

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 5.4 MB/s eta 0:00:00
_____ 72.9/72.9 kB 10.1 MB/s eta 0:00:00
_____ 92.5/92.5 MB 9.0 MB/s eta 0:00:00
_____ 7.4/7.4 MB 53.7 MB/s eta 0:00:00
_____ 77.4/77.4 kB 8.4 MB/s eta 0:00:00
_____ 7.8/7.8 MB 72.7 MB/s eta 0:00:00
_____ 521.2/521.2 kB 43.2 MB/s eta 0:00:00
_____ 115.3/115.3 kB 14.8 MB/s eta 0:00:00
_____ 134.8/134.8 kB 19.5 MB/s eta 0:00:00

```

Step 2: Import All the Required Libraries

```

import os
import torch
from datasets import load_dataset
from transformers import (
    AutoModelForCausalLM,
    AutoTokenizer,
    BitsAndBytesConfig,
    HfArgumentParser,
    TrainingArguments,
    pipeline,
    logging,
)
from peft import LoraConfig, PeftModel
from trl import SFTTrainer

```

In case of Llama 2, the following prompt template is used for the chat models

System Prompt (optional) to guide the model

User prompt (required) to give the instruction

Model Answer (required)

```

<s>[INST] <<SYS>>
System prompt
<</SYS>>

User prompt [/INST] Model answer </s>

```

We will reformat our instruction dataset to follow Llama 2 template.

- Original Dataset: <https://huggingface.co/datasets/timdettmers/openassistant-guanaco>
- Reformat Dataset following the Llama 2 template with 1k sample: <https://huggingface.co/datasets/mlabonne/guanaco-llama2-1k>
- Complete Reformat Dataset following the Llama 2 template: <https://huggingface.co/datasets/mlabonne/guanaco-llama2>

To know how this dataset was created, you can check this notebook.

<https://colab.research.google.com/drive/1Ad7a9zMmkxuXTOh1Z7-rNSICA4dybpM2?usp=sharing>

Note: You don't need to follow a specific prompt template if you're using the base Llama 2 model instead of the chat version.

How to fine tune Llama 2

- Free Google Colab offers a 15GB Graphics Card (Limited Resources --> Barely enough to store Llama 2-7b's weights)
- We also need to consider the overhead due to optimizer states, gradients, and forward activations
- Full fine-tuning is not possible here: we need parameter-efficient fine-tuning (PEFT) techniques like LoRA or QLoRA.
- To drastically reduce the VRAM usage, we must fine-tune the model in 4-bit precision, which is why we'll use QLoRA here.

Step 3

1. Load a llama-2-7b-chat-hf model (chat model)
2. Train it on the mlabonne/guanaco-llama2-1k (1,000 samples), which will produce our fine-tuned model Llama-2-7b-chat-finetune

QLoRA will use a rank of 64 with a scaling parameter of 16. We'll load the Llama 2 model directly in 4-bit precision using the NF4 type and train it for one epoch

```
# The model that you want to train from the Hugging Face hub
model_name = "NousResearch/Llama-2-7b-chat-hf"

