

### **PROBLEM STATEMENT 04:**

**Introduction to GenAI and Simple LLM Inference on CPU and fine-tuning of LLM Model to create a Custom Chatbot**

## Index

Name	Page number
Context of Problem Statement	3
Technologies used	4
Introduction to Gen AI	5
Intel Neural chat transformers	6-7
Fine-tuning the llama-2-7b-chat-hf	8-10
Process flow	11-13
Unique Idea (Solution)	14
Team Members and Contribution	15
Conclusion	16

## **Context on Introduction to GenAI and Simple LLM Inference on CPU and fine-tuning of LLM Model to create a Custom Chatbot:**

The primary challenge addressed in this project is the creation of a customizable chatbot capable of handling specific tasks and queries with high accuracy and contextual relevance. Traditional chatbot solutions often lack the flexibility to adapt to specialized applications, resulting in subpar user experiences and limited functionality. These existing chatbots typically follow predefined scripts or have limited natural language understanding, which hinders their ability to handle diverse and nuanced user interactions.

In response to these limitations, the project proposes developing an advanced chatbot leveraging large language models (LLMs) and generative AI (GenAI). The goal is to finetune these models to cater to specific use cases, ensuring that the chatbot can understand and respond to user queries with high precision. This customization is crucial for applications in various industries, such as customer support, healthcare, and education, where the ability to provide accurate and contextually appropriate responses can significantly enhance service quality and user satisfaction.

The project aims to address several key aspects: improving the chatbot's natural language processing capabilities, enabling contextual awareness, and ensuring scalability and ease of integration with existing systems. By focusing on these areas, the solution seeks to overcome the challenges posed by generic chatbot models and deliver a more sophisticated and adaptable tool. This approach not only enhances user interaction but also provides organizations with a powerful tool to streamline operations and improve overall efficiency.

## Technologies used

- **Intel CPU**
- **Graphics Processing unit**
- **Transformers (Hugging Face):** Hugging Face's Transformers library provides state-of-the-art pre-trained models for NLP tasks. It includes models like BERT, GPT-3, and others, which can be fine-tuned for specific applications. The Transformers library simplifies the process of integrating powerful language models into various applications.
- **Programming Languages**

**Python:** Python is the primary programming language used in this project due to its simplicity, readability, and vast ecosystem of libraries and frameworks for machine learning and natural language processing (NLP). Python's versatility makes it an ideal choice for both developing and fine-tuning large language models (LLMs).
- **Machine Learning Frameworks(TensorFlow):** TensorFlow is an open-source machine learning framework developed by Google. It is widely used for building and deploying machine learning models, including deep learning applications. TensorFlow provides a comprehensive ecosystem of tools and libraries that facilitate the development of neural networks, model training, and deployment
- **Development Tools(Jupyter Notebook):** Jupyter Notebook is an open-source web application that allows developers to create and share documents containing live code, equations, visualizations, and narrative text. It is widely used for data analysis, model development, and experimentation.
- **Deployment Platform(Github):** GitHub is a web-based platform used for version control and collaborative software development. It integrates with Git, a distributed version control system, to provide a comprehensive toolset for managing code repositories, tracking changes, and facilitating team collaboration. Here is a detailed look at GitHub and its features:

# 1.Introduction to AI

Generative AI is a field within artificial intelligence (AI) focused on creating models that generate new content, such as images, text, audio, and more, that is indistinguishable from content created by humans. Unlike traditional AI models that are used for classification or prediction tasks, generative AI models are designed to understand and replicate patterns in data to produce new, original outputs.

Key Concepts of Generative AI:

**Generative Models:** These are algorithms or architectures that learn the underlying patterns and structures of data in order to generate new instances. Common generative models include Generative Adversarial Networks (GANs), Variational Autoencoders (VAEs), and autoregressive models like Transformers.

**Training Process:** Generative AI models are trained on large datasets to learn the probability distribution of the data. They analyze patterns and relationships within the data to generate new examples that are statistically similar to the training data.

