

❑ *Problem Statement:*

Finetuning a pre-trained LLM using a dataset for text generation problem using LoRA technique.

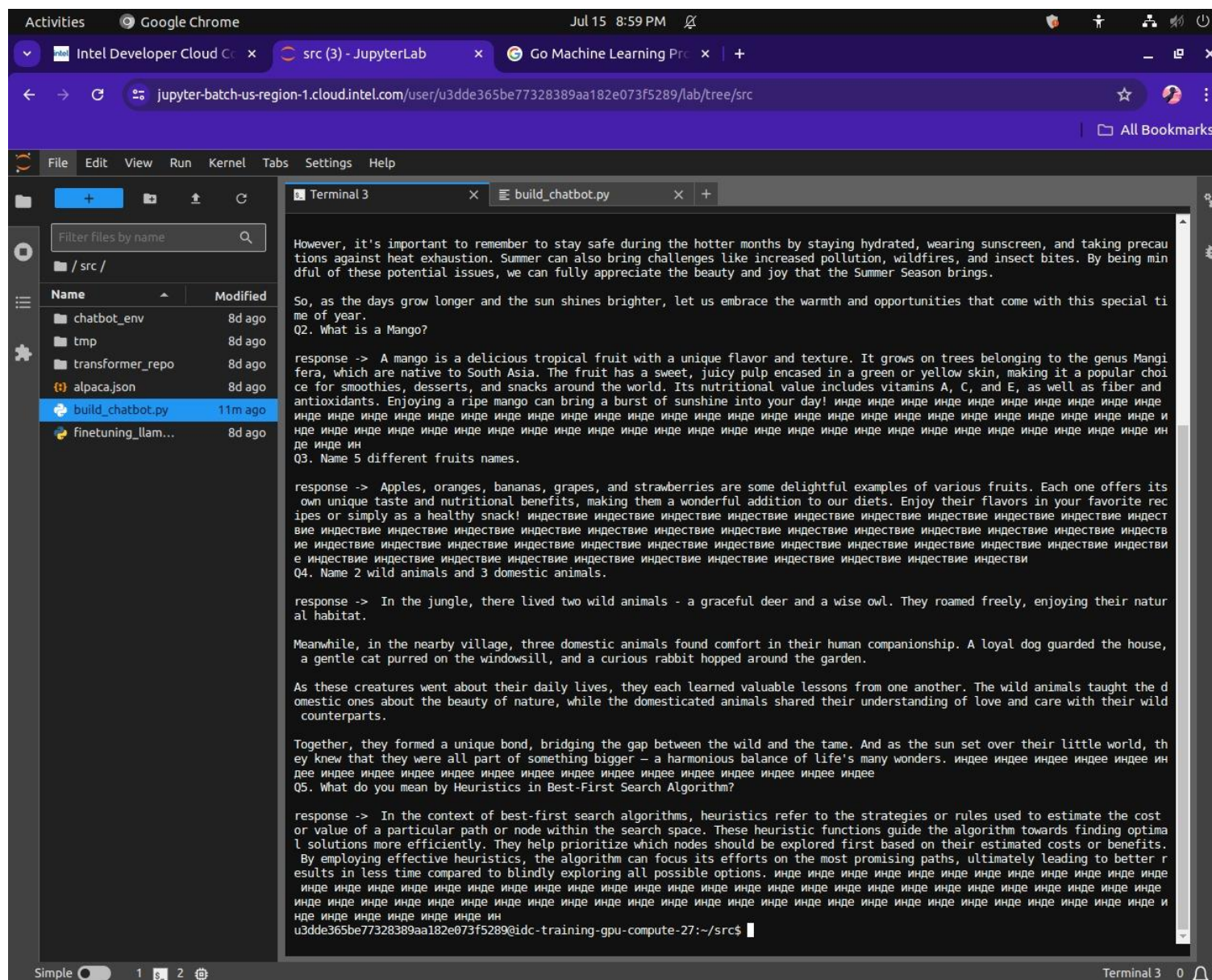
❑ *Solution:*

In deep learning, **fine-tuning** is an approach to transfer learning in which the parameters of a pre-trained model are trained on new data. Fine-tuning can be done on the entire neural network, or on only a subset of its layers, in which case the layers that are not being fine-tuned are "frozen". A model may also be augmented with "adapters" that consist of far fewer parameters than the original model, and fine-tuned in a parameter-efficient way by tuning the weights of the adapters and leaving the rest of the model's weights frozen.

