

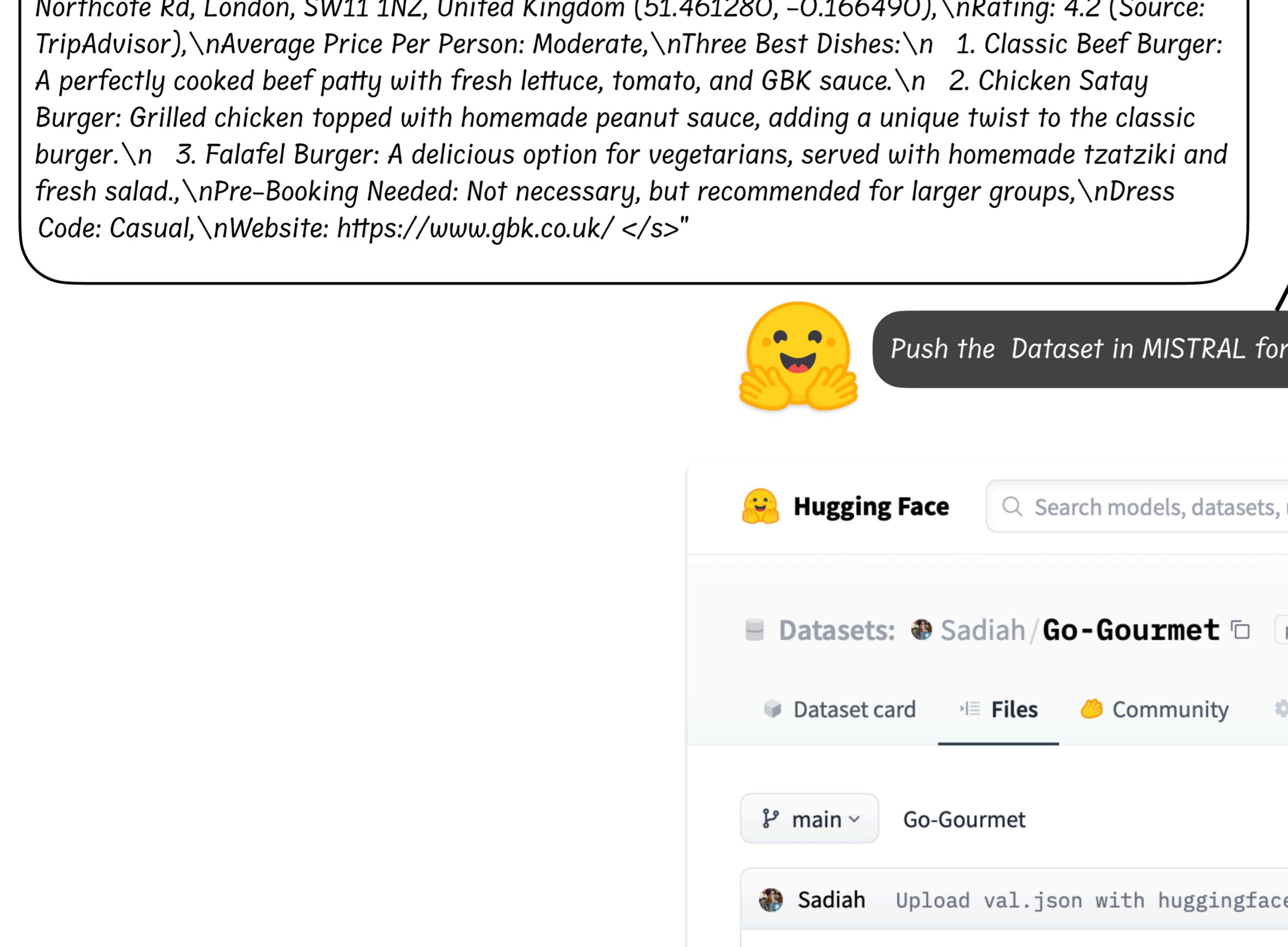