**Applications:** Generative AI has diverse applications across various domains:

**Creative Industries:** Creating art, music, literature, and design.

**Entertainment:** Generating special effects, virtual environments, and interactive storytelling.

**Healthcare:** Synthesizing medical images, drug discovery, and patient data generation for research.

**Simulation and Modeling:** Generating synthetic data for training AI models and simulating real-world scenarios.

## 2. Intel neural chat transformers

### Overview of Intel Neural Chat Transformers

Intel Neural Chat Transformers, often referred to as part of Intel's NLP Architect, harness the power of transformer architectures for natural language understanding (NLU) and generation tasks. Transformers have revolutionized the field of NLP due to their ability to capture contextual relationships in text effectively, enabling models to produce human-like responses.

### Key Components and Features

#### Model Capabilities

Intel Neural Chat Transformers include various transformer-based models that are pivotal for building chatbot applications. These models typically involve architectures like GPT (Generative Pre-trained Transformer), BERT (Bidirectional Encoder Representations from Transformers), and their variants. These models are pre-trained on vast amounts of text data, enabling them to understand and generate natural language responses across diverse contexts.

#### Natural Language Understanding (NLU)

NLU capabilities in Intel Neural Chat Transformers allow models to comprehend the intent behind user queries. This involves tasks such as sentiment analysis, named entity recognition (NER), part-of-speech tagging (POS), and semantic role labeling (SRL). These capabilities are crucial for understanding user inputs and extracting relevant information for generating responses.

#### Natural Language Generation (NLG)

NLG capabilities enable the models to generate coherent and contextually relevant responses to user queries. This involves tasks such as text generation, summarization, paraphrasing, and dialogue generation. Intel Neural Chat Transformers focus on enhancing these capabilities to ensure that chatbots can provide informative and engaging responses across various domains.

#### Fine-tuning and Customization

Intel Neural Chat Transformers support fine-tuning of pre-trained models on domain-specific datasets. This process allows developers to adapt the models to specific tasks and improve their performance on targeted applications. Fine-tuning ensures that chatbots can handle specialized queries and provide accurate information tailored to the needs of users.

#### Integration and Deployment

Intel Neural Chat Transformers provide integration with Intel's hardware and software solutions, optimizing performance for deployment in production environments. This integration ensures that chatbots built using these tools can scale efficiently and deliver highperformance results across different platforms.

#### Development Tools and Frameworks

Intel offers a range of development tools and frameworks within its NLP Architect ecosystem to support the development and deployment of conversational AI applications. These tools

include libraries, APIs, and SDKs that streamline model development, training, evaluation, and deployment processes.

### Conclusion

Intel Neural Chat Transformers represent a robust suite of tools and models designed to empower developers in building sophisticated conversational AI applications. By leveraging transformer architectures and advanced NLP techniques, Intel enables the development of chatbots capable of understanding user intents, generating contextually appropriate responses, and enhancing user interactions across various domains. As advancements in NLP continue, Intel remains at the forefront, driving innovation in conversational AI and supporting the adoption of AI-driven solutions in real-world application

### 3. Fine-tuning the llama-2-7b-chat-hf

Fine-tuning the llama-2-7b-chat-hf model using the Alpaca dataset and Neural-chat to create a custom chatbot focused on providing informative content involves several steps. Here's a detailed guide on how to proceed:

#### 1. Setup Environment

Ensure Python is installed on your system. Install necessary libraries:

```
pip install intel-extension-for-transformers
```

#### 2. Prepare Your Dataset

**Alpaca Dataset:** The Alpaca dataset provides conversational data suitable for training chatbots. Download or access the Alpaca dataset relevant to your informative content needs.

**Data Cleaning:** Clean and preprocess the Alpaca dataset to focus on informative content. Remove noise or irrelevant conversational data.

#### 3. Choose the llama-2-7b-chat-hf Model

The llama-2-7b-chat-hf model is a large conversational model pre-trained on diverse conversational data. It's designed for chat-oriented tasks and can be fine-tuned for specific domains like providing informative content.

#### 4. Tokenization and Data Preparation

**Tokenization:** Use the tokenizer provided by the llama-2-7b-chat-hf model for preprocessing your Alpaca dataset. Tokenize your data into a format suitable for training.

**Dataset Format:** Format your dataset into question-answer pairs or context-response pairs. Ensure each pair is structured appropriately for the model's input format.