Models that are pre-trained on large and general corpora are usually fine-tuned by reusing the model's parameters as a starting point and adding a task-specific layer trained from scratch. Fine-tuning the full model is common as well and often yields better results, but it is more computationally expensive.

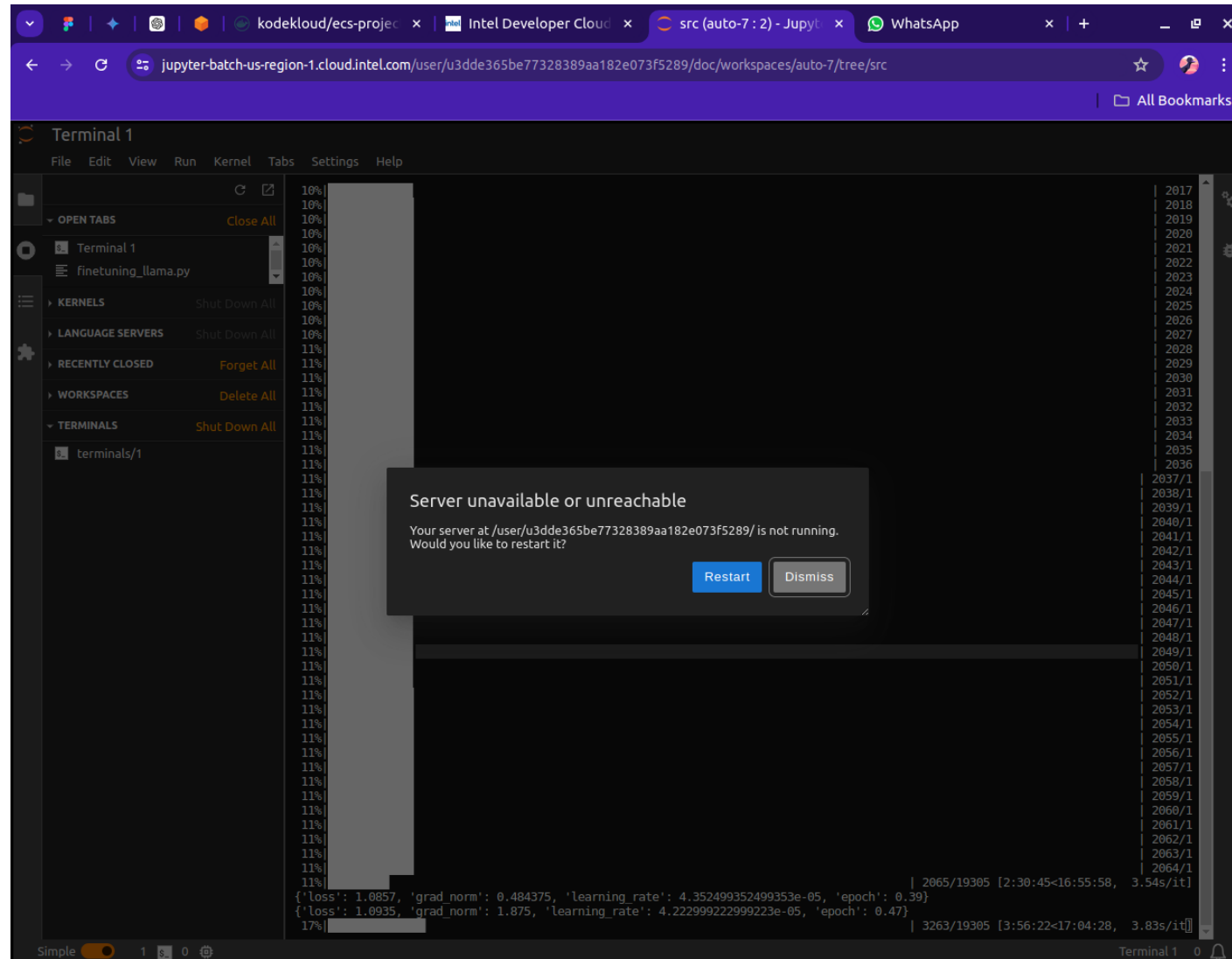
Fine-tuning is typically accomplished with supervised learning, but there are also techniques to fine-tune a model using weak supervision. Fine-tuning can be combined with a reinforcement learning from human feedback-based objective to produce language models like ChatGPT and Sparrow.