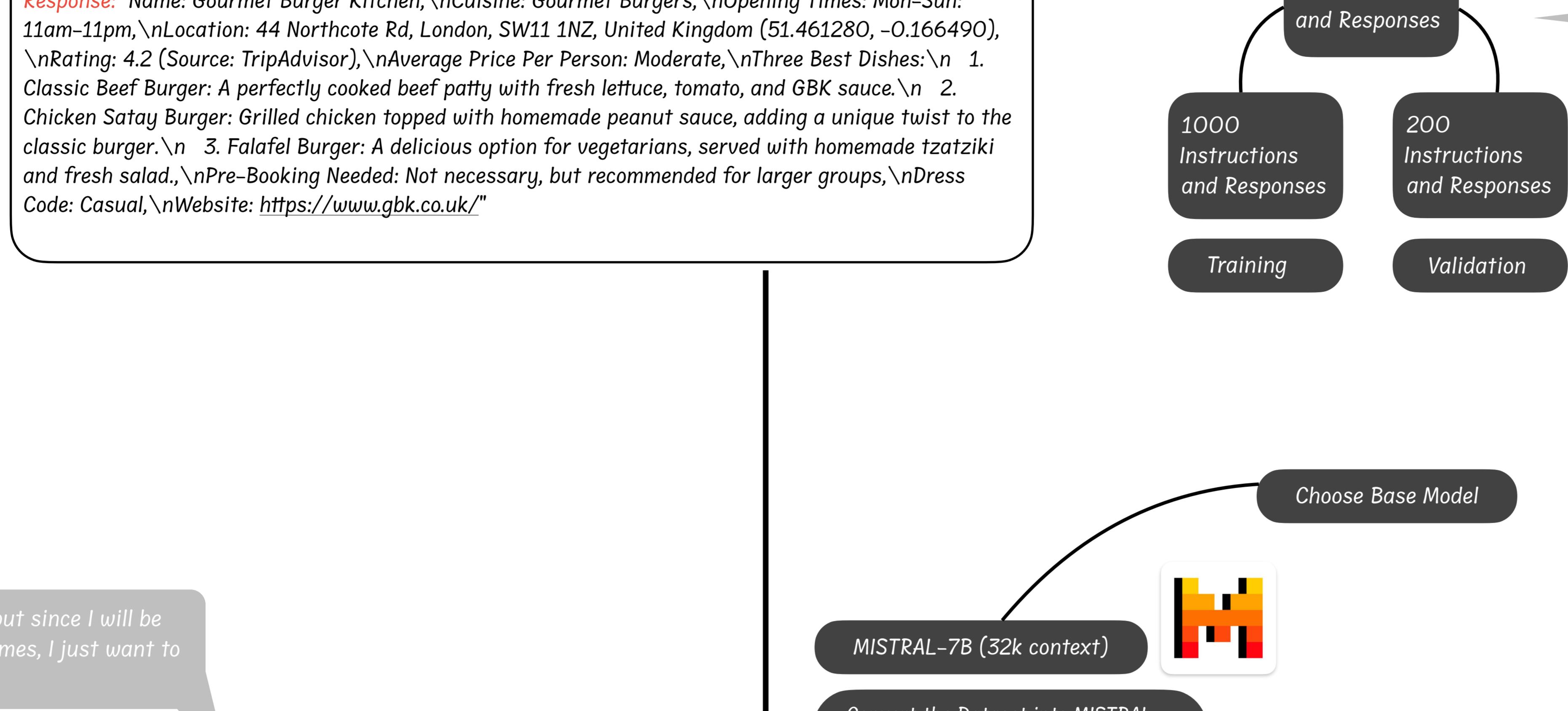