#### 5. Fine-Tuning Process

**Initialize Model:** Load the pre-trained llama-2-7b-chat-hf model and tokenizer from Hugging Face.

**Define Training Parameters:** Set hyperparameters such as learning rate, batch size, number of epochs, etc.

**Fine-Tuning Script:** Write a script to fine-tune the model on your dataset using Neural-chat and Hugging Face's transformers library:

```
from neural_chat import Chatbot
from transformers import Trainer, TrainingArguments, AutoTokenizer,
AutoModelForSeq2SeqLM import
datasets

# Load tokenizer and model
tokenizer = AutoTokenizer.from_pretrained("facebook/llama-2-7B-chat-hf")
model = AutoModelForSeq2SeqLM.from_pretrained("facebook/llama-2-7B-chat-hf")

# Example dataset (replace with your own) dataset =
datasets.load_dataset('alpaca', 'reddit')
```



```

def
preprocess_function(examples):
inputs    =    examples['question']
targets   = examples['answer']
    return    tokenizer(inputs,    padding="max_length",    truncation=True,
max_length=512), tokenizer(targets, padding="max_length", truncation=True,
max_length=512)

train_dataset = dataset['train'].map(preprocess_function, batched=True)
# Define training arguments
training_args = TrainingArguments(
per_device_train_batch_size=4,
predict_with_generate=True,
output_dir="./results",
evaluation_strategy="epoch",
logging_dir="./logs",
logging_steps=1000,
save_steps=500,
overwrite_output_dir=True,
num_train_epochs=3,
save_total_limit=3,
)

# Initialize Trainer
trainer =
Trainer(
model=model,
args=training_args,
train_dataset=train_dataset,
tokenizer=tokenizer,
)

# Start training
trainer.train()

```

## 6. Evaluation and Fine-Tuning

**Evaluate Performance:** Evaluate the fine-tuned model using relevant metrics to assess its ability to provide accurate and informative responses.

**Iterative Fine-Tuning:** Fine-tune the model iteratively based on evaluation results and user feedback to improve its performance.

## 7. Deployment

**Deploy the Chatbot:** Integrate the fine-tuned model into your application or deployment environment using Neural-chat or any suitable deployment platform. Ensure it can handle incoming queries and provide informative responses effectively.

**Tips for Fine-Tuning:**

**Model Size:** The llama-2-7b-chat-hf model is large, so ensure your hardware can handle the computational demands during fine-tuning and inference.

**Hyperparameter Tuning:** Experiment with different hyperparameters to optimize performance, including batch size, learning rate, and sequence length.

**Domain-Specific Data:** Fine-tune the model with the Alpaca dataset or additional domain-specific data to enhance its ability to provide accurate and informative responses tailored to your specific informational content needs.