# The instruction dataset to use
dataset_name = "mlabonne/guanaco-llama2-1k"

# Fine-tuned model name
new_model = "Llama-2-7b-chat-finetune"

#####
# QLoRA parameters
#####

# LoRA attention dimension
lora_r = 64

# Alpha parameter for LoRA scaling
lora_alpha = 16

# Dropout probability for LoRA layers
lora_dropout = 0.1

#####
# bitsandbytes parameters
#####

# Activate 4-bit precision base model loading
use_4bit = True

# Compute dtype for 4-bit base models
bnb_4bit_compute_dtype = "float16"

# Quantization type (fp4 or nf4)
bnb_4bit_quant_type = "nf4"

# Activate nested quantization for 4-bit base models (double quantization)
use_nested_quant = False

#####
# TrainingArguments parameters
#####

# Output directory where the model predictions and checkpoints will be stored
output_dir = "./results"

# Number of training epochs
num_train_epochs = 1

# Enable fp16/bf16 training (set bf16 to True with an A100)
fp16 = False
bf16 = False

# Batch size per GPU for training
per_device_train_batch_size = 4
```

```

# Batch size per GPU for evaluation
per_device_eval_batch_size = 4

# Number of update steps to accumulate the gradients for
gradient_accumulation_steps = 1

# Enable gradient checkpointing
gradient_checkpointing = True

# Maximum gradient normal (gradient clipping)
max_grad_norm = 0.3

# Initial learning rate (AdamW optimizer)
learning_rate = 2e-4

# Weight decay to apply to all layers except bias/LayerNorm weights
weight_decay = 0.001

# Optimizer to use
optim = "paged_adamw_32bit"

# Learning rate schedule
lr_scheduler_type = "cosine"

# Number of training steps (overrides num_train_epochs)
max_steps = -1

# Ratio of steps for a linear warmup (from 0 to learning rate)
warmup_ratio = 0.03

# Group sequences into batches with same length
# Saves memory and speeds up training considerably
group_by_length = True

# Save checkpoint every X updates steps
save_steps = 0

# Log every X updates steps
logging_steps = 25

#####
# SFT parameters
#####

# Maximum sequence length to use
max_seq_length = None

# Pack multiple short examples in the same input sequence to increase efficiency
packing = False

# Load the entire model on the GPU 0
device_map = {"": 0}

```

▼ Step 4: Load everything and start the fine-tuning process

1. First of all, we want to load the dataset we defined. Here, our dataset is already preprocessed but, usually, this is where you would reformat the prompt, filter out bad text, combine multiple datasets, etc.
2. Then, we're configuring bitsandbytes for 4-bit quantization.
3. Next, we're loading the Llama 2 model in 4-bit precision on a GPU with the corresponding tokenizer.
4. Finally, we're loading configurations for QLoRA, regular training parameters, and passing everything to the SFTTrainer. The training can finally start!

```

# Load dataset (you can process it here)
dataset = load_dataset(dataset_name, split="train")

# Load tokenizer and model with QLoRA configuration
compute_dtype = getattr(torch, bnb_4bit_compute_dtype)

bnb_config = BitsAndBytesConfig(
    load_in_4bit=use_4bit,
    bnb_4bit_quant_type=bnb_4bit_quant_type,
    bnb_4bit_compute_dtype=compute_dtype,
    bnb_4bit_use_double_quant=use_nested_quant,
)

```

```

# Check GPU compatibility with bfloat16
if compute_dtype == torch.float16 and use_4bit:
    major, _ = torch.cuda.get_device_capability()
    if major >= 8:
        print("=" * 80)
        print("Your GPU supports bfloat16: accelerate training with bf16=True")
        print("=" * 80)

# Load base model
model = AutoModelForCausalLM.from_pretrained(
    model_name,
    quantization_config=bnb_config,
    device_map=device_map
)
model.config.use_cache = False
model.config.pretraining_tp = 1

# Load LLaMA tokenizer
tokenizer = AutoTokenizer.from_pretrained(model_name, trust_remote_code=True)
tokenizer.pad_token = tokenizer.eos_token
tokenizer.padding_side = "right" # Fix weird overflow issue with fp16 training

# Load LoRA configuration
peft_config = LoraConfig(
    lora_alpha=lora_alpha,
    lora_dropout=lora_dropout,
    r=lora_r,
    bias="none",
    task_type="CAUSAL_LM",
)

# Set training parameters
training_arguments = TrainingArguments(
    output_dir=output_dir,
    num_train_epochs=num_train_epochs,
    per_device_train_batch_size=per_device_train_batch_size,
    gradient_accumulation_steps=gradient_accumulation_steps,
    optim=optim,
    save_steps=save_steps,
    logging_steps=logging_steps,
    learning_rate=learning_rate,
    weight_decay=weight_decay,
    fp16=fp16,
    bf16=bf16,
    max_grad_norm=max_grad_norm,
    max_steps=max_steps,
    warmup_ratio=warmup_ratio,
    group_by_length=group_by_length,
    lr_scheduler_type=lr_scheduler_type,
    report_to="tensorboard"
)

# Set supervised fine-tuning parameters
trainer = SFTTrainer(
    model=model,
    train_dataset=dataset,
    peft_config=peft_config,
    dataset_text_field="text",
    max_seq_length=max_seq_length,
    tokenizer=tokenizer,
    args=training_arguments,
    packing=packing,
)