Go-Gourmet!

Fine-Tuning for uniform structure

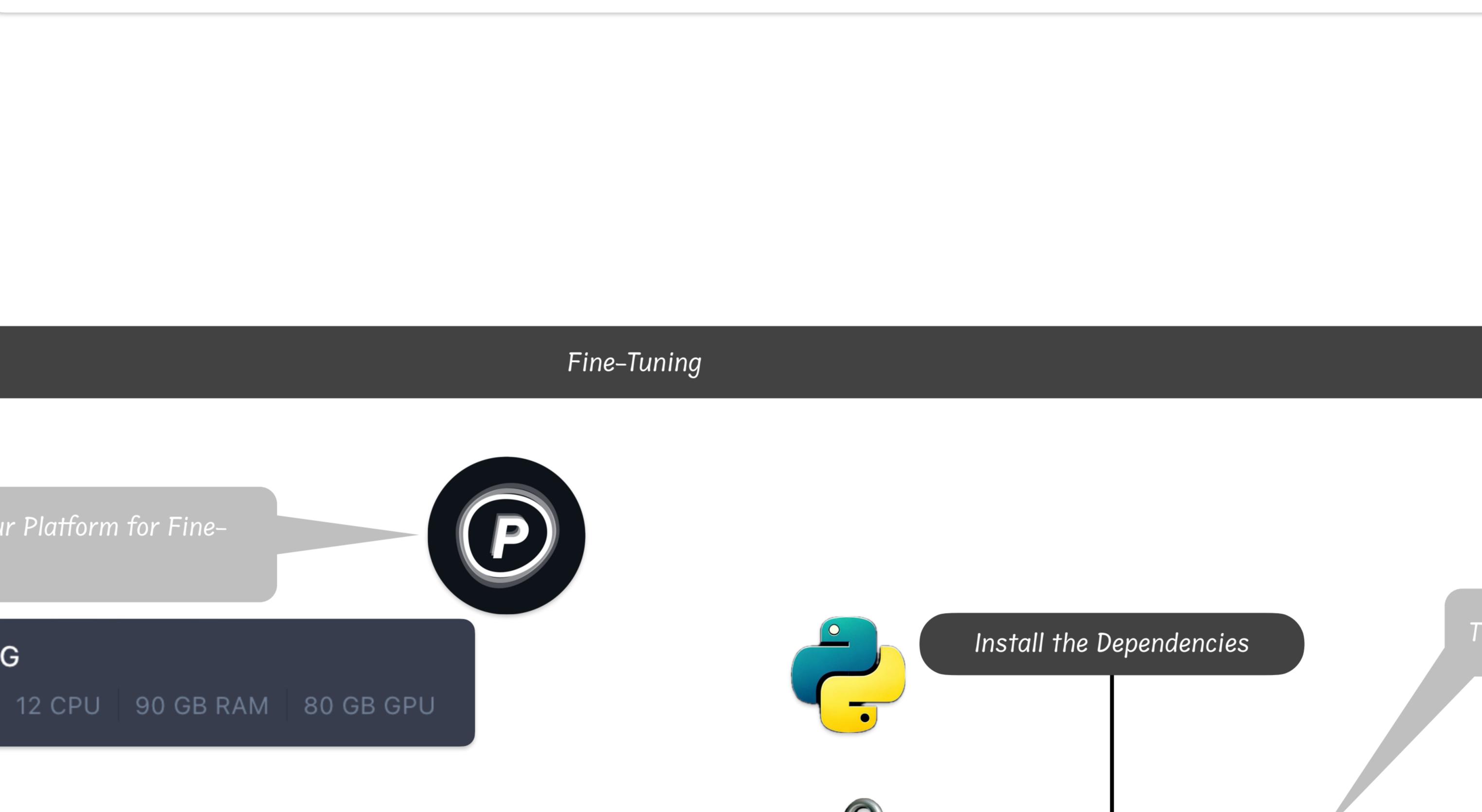
Copyright © 2024 by @sadiyahzahoor.
All rights reserved.

Data-Preparation



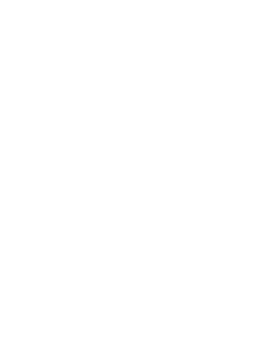
Push the Dataset in MISTRAL format to Hugging Face

This is how your Dataset should look like on Hugging Face.