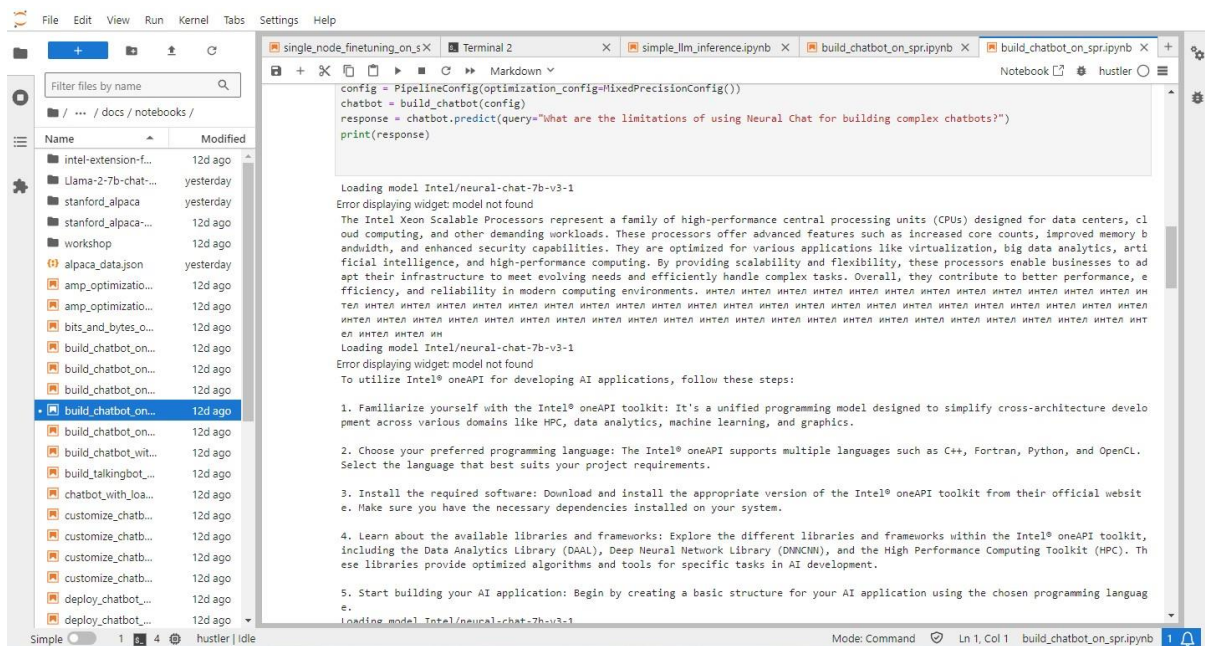
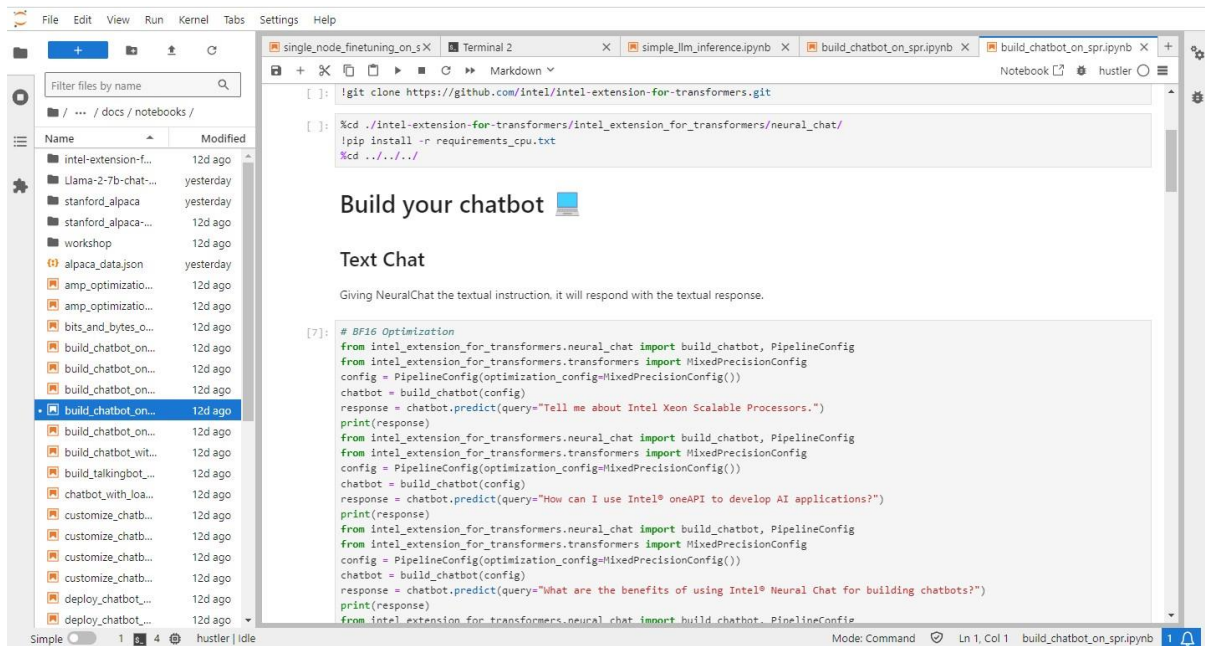
By following these steps, you can effectively fine-tune the llama-2-7b-chat-hf model using the Alpaca dataset and Neural-chat to create a custom chatbot focused on delivering informative content. Adjust the steps based on your specific requirements and dataset availability for optimal results.