# Train model
trainer.train()

```

```
Downloading readme: 100% 1.02k/1.02k [00:00<00:00, 57.7kB/s]
Downloading data files: 100% 1/1 [00:00<00:00, 2.49it/s]
Downloading data: 100% 967k/967k [00:00<00:00, 2.53MB/s]
Extracting data files: 100% 1/1 [00:00<00:00, 49.97it/s]
Generating train split: 100% 1000/1000 [00:00<00:00, 15875.19 examples/s]
config.json: 100% 583/583 [00:00<00:00, 47.3kB/s]
model.safetensors.index.json: 100% 26.8k/26.8k [00:00<00:00, 2.25MB/s]
Downloading shards: 100% 2/2 [01:09<00:00, 31.42s/it]
model-00001-of-00002.safetensors: 100% 9.98G/9.98G [00:54<00:00, 235MB/s]
model-00002-of-00002.safetensors: 100% 3.50G/3.50G [00:14<00:00, 207MB/s]
Loading checkpoint shards: 100% 2/2 [01:10<00:00, 31.87s/it]
generation_config.json: 100% 179/179 [00:00<00:00, 12.7kB/s]
tokenizer_config.json: 100% 746/746 [00:00<00:00, 55.8kB/s]
tokenizer.model: 100% 500k/500k [00:00<00:00, 29.4MB/s]
tokenizer.json: 100% 1.84M/1.84M [00:00<00:00, 7.32MB/s]
added_tokens.json: 100% 21.0/21.0 [00:00<00:00, 1.40kB/s]
special_tokens_map.json: 100% 435/435 [00:00<00:00, 32.1kB/s]
/usr/local/lib/python3.10/dist-packages/peft/utils/other.py:102: FutureWarning: prepare_model_for_int8_training is deprecated.
warnings.warn(
/usr/local/lib/python3.10/dist-packages/trl/trainer/sft_trainer.py:159: UserWarning: You didn't pass a `max_seq_length` argument, we will use the default value of 1024.
warnings.warn(
Map: 100% 1000/1000 [00:01<00:00, 897.76 examples/s]
You're using a LlamaTokenizerFast tokenizer. Please note that with a fast tokenizer, using the `__call__` method is discouraged as it will tokenize the prompts on the CPU.
/usr/local/lib/python3.10/dist-packages/torch/utils/checkpoint.py:429: UserWarning: torch.utils.checkpoint: please use `torch.utils.checkpoint.checkpoint` instead of `torch.utils.checkpoint.checkpoint_no_grad`
warnings.warn(
[250/250 26:57, Epoch 1/1]
```

Step Training Loss

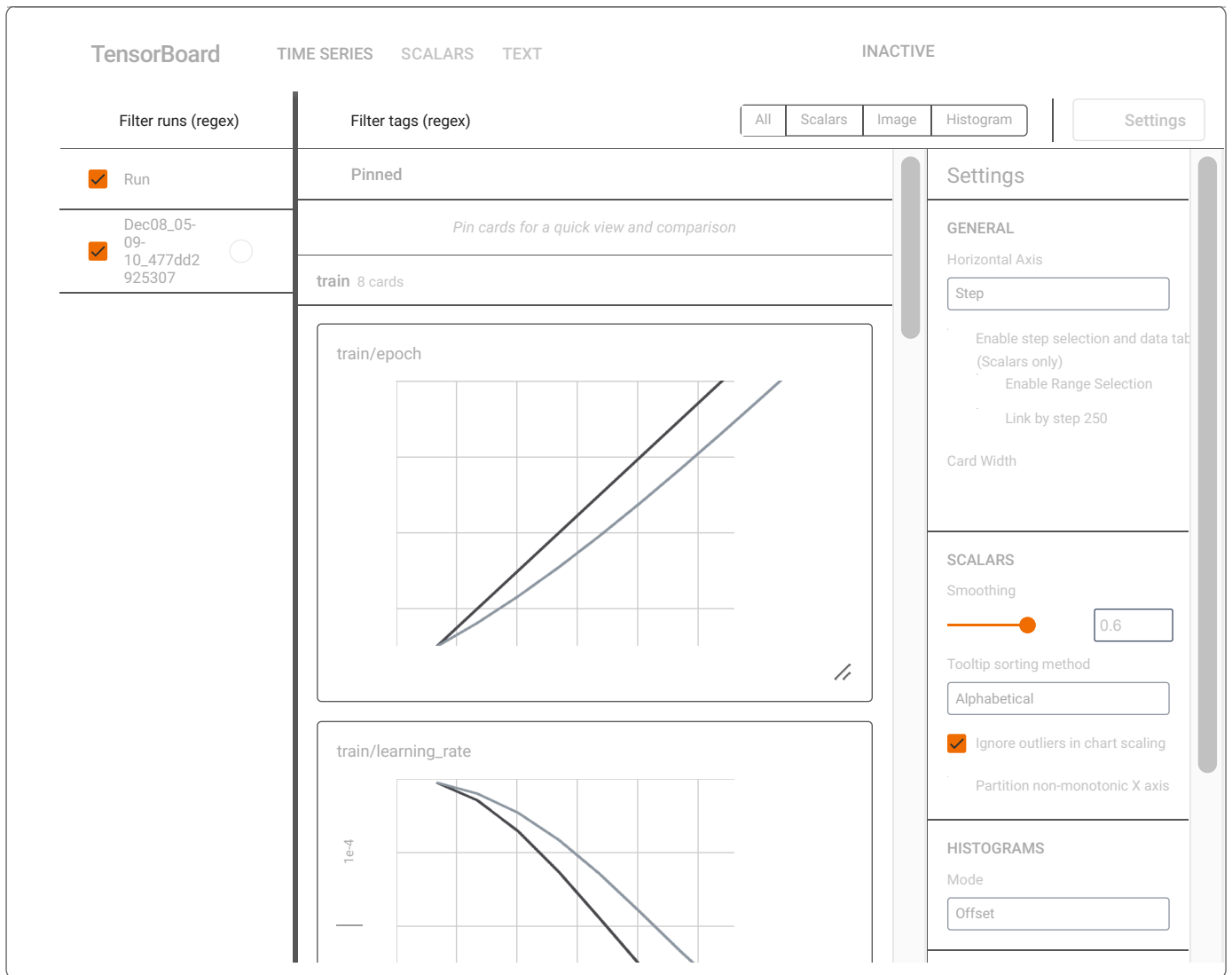
25	1.408600
50	1.662100
75	1.215200
100	1.444100
125	1.176900
150	1.367300
175	1.174000
200	1.468100
225	1.158200
250	1.542700

