- Install the required libraries/frameworks

```
!pip install --no-deps bitsandbytes accelerate xformers==0.0.29.post3 peft trl==0.15.2 triton cut_cross_entropy unsloth_zoo
!pip install sentencepiece protobuf datasets huggingface_hub hf_transfer
!pip install --no-deps unsloth
```

```
Requirement already satisfied: aiohttp in /usr/local/lib/python3.11/dist-packages (from datasets) (3.11.15)
Requirement already satisfied: packaging in /usr/local/lib/python3.11/dist-packages (from datasets) (24.2)
Requirement already satisfied: pyyaml>=5.1 in /usr/local/lib/python3.11/dist-packages (from datasets) (6.0.2)
Requirement already satisfied: typing-extensions>=3.7.4.3 in /usr/local/lib/python3.11/dist-packages (from huggingface_h
Requirement already satisfied: aiohappyeyeballs>=2.3.0 in /usr/local/lib/python3.11/dist-packages (from aiohttp->dataset
Requirement already satisfied: aiosignal>=1.1.2 in /usr/local/lib/python3.11/dist-packages (from aiohttp->datasets) (1.3
Requirement already satisfied: attrs>=17.3.0 in /usr/local/lib/python3.11/dist-packages (from aiohttp->datasets) (25.3.0
Requirement already satisfied: frozenlist>=1.1.1 in /usr/local/lib/python3.11/dist-packages (from aiohttp->datasets) (1.
Requirement already satisfied: multidict<7.0,>=4.5 in /usr/local/lib/python3.11/dist-packages (from aiohttp->datasets) (
Requirement already satisfied: propcache>=0.2.0 in /usr/local/lib/python3.11/dist-packages (from aiohttp->datasets) (0.3
Requirement already satisfied: yarl<2.0,>=1.17.0 in /usr/local/lib/python3.11/dist-packages (from aiohttp->datasets) (1.
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.11/dist-packages (from requests>=2.32.
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.11/dist-packages (from requests>=2.32.2->datasets)
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.11/dist-packages (from requests>=2.32.2->dat
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.11/dist-packages (from requests>=2.32.2->dat
Requirement already satisfied: python-dateutil>=2.8.2 in /usr/local/lib/python3.11/dist-packages (from pandas->datasets)
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.11/dist-packages (from pandas->datasets) (2025.2)
Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/python3.11/dist-packages (from pandas->datasets) (2025.2
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.11/dist-packages (from python-dateutil>=2.8.2->pandas-
Downloading datasets-3.5.0-py3-none-any.whl (491 kB)
                                            491.2/491.2 kB 17.4 MB/s eta 0:00:00
Downloading hf_transfer-0.1.9-cp38-abi3-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (3.6 MB)
                                            3.6/3.6 MB 81.5 MB/s eta 0:00:00
Downloading dill-0.3.8-py3-none-any.whl (116 kB)
                                            116.3/116.3 kB 12.2 MB/s eta 0:00:00
Downloading fsspec-2024.12.0-py3-none-any.whl (183 kB)
                                            183.9/183.9 kB 19.9 MB/s eta 0:00:00
Downloading multiprocess-0.70.16-py311-none-any.whl (143 kB)
                                            143.5/143.5 kB 16.4 MB/s eta 0:00:00
Downloading xxhash-3.5.0-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (194 kB)
                                            194.8/194.8 kB 20.9 MB/s eta 0:00:00
Installing collected packages: xxhash, hf_transfer, fsspec, dill, multiprocess, datasets
  Attempting uninstall: fsspec
    Found existing installation: fsspec 2025.3.2
    Uninstalling fsspec-2025.3.2:
      Successfully uninstalled fsspec-2025.3.2
ERROR: pip's dependency resolver does not currently take into account all the packages that are installed. This behaviou
unsloth-zoo 2025.3.17 requires tyro, which is not installed.
unsloth-zoo 2025.3.17 requires protobuf<4.0.0, but you have protobuf 5.29.4 which is incompatible.
torch 2.6.0+cu124 requires nvidia-cublas-cu12==12.4.5.8; platform_system == "Linux" and platform_machine == "x86_64", bu
torch 2.6.0+cu124 requires nvidia-cuda-cupti-cu12==12.4.127; platform_system == "Linux" and platform_machine == "x86_64"
torch 2.6.0+cu124 requires nvidia-cuda-nvrtc-cu12==12.4.127; platform_system == "Linux" and platform_machine == "x86_64"
torch 2.6.0+cu124 requires nvidia-cuda-runtime-cu12==12.4.127; platform_system == "Linux" and platform_machine == "x86_6
torch 2.6.0+cu124 requires nvidia-cudnn-cu12==9.1.0.70; platform_system == "Linux" and platform_machine == "x86_64", but
torch 2.6.0+cu124 requires nvidia-cufft-cu12==11.2.1.3; platform_system == "Linux" and platform_machine == "x86_64", but
torch 2.6.0+cu124 requires nvidia-curand-cu12==10.3.5.147; platform_system == "Linux" and platform_machine == "x86_64",
torch 2.6.0+cu124 requires nvidia-cusolver-cu12==11.6.1.9; platform_system == "Linux" and platform_machine == "x86_64",
torch 2.6.0+cu124 requires nvidia-cusparse-cu12==12.3.1.170; platform_system == "Linux" and platform_machine == "x86_64"
torch 2.6.0+cu124 requires nvidia-nvjitlink-cu12==12.4.127; platform_system == "Linux" and platform_machine == "x86_64",
gcsfs 2025.3.2 requires fsspec==2025.3.2, but you have fsspec 2024.12.0 which is incompatible.
Successfully installed datasets-3.5.0 dill-0.3.8 fsspec-2024.12.0 hf_transfer-0.1.9 multiprocess-0.70.16 xxhash-3.5.0
Collecting unsloth
  Downloading unsloth-2025.3.19-py3-none-any.whl.metadata (46 kB)
                                            46.2/46.2 kB 3.0 MB/s eta 0:00:00
Downloading unsloth-2025.3.19-py3-none-any.whl (192 kB)
                                            192.7/192.7 kB 8.9 MB/s eta 0:00:00
Installing collected packages: unsloth
Successfully installed unsloth-2025.3.19
```