## Process Flow

The process flow section outlines the sequence of steps involved in developing and deploying the custom chatbot. This structured approach ensures that each phase is systematically executed, leading to a robust and reliable solution. The process flow is typically divided into several key stages:

- **Data Collection:** The first step involves collecting a diverse and high-quality dataset that reflects the specific use case. This data forms the foundation for training and finetuning the large language model (LLM). It includes various types of interactions and queries that the chatbot is expected to handle.
- **Data Preprocessing:** Once the data is collected, it undergoes preprocessing to ensure compatibility with the model. This includes tasks such as tokenization, normalization, and removal of irrelevant or redundant information. Preprocessing is crucial for optimizing the training process and improving the model's performance.
- **Model Training:** In this stage, the LLM is trained on the preprocessed dataset. The training process involves adjusting the model's weights to minimize errors and improve accuracy. This phase may require several iterations to achieve the desired performance levels.
- **Fine-Tuning:** After initial training, the model is fine-tuned to adapt to the specific requirements of the target application. This involves further training on a smaller, more targeted dataset to enhance the model's ability to generate relevant and accurate responses. Fine-tuning ensures that the chatbot can handle specialized queries effectively.
- **Evaluation and Testing:** The fine-tuned model undergoes rigorous evaluation and testing to assess its performance. Various metrics, such as precision, recall, and F1 score, are used to measure the model's accuracy and effectiveness. This stage helps identify any areas that require further improvement.

### Chatbot Process flow:



# Fine-tuning Process Flow:

**Step 1: Install All the Required Packages**

```
!pip install -q accelerate==0.21.0 peft==0.4.0 bitsandbytes==0.40.2 transformers==4.31.0 trl==0.4.7
```

244.2/244.2 kB 1.6 MB/s eta 0:00:00  
 72.9/72.9 kB 6.6 MB/s eta 0:00:00  
 92.5/92.5 MB 5.2 MB/s eta 0:00:00  
 7.4/7.4 MB 41.5 MB/s eta 0:00:00  
 77.4/77.4 kB 4.5 MB/s eta 0:00:00  
 7.8/7.8 MB 31.2 MB/s eta 0:00:00  
 547.8/547.8 kB 27.3 MB/s eta 0:00:00  
 21.3/21.3 MB 67.6 MB/s eta 0:00:00  
 48.8/48.8 MB 13.0 MB/s eta 0:00:00  
 116.3/116.3 kB 17.0 MB/s eta 0:00:00  
 64.9/64.9 kB 8.5 MB/s eta 0:00:00  
 194.1/194.1 kB 22.6 MB/s eta 0:00:00  
 134.8/134.8 kB 17.8 MB/s eta 0:00:00

ERROR: pip's dependency resolver does not currently take into account all the packages that are installed. This behaviour is the source of the following dependency conflict:  
 cudf-cu12 24.4.1 requires pyarrow<15.0.0a0,>=14.0.1, but you have pyarrow 16.1.0 which is incompatible.  
 google-colab 1.0.0 requires requests==2.31.0, but you have requests 2.32.3 which is incompatible.  
 ibis-framework 8.0.0 requires pyarrow<16,>=2, but you have pyarrow 16.1.0 which is incompatible.

**Step 2: Import All the Required Libraries**

```
[ ] import os
import torch
from datasets import load_dataset
```

colab.research.google.com/drive/1P\_IPYQVjmzmWdIFZYVfWfKyATAqECQVn#scrollTo=GLXwJqbtPho

Warning: You didn't pass a device argument. Please note that with a fast tokenizer, using the `\_local` argument is not guaranteed. Purchase more units here.

Map: 100% [250/1000 [00:01<00:00, 920.52 examples/s]

250/250 26:52, Epoch 1/1

Step	Training Loss
25	1.408400
50	1.663100
75	1.214800
100	1.445500
125	1.176600
150	1.365900
175	1.173400
200	1.467600
225	1.158000
250	1.542000

TrainOutput(global\_step=250, training\_loss=1.3615322799682616, metrics={'train\_runtime': 1627.8174, 'train\_samples\_per\_second': 0.614, 'train\_steps\_per\_second': 0.154, 'total\_flos': 8755214198673920.0, 'train\_loss': 1.3615322799682616, 'epoch': 1.0})

[5] # Save trained model

Automatic saving failed. This file was updated remotely or in another tab. Show diff

Resources

You are not subscribed. Learn more.  
 You currently have zero compute units available. Resources offered free of charge are not guaranteed. Purchase more units here.  
 At your current usage level, this runtime may last up to 2 hours 20 minutes.

