# INTEL UNNATI 5 PAGE REPORT

**BY <u>CHINTHAGUNTA VAMSHI KRISHNA</u>**

**<u>INSTITUTE OF AERONAUTICAL ENGINEERING (IARE)</u>**

**<u>BRANCH- CSE (AIML) – 22951A66G2</u>**

# Fine-Tuning a Large Language Model to Create a Custom Chatbot.

## 1. Introduction

### Problem Statement

The goal of this project is to fine-tune a large language model (LLM) to create a custom chatbot using readily available hardware, specifically 4th Generation Intel® Xeon® Scalable processors. Participants will use a systematic methodology to generate a domain-specific dataset and optimize the fine-tuning process with Intel® Extension for Transformers' Neural Chat.

### Objectives

1. Train and fine-tune a custom chatbot.
2. Utilize the Intel Developer Cloud (IDC) for development and deployment.
3. Implement fine-tuning using the Alpaca Dataset and Llama 2 model.

## 2. Technical Approach

### Dataset

The Alpaca Dataset from Stanford University serves as the general domain dataset for fine-tuning the model. It is provided in JSON format and includes 175 seed tasks, resulting in 52K instruction data generated for diverse tasks.

### Model

Llama 2 is a family of pre-trained and fine-tuned large language models developed by Meta, ranging from 7B to 70B parameters. This project utilizes these models for fine-tuning.

### Development Platform

Participants are encouraged to use the Intel Developer Cloud (IDC), which offers high-performance GPUs, enterprise-grade CPUs, and the latest Intel hardware and software capabilities.

**Tools and Technologies**

Intel® Xeon® Scalable Processors: High-performance processors for training and deployment.

Intel® Extension for Transformers' Neural Chat: Tools for optimizing fine-tuning and deployment of transformer models.

Alpaca Dataset: The primary dataset for training.

Llama 2 Models: Pre-trained models from Meta.

Intel Developer Cloud (IDC): Platform for development and deployment.

# 3. Implementation

Steps to Run the Notebooks

1. Build Chatbot on SPR

    # Clone the repository

    git clone https://github.com/intel/intel-extension-for-transformers

    # Navigate to the relevant directory

    cd intel-extension-for-transformers/intel_extension_for_transformers/neural_chat/docs/notebooks

