Finetuning Qwen2.5-7B-Instruct model using LoRA and QLoRA

```
# pip installs
!pip install -q --upgrade torch==2.5.1+cu124 torchvision==0.20.1+cu124 torchaudio==2.5.1+cu124 --index-u
!pip install -q requests bitsandbytes==0.46.0 transformers==4.48.3 accelerate==1.3.0
!pip install -q datasets requests peft
                                           -- 908.2/908.2 MB 1.2 MB/s eta 0:00:00
                                             - 7.3/7.3 MB 35.1 MB/s eta 0:00:00
                                          ---- 3.4/3.4 MB 43.5 MB/s eta 0:00:00
                                          ---- 24.6/24.6 MB 71.1 MB/s eta 0:00:00
                                           — 883.7/883.7 kB 40.4 MB/s eta 0:00:00
                                             - 13.8/13.8 MB 73.0 MB/s eta 0:00:00
                                             - 664.8/664.8 MB ? eta 0:00:00
                                           --- 363.4/363.4 MB 1.2 MB/s eta 0:00:00
                                            -- 211.5/211.5 MB 6.7 MB/s eta 0:00:00
                                              - 56.3/56.3 MB 11.7 MB/s eta 0:00:00
                                         ---- 127.9/127.9 MB <mark>8.9 MB/s</mark> eta 0:00:00
                                             - 207.5/207.5 MB 1.6 MB/s eta 0:00:00
                                            - 188.7/188.7 MB 6.8 MB/s eta 0:00:00
                                             - 99.1/99.1 kB <mark>9.2 MB/s</mark> eta 0:00:00
                                            - 21.1/21.1 MB 112.0 MB/s eta 0:00:00
                                              - 209.6/209.6 MB ? eta 0:00:00
                                          ---- 6.2/6.2 MB 81.4 MB/s eta 0:00:00
                                             - 44.4/44.4 kB 3.2 MB/s eta 0:00:00
                                           - 67.0/67.0 MB 12.9 MB/s eta 0:00:00
                                           - 9.7/9.7 MB 138.4 MB/s eta 0:00:00
                                           - 336.6/336.6 kB 29.6 MB/s eta 0:00:00
                                           - 3.1/3.1 MB 98.2 MB/s eta 0:00:00
```

```
# imports

import os
import re
import math
from tqdm import tqdm
from google.colab import userdata
from huggingface_hub import login
import torch
import transformers
from transformers
from transformers import AutoModelForCausalLM, AutoTokenizer, BitsAndBytesConfig, TrainingArguments, se
from peft import LoraConfig, PeftModel
from datetime import datetime
```

```
# Constants

BASE_MODEL = "Qwen/Qwen2.5-7B-Instruct"

# Hyperparameters for QLoRA Fine-Tuning

LORA_R = 32
LORA_ALPHA = 64
TARGET_MODULES = ["q_proj", "v_proj", "k_proj", "o_proj"]
```

Log in to HuggingFace

```
# Log in to HuggingFace

hf_token = userdata.get('HF_TOKEN')
login(hf_token, add_to_git_credential=True)
```

Analyzing the base model

```
# Load the Base Model without quantization
base model = AutoModelForCausalLM.from pretrained(BASE MODEL, device map="auto")
config.json: 100%
                                                               663/663 [00:00<00:00, 77.7kB/s]
model.safetensors.index.json:
                                27.8k/? [00:00<00:00, 2.17MB/s]
Downloading shards: 100%
                                                                        4/4 [04:32<00:00, 67.03s/it]
model-00001-of-00004.safetensors: 100%
                                                                                    3.95G/3.95G [01:03<00:00, 54.8MB/s]
model-00002-of-00004.safetensors: 100%
                                                                                    3.86G/3.86G [01:22<00:00, 75.0MB/s]
model-00003-of-00004.safetensors: 100%
                                                                                    3.86G/3.86G [01:01<00:00, 120MB/s]
model-00004-of-00004.safetensors: 100%
                                                                                    3.56G/3.56G [01:03<00:00, 60.8MB/s]
Loading checkpoint shards: 100%
                                                                              4/4 [00:55<00:00, 18.24s/it]
generation_config.json: 100%
                                                                         243/243 [00:00<00:00, 26.4kB/s]
WARNING:accelerate.big modeling:Some parameters are on the meta device because they were offloaded to the
```

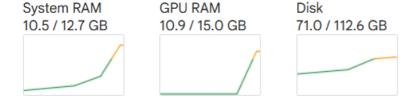
Double-click (or enter) to edit

```
print(f"Memory footprint: {base_model.get_memory_footprint() / 1e9:,.1f} GB")

Memory footprint: 30.5 GB
```

```
base_model
Qwen2ForCausalLM(
  (model): Qwen2Model(
    (embed_tokens): Embedding(152064, 3584)
    (layers): ModuleList(
      (0-27): 28 x Qwen2DecoderLayer(
        (self_attn): Qwen2Attention(
          (q_proj): Linear(in_features=3584, out_features=3584, bias=True)
          (k_proj): Linear(in_features=3584, out_features=512, bias=True)
(v_proj): Linear(in_features=3584, out_features=512, bias=True)
           (o_proj): Linear(in_features=3584, out_features=3584, bias=False)
        (mlp): Qwen2MLP(
           (gate_proj): Linear(in_features=3584, out_features=18944, bias=False)
           (up_proj): Linear(in_features=3584, out_features=18944, bias=False)
           (down_proj): Linear(in_features=18944, out_features=3584, bias=False)
          (act_fn): SiLU()
        (input_layernorm): Qwen2RMSNorm((3584,), eps=1e-06)
        (post_attention_layernorm): Qwen2RMSNorm((3584,), eps=1e-06)
    (norm): Qwen2RMSNorm((3584,), eps=1e-06)
    (rotary_emb): Qwen2RotaryEmbedding()
  (lm_head): Linear(in_features=3584, out_features=152064, bias=False)
```