Fine-Tuning

Choose your Platform for Fine-Tuning.



A100-80G
\$3.18/hr 12 CPU – 90 GB RAM – 80 GB GPU

Install the Dependencies



Pass on the credentials



This is your Hugging Face Token

Load the Base-Transformer model and its Tokeniser from Hugging Face with Quantization Configuration

mistral/Mistral-7B-Instruct-v0.2

load_in_4bit,...
bb_4bit_quant_type: ...
bnb_4bit_compute_dtype, ...
bnb_4bit_use_double_quant, ...
....add more

from transformers import AutoModelForCausalLM, AutoTokenizer

Load the base model

model = AutoModelForCausalLM.from_pretrained(BASE_MODEL, trust_remote_code=True, device_map="auto")

Load the entire model on the GPU

model.config.use_cache = False # Disable caching of model outputs

model.config.pretraining_tp = 1 # Tensor parallelism settings

Load tokenizer

tokenizer = AutoTokenizer.from_pretrained(BASE_MODEL, trust_remote_code=True)

tokenizer.padding_side = "right" # Sequences are left indented

Token has been saved to ./root/.cache/huggingface/token

Login successful

config.json: 100% 596/696 [00:00:00, 41.0KB/s]

model.safetensors.index.json: 100% 25.4/25.4K [00:00:00, 1.97MB/s]

Downloading shards: 0% 2/3 [00:41<00:20, 20.07s/R]

model-00001-of-00003.safetensors: 100% 4.94/4.94G [00:24<00:00, 320MB/s]

model-00002-of-00003.safetensors: 100% 5.00/5.00G [00:16<00:00, 305MB/s]

model-00003-of-00003.safetensors: 78% 3.53/4.54G [00:53<00:31, 32.4MB/s]

Load PEFT configuration for training

PEFT used: QLoRA

Load the Go-Gourmet Dataset from Hugging Face

Add Training Configuration

Batch Size: ...
Learning Rate: ...
Optimizer: ...
Epochs: ...
.....add more

Start training loop

Time taken - 2 hrs (approx)