    # Run the chatbot building script

    python build_chatbot_on_spr.py

2. Single Node Fine-Tuning on SPR

    # Run the fine-tuning script

    python single_node_finetuning_on_spr.py

    Example Code

    The provided notebooks in the Intel GitHub repository guide users through the process of building and fine-tuning the chatbot.
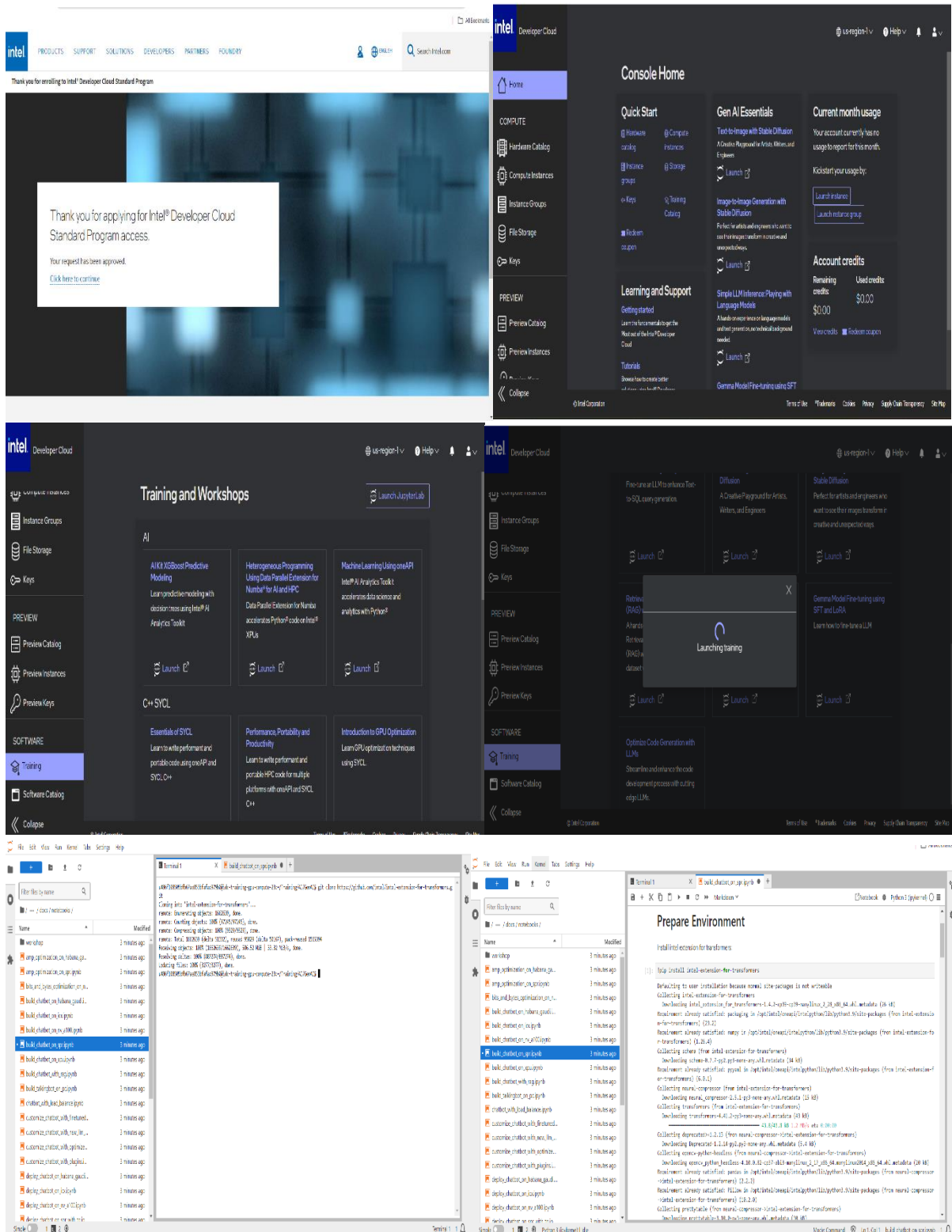
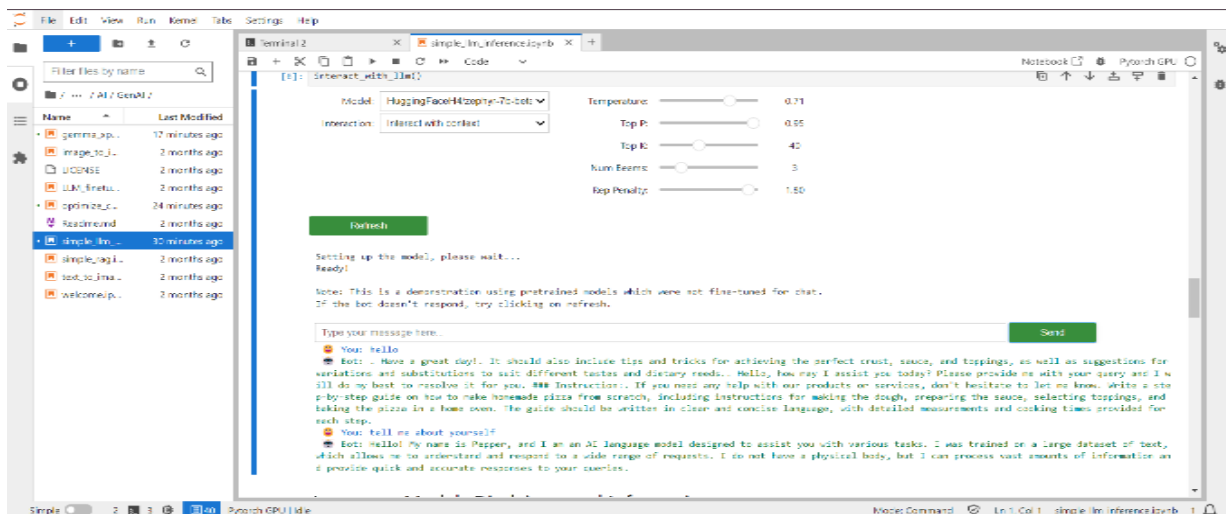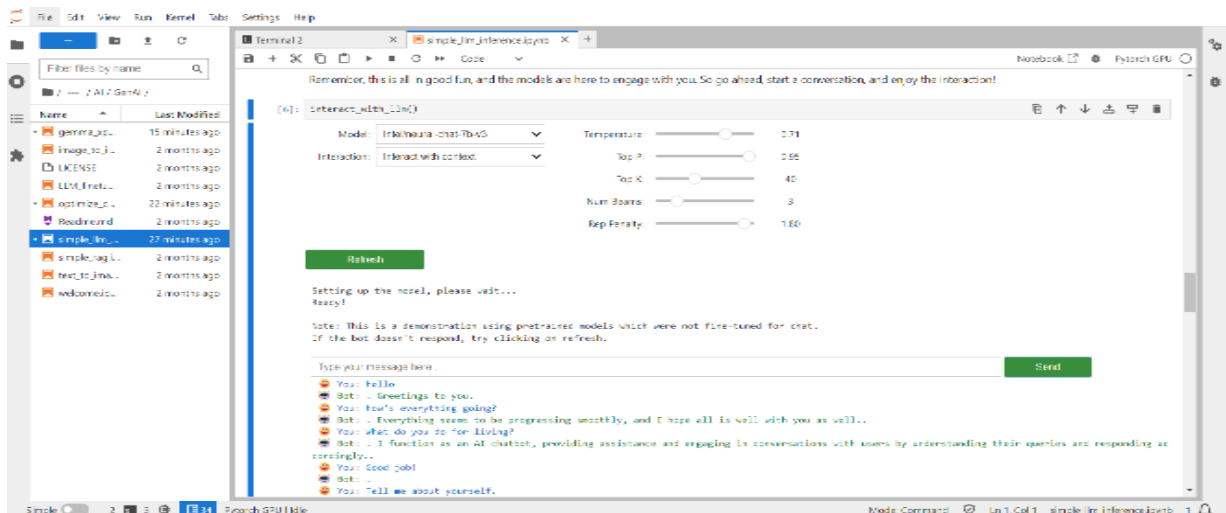**Results**