```
TrainOutput(global_step=250, training_loss=1.361721710205078, metrics={'train_runtime': 1634.449, 'train_samples_per_second': 0.612, 'train_steps_per_second': 0.153, 'total_flos': 8755214190673920.0, 'train_loss': 1.361721710205078, 'epoch': 1.0})
```

```
# Save trained model
trainer.model.save_pretrained(new_model)
```

Step 5: Check the plots on tensorboard, as follows

```
%load_ext tensorboard
%tensorboard --logdir results/runs
```



Step 6: Use the text generation pipeline to ask questions like “What is a large language model?” Note that I’m formatting the input to match Llama 2 prompt template.

```
# Ignore warnings
logging.set_verbosity(logging.CRITICAL)

# Run text generation pipeline with our next model
prompt = "What is a large language model?"
pipe = pipeline(task="text-generation", model=model, tokenizer=tokenizer, max_length=200)
result = pipe(f"<s>[INST] {prompt} [/INST]")
print(result[0]['generated_text'])
```

```
<s>[INST] What is a large language model? [/INST] A large language model is a type of artificial intelligence (AI) model. Large language models are neural networks that are trained on vast amounts of text data to predict the next word in a sequence. Large language models are trained on vast amounts of text data, such as books, articles, and websites. They are designed to generate human-like text. Large language models are used in a variety of applications, such as chatbots, text summarization, and machine translation.
```

```
# Empty VRAM
del model
del pipe
del trainer
import gc
gc.collect()
gc.collect()
```

```
20933
```

You can train a Llama 2 model on the entire dataset using [mlabonne/guanaco-llama2](https://github.com/mlabonne/guanaco-llama2)

Unfortunately, as far as I know, there is no straightforward way to do it: we need to reload the base model in FP16 precision and use the peft library to merge everything.

2/2 [01:01<00:00, 28.04s/it]