- Unsloth configuration and Model

```
from unsloth import FastLanguageModel
import torch

model, tokenizer = FastLanguageModel.from_pretrained(
    model_name = "unsloth/Meta-Llama-3.1-8B-bnb-4bit",
    max_seq_length=2048,
    dtype=None,
    load_in_4bit=True
)
```

```
==((====))==  Unsloth 2025.3.19: Fast Llama patching. Transformers: 4.51.3.
   \\   /|    Tesla T4. Num GPUs = 1. Max memory: 14.741 GB. Platform: Linux.
O^O/ \_/ \    Torch: 2.6.0+cu124. CUDA: 7.5. CUDA Toolkit: 12.4. Triton: 3.2.0
\        /    Bfloat16 = FALSE. FA [Xformers = 0.0.29.post3. FA2 = False]
 "-____-"     Free license: http://github.com/unslothai/unsloth
Unsloth: Fast downloading is enabled — ignore downloading bars which are red colored!
/usr/local/lib/python3.11/dist-packages/transformers/quantizers/auto.py:212: UserWarning: You passed `quantization_confi
  warnings.warn(warning_msg)
```

model.safetensors: 100%                                      5.70G/5.70G [00:39<00:00, 223MB/s]

generation_config.json: 100%                                 235/235 [00:00<00:00, 20.6kB/s]

tokenizer_config.json: 100%                                  50.6k/50.6k [00:00<00:00, 4.20MB/s]

special_tokens_map.json: 100%                                459/459 [00:00<00:00, 49.8kB/s]

tokenizer.json: 100%                                         17.2M/17.2M [00:00<00:00, 49.4MB/s]

- LoRA Config

```
model = FastLanguageModel.get_peft_model(
    model,
    r = 8,
    target_modules = ["q_proj", "k_proj", "v_proj", "o_proj", "gate_proj", "up_proj", "down_proj"],
    lora_alpha = 16,
    lora_dropout = 0,
    bias = "none",
    use_gradient_checkpointing = "unsloth",
    random_state = 6666,
    use_rslora = False,
    loftq_config = None
)
```

```
Unsloth 2025.3.19 patched 32 layers with 32 QKV layers, 32 O layers and 32 MLP layers.
```

- Load and prepare Dataset

```
alpaca_prompt = """Below is an instruction that describes a task, paired with an input that provides further context. Write

### Instruction:
{}

### Input:
{}

### Response:
{}"""

EOS_TOKEN = tokenizer.eos_token
def formatting_prompts_func(examples):
    instructions = examples["instruction"]
    inputs       = examples["input"]
    outputs      = examples["output"]
    texts = []
    for instruction, input, output in zip(instructions, inputs, outputs):
        text = alpaca_prompt.format(instruction, input, output) + EOS_TOKEN
        texts.append(text)
    return { "text" : texts, }
pass

from datasets import load_dataset
dataset = load_dataset("Vezora/Tested-143k-Python-Alpaca", split = "train")
dataset = dataset.map(formatting_prompts_func, batched = True,)
```

README.md: 100%                                              4.16k/4.16k [00:00<00:00, 93.4kB/s]

143k-Tested-Python-Alpaca-Vezora.json: 100%                  295M/295M [00:04<00:00, 29.2MB/s]

Generating train split: 100%                                 143327/143327 [00:03<00:00, 39514.02 examples/s]

Map: 100%                                                    143327/143327 [00:02<00:00, 63089.14 examples/s]

