

Overall Goal:

The primary purpose of this notebook is to **load and test the model that was fine-tuned** in the previous notebook (File 2). It loads the base Qwen2.5-VL model, applies the saved fine-tuned LoRA adapters, runs inference on the same set of 200 processed images from the EXAMS-V dataset, extracts the predicted answers, and saves these predictions into a new file (output_fine_tune.jsonl). This allows for evaluating how well the fine-tuned model performs compared to the original model (from File 1).

Key Steps:

1. Setup & Environment:

- Installs unsloth again.
- Installs a *specific version* (v4.49.0) of the transformers library from GitHub using `--no-deps` (to avoid reinstalling other dependencies). This rollback might be necessary for compatibility reasons when loading the fine-tuned PEFT/LoRA model saved by the trl library in the previous notebook.
- Prints the transformers version to confirm it's 4.49.0.
- Confirms torch is installed.

2. Loading the Fine-Tuned Model:

- **Unzipping:** It expects a zip file named `vinayak_model.zip` to be present in `/content/`. It extracts the contents of this zip file. This zip file presumably contains the fine-tuned LoRA adapters and configuration saved in the previous notebook (which were originally saved to a directory named `./vinayak_model_new` or similar). The extracted folder is likely named `/content/vinayak_model`.
- **Loading Base Model:** It loads the *base* 4-bit quantized Qwen2.5-VL model from Unsloth (`unsloth/Qwen2.5-VL-7B-Instruct-unsloth-bnb-4bit`), similar to the previous notebooks.
- **Applying LoRA Adapters:** This is the crucial step. It uses `PeftModel.from_pretrained` to load the LoRA adapters from the extracted `/content/vinayak_model` directory and *apply them* to the already loaded base model. The model variable now represents the **fine-tuned** version.
- **Loading Processor:** It loads the standard processor from the *original* Qwen checkpoint (`Qwen/Qwen2.5-VL-7B-Instruct`) to handle input formatting correctly.

3. Data Preparation (Repetition):

- Similar to the previous notebooks, it loads the EXAMS-V dataset, likely from the locally saved `/content/dataset` directory created before.
- It **re-processes** the first 200 images: resizing them and saving them locally to `/content/downloaded_images` (overwriting if they already exist).

- It re-creates the results.json file mapping image paths to *ground truth* answers. Again, this ground truth seems mainly for reference during the image saving step, not directly used in the inference prompt in this notebook. **Note:** This time, sort_keys=True is used when saving results.json, ensuring the keys (image paths) are alphabetically ordered in the file.
- 4. **Memory Management:** Includes explicit calls to gc.collect() and torch.cuda.empty_cache() to free up RAM and GPU memory, which is good practice.
- 5. **Inference with Fine-Tuned Model:**
 - It defines the same parse_model_output and generate_reasoning helper functions as used in *File 1* (the original prediction notebook). The prompt used in generate_reasoning asks the model to *determine* the correct answer and provide reasoning, **not** relying on a provided ground truth answer.
 - It iterates through the 200 locally saved images.
 - For each image:
 - It calls generate_reasoning, passing the image path to the **fine-tuned model**.
 - It gets the model's raw text output (which should be a JSON string with reasoning and the predicted answer).
 - It parses this output using parse_model_output.
 - It extracts the correct_answer *predicted by the fine-tuned model*.
 - It saves a simple JSON object containing the image_path and the correct_answer (predicted by the fine-tuned model) to a new file: **output_fine_tune.jsonl**.
- 6. **Download Results:**
 - Finally, it uses Colab's files.download to download the output_fine_tune.jsonl file, which contains the predictions made by the *fine-tuned* model.

In Essence:

This notebook takes the specialized LoRA weights created during the fine-tuning process (File 2), applies them to the base VLM, and then runs inference on the test images *using this enhanced model*. The goal is to see if the fine-tuning improved the model's ability to correctly predict the answers compared to the original model's performance (results from File 1). The output file output_fine_tune.jsonl contains the predictions from this fine-tuned model. The cleared cell outputs indicate the script was run, but the visual logs were not saved in the notebook file.