

Using Google Colab and Google Drive for Training and Fine-Tuning LLaMA 3.1 M

This document summarizes the steps to use Google Colab, Google Drive, and a quantized LLaMA 3.1 model for training and fine-tuning.

1. **Overview**:

- Using Google Colab for training/fine-tuning models on free GPUs.
- Mounting Google Drive to persist models, training data, and checkpoints.
- Working with a quantized LLaMA 3.1 model (INT8 quantization) for efficient usage of resources.

2. **Step-by-Step Guide**:

a) **Clone Your Git Repository in Google Colab**:

- Clone your Git repository into Colab using:

```
```python
!git clone https://github.com/your-username/your-repository.git

%cd your-repository
```
```

b) **Mount Google Drive**:

- Mount Google Drive to access storage for models and data:

```
```python
from google.colab import drive

drive.mount('/content/drive')
```
```

c) ****Set Up Your Environment****:

- Install dependencies from your Git repo (usually via `requirements.txt`):

```
```python  

!pip install -r requirements.txt
...`
```

d) **\*\*Load and Quantize Model\*\***:

- Load pre-trained LLaMA model and quantize it to INT8:

```
```python  
  
from transformers import AutoModelForCausalLM, AutoTokenizer  
  
model = AutoModelForCausalLM.from_pretrained('your-model-name')  
  
model = model.quantize(bits=8) # Example quantization  
...`
```

e) ****Prepare Data****:

- Load your training data from Google Drive:

```
```python  

data_path = '/content/drive/MyDrive/data/train_data.txt'
...`
```

f) **\*\*Training or Fine-Tuning\*\***:

- Run the training script on Colab:

```
```python  
  
!python train.py --model_name_or_path 'your-model-path' --train_data 'your-train-data-path'  
...`
```

g) ****Persist Models****:

- Save model weights and checkpoints to Google Drive:

```
```python  

model.save_pretrained('/content/drive/MyDrive/your-models')

tokenizer.save_pretrained('/content/drive/MyDrive/your-models')

```
```

3. **Key Notes**:

- **Google Colab** provides limited session time (12 hours), so persistent storage on **Google Drive** is essential.
- **Quantized models** (INT8) help save memory and disk space, making them ideal for training and experimentation.
- **Google Drive** offers sufficient space for smaller models but may require an upgraded plan for larger models.
- Using a **Git repository** ensures that the training scripts and configurations are version-controlled and reproducible.

4. **Recommendations**:

- For large models like LLaMA 7B, consider using quantized versions (INT8) to reduce resource consumption.
- Experiment with training and fine-tuning on **Google Colab** with smaller datasets before moving to larger datasets.