- Model Training

```
from trl import SFTTrainer
from transformers import TrainingArguments
from unsloth import is_bfloat16_supported
```

```python
trainer = SFTTrainer(
    model = model,
    tokenizer = tokenizer,
    train_dataset = dataset,
    dataset_text_field = "text",
    max_seq_length = 2048,
    dataset_num_proc = 2,
    packing = False,
    args = TrainingArguments(
        per_device_train_batch_size = 2,
        gradient_accumulation_steps = 4,
        warmup_steps = 5,
        num_train_epochs = 1,
        learning_rate = 2e-4,
        fp16 = not is_bfloat16_supported(),
        bf16 = is_bfloat16_supported(),
        logging_steps = 1,
        optim = "adamw_8bit",
        weight_decay = 0.01,
        lr_scheduler_type = "linear",
        seed = 3407,
        output_dir = "outputs",
        report_to = "none",
    ),
)
```

- VRam Stats

```python
gpu_stats = torch.cuda.get_device_properties(0)
start_gpu_memory = round(torch.cuda.max_memory_reserved() / 1024 / 1024 / 1024, 3)
max_memory = round(gpu_stats.total_memory / 1024 / 1024 / 1024, 3)
print(f"GPU = {gpu_stats.name}. Max memory = {max_memory} GB.")
print(f"{start_gpu_memory} GB of memory reserved.")
```

```
GPU = Tesla T4. Max memory = 14.741 GB.
7.137 GB of memory reserved.
```

```python
trainer.train()
```

```
==((====))==  Unsloth - 2x faster free finetuning | Num GPUs used = 1
   \\   /|    Num examples = 143,327 | Num Epochs = 1 | Total steps = 17,916
O^O/ \_/ \    Batch size per device = 2 | Gradient accumulation steps = 4
\        /    Data Parallel GPUs = 1 | Total batch size (2 x 4 x 1) = 8
 "-____-"     Trainable parameters = 20,971,520/8,000,000,000 (0.26% trained)
Unsloth: Will smartly offload gradients to save VRAM!
```

[ 68/17916 14:56 < 67:20:16, 0.07 it/s, Epoch 0.00/1]

| Step | Training Loss |
|------|---------------|
| 1 | 1.053400 |
| 2 | 1.292700 |
| 3 | 1.163700 |
| 4 | 1.066800 |
| 5 | 1.010600 |
| 6 | 0.852200 |
| 7 | 0.684800 |
| 8 | 0.875200 |
| 9 | 0.986400 |
| 10 | 0.950300 |
| 11 | 0.760300 |
| 12 | 0.836300 |
| 13 | 0.995200 |
| 14 | 0.786700 |
| 15 | 0.603200 |
| 16 | 0.704100 |
| 17 | 1.061300 |
| 18 | 0.788600 |
| 19 | 0.828200 |
| 20 | 0.838500 |
| 21 | 0.600900 |
| 22 | 0.953100 |
| 23 | 0.780000 |
| 24 | 0.678000 |
| 25 | 0.739300 |
| 26 | 0.830300 |
| 27 | 0.961600 |
| 28 | 0.673400 |
| 29 | 0.536000 |
| 30 | 0.789100 |
| 31 | 0.639200 |
| 32 | 0.588000 |
| 33 | 0.750800 |
| 34 | 0.790900 |
| 35 | 0.665600 |
| 36 | 0.666100 |
| 37 | 0.560100 |
| 38 | 0.772200 |
| 39 | 0.670800 |
| 40 | 0.785900 |
| 41 | 0.826600 |
| 42 | 0.754300 |
| 43 | 0.757600 |
| 44 | 0.666500 |
| 45 | 0.631700 |

| | |
|---|---|
| 46 | 0.672500 |
| 47 | 0.729200 |
| 48 | 0.576600 |
| 49 | 1.155200 |
| 50 | 0.539700 |
| 51 | 0.734500 |
| 52 | 0.750100 |
| 53 | 0.732800 |
| 54 | 0.772800 |
| 55 | 0.793300 |
| 56 | 0.716900 |
| 57 | 0.550300 |
| 58 | 0.648900 |
| 59 | 0.842800 |
| 60 | 0.558900 |
| 61 | 0.786900 |
| 62 | 0.559000 |
| 63 | 0.596100 |
| 64 | 0.677200 |
| 65 | 0.670000 |
| 66 | 0.618400 |

[ 135/17916 30:22 < 67:41:58, 0.07 it/s, Epoch 0.01/1]

| Step | Training Loss |
|---|---|
| 1 | 1.053400 |
| 2 | 1.292700 |
| 3 | 1.163700 |
| 4 | 1.066800 |
| 5 | 1.010600 |
| 6 | 0.852200 |
| 7 | 0.684800 |
| 8 | 0.875200 |
| 9 | 0.986400 |
| 10 | 0.950300 |

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 50 minutes.

**Manage sessions**

Want more memory and disk space? Upgrade to Colab Pro     ✕

Not connected to runtime.

**Change runtime type**