Observations:




```
2024-07-06 17:20:52,639 - finetuning.py - intel_extension_for_transformers.transformers.llm.finetuning.finetuning - INFO - Using data
collator of type DataCollatorForSeq2Seq
trainable params: 4,194,304 || all params: 6,742,609,920 || trainable%: 0.06220594176090199
[INFO|trainer.py:641] 2024-07-06 17:20:53,748 >> Using cpu_amp half precision backend
[INFO|trainer.py:2078] 2024-07-06 17:20:54,320 >> ***** Running training *****
[INFO|trainer.py:2079] 2024-07-06 17:20:54,320 >>   Num examples = 51,482
[INFO|trainer.py:2080] 2024-07-06 17:20:54,320 >>   Num Epochs = 3
[INFO|trainer.py:2081] 2024-07-06 17:20:54,320 >>   Instantaneous batch size per device = 4
[INFO|trainer.py:2084] 2024-07-06 17:20:54,320 >>   Total train batch size (w. parallel, distributed & accumulation) = 8
[INFO|trainer.py:2085] 2024-07-06 17:20:54,321 >>   Gradient Accumulation steps = 2
[INFO|trainer.py:2086] 2024-07-06 17:20:54,321 >>   Total optimization steps = 19,305
[INFO|trainer.py:2087] 2024-07-06 17:20:54,324 >>   Number of trainable parameters = 4,194,304
0%|
```

| 0/19305 [00:00<?, ?it/s]



Time to complete finetuning using LoRA is 17 hours.

Advantages:

- Efficiency in training and adaptation of large language models
- Reduced computational resources requirement
- Preservation of pre-trained model weights
- Reduced parameter overhead
- Efficient task-switching
- Faster iterations
- Scalability
- Memory efficiency

what will happen if we don't use LoRA?

If we don't use LoRa several potential drawbacks and limitations could arise:

1. Increased Computational Resources
2. Longer Training Time
3. Higher Costs
4. Overfitting Risks
5. Scalability Issues
6. Difficulty in Transfer Learning

❑ *Features Offered:*

Currently, we've fine-tuned the LLM for text generation using the Alpaca dataset. To extend its capabilities for generating project reports, we'll incorporate datasets from technical, business, and academic reports and implement advanced fine-tuning strategies.

❑ *Technologies Used:*

- Intel Developer Cloud Platform
- Python
- Huggingface

❑ *Team Contribution*

- **Vaishnavi Ankatwar** - Undertook the comprehensive research and meticulously prepared the project report.
- **Tanishq Palandurkar** - Provided an insightful presentation on LoRA and its advantages, drawing from his observations during the fine-tuning and execution processes on IDC.
- **Vedanti Dhage** – Committed to assisting Vaishnavi with the project reports and research efforts.
- **Sanjana Sanjai** – Created the IDC account.

❑ *Conclusion:*

Consequently, we gained valuable insights into the process and workflow of building a chatbot and fine-tuning it to address specific problems on Intel's cloud development platform.