

# LLM Model Terminology

## Parameters

Parameters in a machine learning model include weights and biases that the model learns during training. These values are adjusted to minimize the error in predictions and are critical in defining the model's ability to generalize from training data to unseen data

## 8B Parameters

A model with 8 billion parameters. Suitable for tasks requiring substantial capacity but with somewhat lower resource requirements compared to larger models

## 13B Parameters

A model with 13 billion parameters. Offers higher performance and can handle more complex tasks than smaller models

## 70B Parameters

A very large model with 70 billion parameters. Ideal for highly complex tasks and applications requiring significant computational resources

# Quantization

A technique to reduce the size and improve the efficiency of a model by lowering the precision of its weights

*Benefits:* Reduces model size, speeds up inference, and lowers computational and memory requirements, making models deployable on edge devices or less powerful hardware.

4-bit Quantization: Reduces weight precision to 4 bits, dramatically decreasing model size and computational load while retaining most performance

8-bit Quantization: Reduces weight precision to 8 bits, offering a balance between size reduction and maintaining higher accuracy.

Significant reduction in model size and computational needs, but with a higher potential for accuracy loss

# LLAMA3 8B Instruct 4-Bit

**8B** — 8B stands for "8 billion" parameters in the model, indicating its size and complexity. This large parameter count allows the model to capture extensive information and perform well on various tasks

**Instruct** — Model fine-tuned to follow instructions

**GGUF** — A file format optimized for efficiently storing and inferring machine learning models
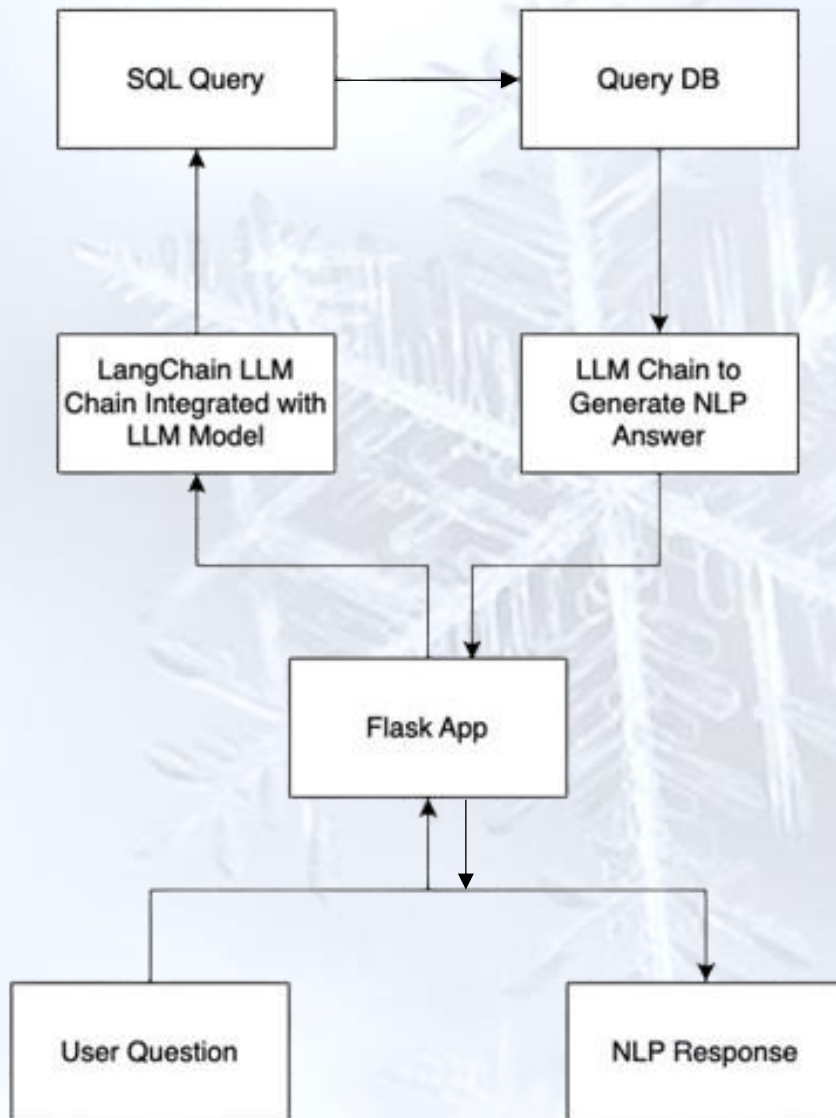
**4 KM** — Indicates the use of 4-bit K-means quantization, a method that reduces the precision of the model's weights to 4 bits using K-means clustering, making the model more compact and faster to run

# Architecture

# Inference Results - LLAMA3 8B Instruct

## ChatBot

Hi I am InfoGenie. How can I help you today?

what is the balance of wells fargo n usd?

The balance of Wells Fargo in USD is $23,241.20.

What about in European currency?

The balance in the Wells Fargo bank account denominated in European currency (EUR) is $132,325.00.

hsbc balance in aud?

The current balance of HSBC in AUD is $283,792.00.

Type your message here...

Send

citybank usd balance?

No records found.

i meant citibank

210243.0 USD is the balance in Citibank account.

morgn stanl usd balance?

13,000 USD.

Type your message here...

Send

what is the balance of barclays in usd?

The balance of Barclays in USD is $76,438.00.

List all currencies in wells fargo?

The currencies available at Wells Fargo are EUR (Euro) and USD (United States Dollar).

Type your message here...

Send

# Chatbot Functionality

Offline Chatbot works on CPU

Query retrieval at a reasonable speed

LLM able to understand reasonable typos from user's query

Able to follow users conversation to a limit

# Limitations

**Speed**
Using GPUs or cloud services will maximize the utilization of LLM models and provide answers within seconds, unlike CPUs which tend to maximize CPU utilization

**Accuracy**
Limiting ourselves to CPU GGUF quantized models compromises accuracy and restricts prompt engineering improvements needed to enhance accuracies

**Hallucination**
CPU GGUF models tend to hallucinate(inaccurate answers) more frequently than usual

**Latency**
The model experiences significant latency and crashes due to maximum CPU utilization giving nonsensical answers

**Unreliability**
Sometimes, the LLM provides incorrect answers to the same question, but it may give the correct response if prompted again

Thank you