Epoch Training Loss Validation Loss

0 1.990300 1.886203

1 1.578500 1.469587

2 1.193700 1.102770

4 0.907400 0.907790

5 0.864800 0.859449

6 0.859400 0.876061

7 0.853900 0.853939

9 0.737900 0.802890

10 0.767200 0.762369

12 0.746000 0.763486

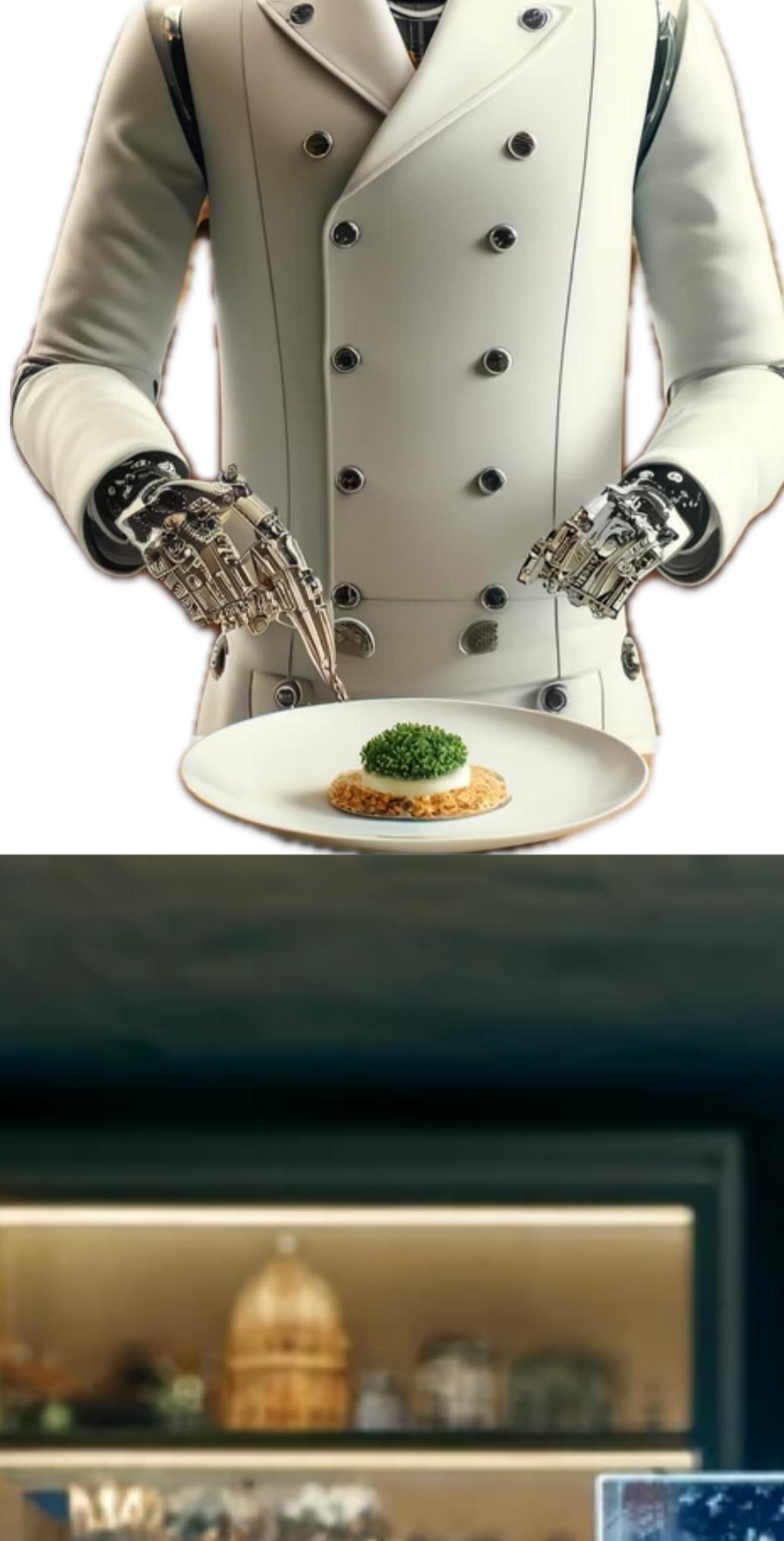
13 0.787600 0.758479

14 0.732400 0.754083

We will use Parameter Efficient Fine-Tuning which reduces the number of parameters we need to train.

This is the Dataset that we just pushed above to the Hugging Face.

This is the Dataset that we just pushed above to the Hugging Face.



Name: Duck and Waffles,
Cuisine: Modern British,
Opening Times: Mon-Sun: 24 Hours,
Location: 110 Borough High St, London SE1 1LB, United Kingdom
(51.5051, -0.0793),
Rating: 4.1 (Source: TripAdvisor),
Average Price Per Person: Moderate,
Three Best Dishes:
1. Duck and Waffles: Crispy duck leg confit served with a waffle and maple syrup.
2. Fried Chicken Wings: Spicy and crispy, served with a side of coleslaw.
3. Churros with Chocolate Sauce: A sweet treat to finish the meal.
Pre-Booking Needed: Recommended, especially for weekend brunch,
Dress Code: Casual,
Website: https://duckandwaffles.com/

Inference



Copyright © 2024 by @sadiyahzahoor.
All rights reserved.

Adapter model weights in techniques like LoRA (Low-Rank Adaptation) are generally incomplete compared to the full weights of a large language model (Base LLM).

So we must load both the base model and adapter but doing so separately requires loading both components during inference, which can add latency. Merging simplifies the model by combining everything into a single entity.

Here is a link to my trained Model on Hugging Face.

Merge Adapter with Base Model

Save & Push Merged Model & Tokenizer to Hub

Sadiyah /Go-Gourmet

huggingface.co

Sadiyah /Go-Gourmet - Hugging Face

huggingface.co