**Overall Goal:**

The primary goal of this notebook is to create a **live, interactive web application (frontend)** using **Streamlit** to showcase the capabilities of the **fine-tuned Qwen2.5-VL model** (which was trained in File 2 and potentially tested in File 3). This app allows users to either select images from the pre-processed dataset or upload their own image, interact with the fine-tuned model to get explanations or predictions, and see the results directly in a web interface.

**How it Works - A Consolidated View:**

1. **Environment Setup & Model Loading:**

   o **Installs:** It first installs necessary libraries: unsloth (for efficient model handling), a specific version of transformers (v4.49.0, likely for compatibility with the saved fine-tuned model), streamlit (to build the web app), pyngrok (to make the Colab-hosted app publicly accessible), nest_asyncio (often needed for pyngrok in notebooks), and gdown (to download files from Google Drive).

   o **Downloads Fine-Tuned Model:** It uses gdown to download a vinayak_model.zip file from a Google Drive link.

   o **Unzips Model:** It unzips this file, extracting the fine-tuned LoRA adapter weights and associated configuration files into a directory likely named /content/vinayak_model/.

   o **Loads Base Model:** It loads the base 4-bit quantized Qwen2.5-VL model using unsloth.FastVisionModel.from_pretrained, just like in the previous notebooks.

   o **Applies Fine-Tuned Adapters:** It then uses PeftModel.from_pretrained to load the LoRA adapter weights from the unzipped /content/vinayak_model directory and merge them into the base model. The model variable now holds the **fine-tuned** version.

   o **Loads Processor:** It loads the standard processor for the Qwen model to handle image/text inputs.

2. **Data Preparation (Repetition):**

   o It re-loads the EXAMS-V dataset and re-processes the first 200 images (resizing, saving locally). This is likely done to ensure the images are available for the Streamlit app to display when running in "Dataset Mode". It also regenerates the results.json containing ground truth answers for these images.

3. **Creating the Streamlit App (app.py):**

   o The core of this notebook is a large Python code block that writes a complete Streamlit application script to a file named app.py.

   o **Key Features of app.py:**

      ▪ **Caching:** Uses Streamlit's @st.cache_resource to load the model/processor only once and @st.cache_data to load the dataset's ground truth answers (results.json) once, significantly speeding up interactions after the initial load.

- **Mode Selection:** Provides radio buttons for the user to choose between "Use Image from Dataset" or "Upload Your Own Image".
- **Dataset Mode:**
    - If selected, it loads the image paths and ground truth answers from results.json.
    - It displays one image at a time.
    - Provides "Previous", "Next", and "Random" buttons for navigation.
    - Includes a dropdown (st.selectbox) to jump directly to a specific image.
    - When the user selects an answer (A/B/C/D):
        - It compares the user's choice to the known correct answer from results.json.
        - It provides feedback (Correct/Incorrect).
        - If the user is *incorrect*, it calls the generate_explanation_for_correct_answer function. This function prompts the *fine-tuned model* to explain *why the actual correct answer is right*.
- **Upload Mode:**
    - Provides a file uploader (st.file_uploader) for the user to select their own image.
    - When the user selects an answer (A/B/C/D):
        - It calls the get_model_answer_and_explanation_upload function. This function prompts the *fine-tuned model* to both *determine the most likely correct answer* and *provide its reasoning*.
        - It compares the user's choice to the *model's predicted answer*.
        - It provides feedback based on whether the user agreed with the model.
        - It *always* displays the model's generated reasoning for its chosen answer.
- **State Management:** Uses st.session_state extensively to keep track of the current mode, image, user choices, results, etc., across user interactions and app reruns.
- **UI Layout:** Uses st.columns to arrange the image, options, and results neatly.

4. **Running the App and Tunneling:**

- o **Kill Previous Processes:** It tries to find and kill any existing processes using the default Streamlit port (8501) to ensure a clean start.

- o **Start Streamlit:** It runs the app.py script using streamlit run in the background (nohup) and logs output to streamlit_log.txt.

- o **Start Ngrok:** After a delay (to allow Streamlit and the model to load), it uses pyngrok to create a secure public URL that tunnels traffic to the Streamlit application running locally inside the Colab environment.

- o **Display URL:** It prints the public ngrok URL, which the user can click to access the live web application from their browser.

**In Summary:**

This notebook takes the fine-tuned model from the previous step, packages it into an interactive Streamlit web application (app.py), and then launches that application within the Colab environment, making it accessible via a public ngrok URL. The web app provides a user-friendly interface ("Frontend") to test the fine-tuned model either on images from the original dataset (with ground truth comparison and model explanations) or on user-uploaded images (where the model predicts and explains). It essentially deploys the fine-tuned model for live interaction and demonstration.