Python 3 Google Compute Engine backend (GPU) Showing resources from 23:09 to 23:23



Quantizing with 8 bits

```
# Load the Base Model using 8 bit

quant_config = BitsAndBytesConfig(load_in_8bit=True)

base_model = AutoModelForCausalLM.from_pretrained(
    BASE_MODEL,
    quantization_config=quant_config,
    device_map="auto",
)

Loading checkpoint shards: 100%

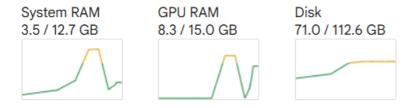
4/4 [01:12<00:00, 17.54s/it]</pre>
```

Double-click (or enter) to edit

```
print(f"Memory footprint: {base_model.get_memory_footprint() / 1e9:,.1f} GB")
Memory footprint: 8.7 GB
```

```
base_model
Qwen2ForCausalLM(
  (model): Owen2Model(
    (embed_tokens): Embedding(152064, 3584)
    (layers): ModuleList(
      (0-27): 28 x Qwen2DecoderLayer(
        (self_attn): Qwen2Attention(
          (q_proj): Linear8bitLt(in_features=3584, out_features=3584, bias=True)
          (k_proj): Linear8bitLt(in_features=3584, out_features=512, bias=True)
          (v_proj): Linear8bitLt(in_features=3584, out_features=512, bias=True)
          (o_proj): Linear8bitLt(in_features=3584, out_features=3584, bias=False)
        (mlp): Qwen2MLP(
          (gate_proj): Linear8bitLt(in_features=3584, out_features=18944, bias=False)
          (up_proj): Linear8bitLt(in_features=3584, out_features=18944, bias=False)
          (down_proj): Linear8bitLt(in_features=18944, out_features=3584, bias=False)
          (act_fn): SiLU()
        (input layernorm): Qwen2RMSNorm((3584,), eps=1e-06)
        (post_attention_layernorm): Qwen2RMSNorm((3584,), eps=1e-06)
    (norm): Qwen2RMSNorm((3584,), eps=1e-06)
    (rotary_emb): Qwen2RotaryEmbedding()
  (lm_head): Linear(in_features=3584, out_features=152064, bias=False)
)
```

Python 3 Google Compute Engine backend (GPU) Showing resources from 23:09 to 23:30



Double-click (or enter) to edit

Double Quantization with 4 bits

```
# Load the Tokenizer and the Base Model using 4 bit

quant_config = BitsAndBytesConfig(
    load_in_4bit=True,
    bnb_4bit_use_double_quant=True,
    bnb_4bit_compute_dtype=torch.bfloat16,
    bnb_4bit_quant_type="nf4")

base_model = AutoModelForCausalLM.from_pretrained(
    BASE_MODEL,
    quantization_config=quant_config,
    device_map="auto",
)

Loading checkpoint shards: 100%

4/4 [01:15<00:00, 18.41s/it]</pre>
```

```
print(f"Memory footprint: {base_model.get_memory_footprint() / 1e9:,.2f} GB")
Memory footprint: 5.44 GB
```

```
base_model
Qwen2ForCausalLM(
  (model): Qwen2Model(
    (embed_tokens): Embedding(152064, 3584)
    (layers): ModuleList(
      (0-27): 28 x Qwen2DecoderLayer(
        (self_attn): Qwen2Attention(
           (q_proj): Linear4bit(in_features=3584, out_features=3584, bias=True)
          (k_proj): Linear4bit(in_features=3584, out_features=512, bias=True)
(v_proj): Linear4bit(in_features=3584, out_features=512, bias=True)
           (o_proj): Linear4bit(in_features=3584, out_features=3584, bias=False)
        (mlp): Qwen2MLP(
          (gate_proj): Linear4bit(in_features=3584, out_features=18944, bias=False)
           (up proj): Linear4bit(in features=3584, out features=18944, bias=False)
           (down_proj): Linear4bit(in_features=18944, out_features=3584, bias=False)
           (act_fn): SiLU()
        (input_layernorm): Qwen2RMSNorm((3584,), eps=1e-06)
        (post_attention_layernorm): Qwen2RMSNorm((3584,), eps=1e-06)
    (norm): Qwen2RMSNorm((3584,), eps=1e-06)
    (rotary_emb): Qwen2RotaryEmbedding()
```

(lm_head): Linear(in_features=3584, out_features=152064, bias=False)
)

Python 3 Google Compute Engine backend (GPU) Showing resources from 23:09 to 23:34