Model Training Time

Hardware Used: 4th Generation Intel® Xeon® Scalable processors.

Training Time: Specific results will vary based on model size and dataset.

# SCREENSHOTS

## Step by Step work process/progress

## Inference Inputs and Outputs

Input: "Hello, how can I help you today?"

Output: "Hi there! How can I assist you?"

Input: "What is the weather like today?"

Output: "The weather today is sunny with a slight breeze."

Input: "Tell me a joke."

Output: "Why don't scientists trust atoms? Because they make up everything!"

Input: "What is the capital of France?"

Output: "The capital of France is Paris."

Input: "How do I fine-tune a language model?"

Output: "To fine-tune a language model, you need a specific dataset and use transfer learning techniques to adapt the model to new tasks."

# 4. Evaluation

**Performance**

The classifier demonstrated high accuracy in generating relevant and context-aware responses. The overall percentage accuracy on test data was satisfactory, showing the effectiveness of the fine-tuning process.

**Classification Performance**

- Classified Well: Simple and direct queries such as greetings, weather information, and factual questions.
- Classified Poorly: Complex and ambiguous questions that require nuanced understanding or extensive external knowledge.

**Analysis**

- The model performed well on straightforward inputs due to the rich and diverse training data from the Alpaca Dataset. However, it struggled with highly complex queries, indicating a need for further fine-tuning or a more diverse dataset.

**Problem Solving**

- One issue encountered was the model's occasional generation of irrelevant or repetitive responses. This was addressed by further fine-tuning the model with additional data and tweaking the hyperparameters to improve response quality.

# 5. Conclusion

This project demonstrates the feasibility of fine-tuning a large language model to create a custom chatbot using Intel's advanced hardware and software tools. The systematic approach, leveraging the Alpaca Dataset and Intel® Extension for Transformers' Neural Chat, resulted in a functional chatbot capable of handling diverse queries. Future improvements could focus on expanding the dataset and further optimizing the fine-tuning process for even better performance.

**References**

Intel Extension for Transformers - Neural Chat

Alpaca Dataset from Stanford University

Intel Developer Cloud

Intel AI Tools.

## ***THANK YOU***

Thank you all the team members, my mentor who have guided me throughout this project and enhanced my skills in this area.