Manage sessions

Want more memory and disk space? Upgrade to Colab Pro

Python 3 Google Compute Engine backend (GPU)  
 Showing resources from 4:50 PM to 5:22 PM

System RAM: 3.2 / 12.7 GB  
 GPU RAM: 10.8 / 15.0 GB  
 Disk: 43.5 / 78.2 GB

Change runtime type

Activate Windows  
 Go to Settings to activate Windows.

0s completed at 5:22 PM

## Unique Idea Brief (Solution)

The unique idea brief outlines an innovative approach to developing a chatbot that stands out from existing solutions. This project leverages the power of large language models (LLMs) and generative AI (GenAI) to create a highly customizable and adaptable chatbot. The core idea is to fine-tune these advanced models to cater to specific applications, addressing the limitations of generic chatbots that often fail to meet specialized needs.

The proposed solution aims to offer a chatbot that can understand and generate human-like responses with high accuracy and contextual relevance. This adaptability is crucial for various industries, such as customer support, healthcare, education, and more, where tailored interactions can significantly enhance user experience and operational efficiency.

One of the key aspects of this solution is its ability to learn and adapt over time. By continuously fine-tuning and updating the model with new data, the chatbot can improve its performance and stay relevant to changing user requirements and industry trends. This dynamic learning capability ensures that the chatbot remains effective and efficient, providing a superior user experience compared to static, rule-based systems.

The document likely highlights the technical innovations and methodologies employed to achieve this level of customization and adaptability. This includes leveraging advanced machine learning techniques, optimizing inference processes, and ensuring seamless integration with existing platforms and systems. The unique idea brief emphasizes the potential impact of this solution, showcasing its ability to transform how organizations interact with their users and streamline various operations.

By focusing on customization, adaptability, and continuous improvement, the project proposes a solution that not only addresses current limitations but also anticipates future needs. This forward-thinking approach positions the chatbot as a valuable tool for a wide range of applications, offering significant benefits in terms of user engagement, satisfaction, and overall efficiency.

## Team Members and Contribution

- **Atharva Ahire ( Team Leader ):** Observing and maintaining coordination between teammates as well as analyzing all the suggestions from the teammates and deciding how it can be implemented in efficient.
- **Vaishnavi Deshmukh:** Resolving the errors in the jupyter notebook for better results especially the configuration of pipeline of the chatbot code
- **Snehal Pawar :** Implementation of the fine-tuning model as well as Fine-tuning different models on trial and error basis to reduce the training time.
- **Rohit Gaikwad :** Found and tested alternative solutions to the errors in an efficient way like using google colab instead of jupyter lab for fine-tuning the model.
- **Vrushabh Patil :** Created a custom dataset of a specific domain and implemented it into the chatbot.

## Conclusion

The custom chatbot leverages advanced AI technologies, including large language models (LLMs) and generative AI (GenAI), to provide a highly adaptable and effective solution for various applications. By addressing the limitations of traditional chatbots, the project aims to deliver a tool that can handle specialized queries with high accuracy and contextual relevance.

The document highlights the successful implementation of key features, such as natural language understanding, contextual awareness, and continuous learning. These features ensure that the chatbot can provide meaningful and accurate responses, enhancing user engagement and satisfaction. The chatbot's ability to integrate with various platforms and systems further extends its utility, making it a versatile tool for diverse industries.

Looking forward, the project envisions ongoing refinement and improvement of the chatbot. This includes updating the model with new data, incorporating user feedback, and leveraging advancements in AI research to enhance performance. The document also emphasizes the potential for expanding the chatbot's capabilities, such as adding support for more languages, improving scalability, and enhancing security measures.

In summary, the custom chatbot represents a significant advancement in AI-driven interactions, offering a powerful tool for organizations to improve their services and user experience. The project's systematic approach to development and deployment, combined with continuous learning and adaptation, ensures that the chatbot remains relevant and effective in addressing evolving user needs and industry trends.