# Steps to Start This project Locally

## Step 1: Set Up Environment

Install Python: Make sure you have Python installed on your system. You can download it from the official Python website.

Create a Virtual Environment (Optional): It's good practice to create a virtual environment for your project to isolate dependencies. You can create one using virtualenv or venv module:

```
# Using virtualenv
pip install virtualenv
virtualenv venv
source venv/bin/activate


# Using venv (Python 3)
python3 -m venv venv
source venv/bin/activate
```

## Step 2: Install Dependencies

Navigate to your project directory in the terminal and install the required Python dependencies:

pip install flask torch torchvision transformers Pillow

## Step 3: Prepare Model Files

Ensure that you have the model files accessible. For the ResNet-50 and Faster R-CNN models, you don't need to download any additional files as they are provided by torchvision. However, you may need to download the tokenizer and model files for the VisionEncoderDecoderModel and ViTFeatureExtractor.

You can download the required model files and tokenizer from the Hugging Face model hub (VisionEncoderDecoderModel and ViTFeatureExtractor).

Step 4: Project Structure

```
project_folder/
|
├── static/
|     └── index.html
|
├── final.py
└── imagenet-simple-labels.json
```

final.py, you just need to run it with Python:

python final.py

After executing final.py, the Flask server starts, directing you to a webpage stored in index.html. Clients have the option to upload files directly to the server or use the "Choose File" button to select an image they wish to analyze for object detection. Upon selecting an image, clients can click "Detect Object". This action triggers the server to classify the object, generate captions, and mark grid lines around the detected objects for easy identification.

## Detailed Explanation of the Solution

The Flask application (final.py) provides routes for handling image uploads, object detection, classification, and captioning.

Upon uploading an image, the application preprocesses it and passes it through a pre-trained Faster R-CNN model for object detection. Detected objects are marked with bounding boxes on the image.

The application then uses a pre-trained ResNet-50 model for image classification to classify the entire image.

Additionally, the application utilizes a VisionEncoderDecoder model to generate captions for the uploaded images.

The detected objects, along with their classifications and captions, are overlaid onto the original image, providing users with a comprehensive analysis of the uploaded image.

Conclusion

This project demonstrates the integration of object detection, classification, and captioning functionalities into a web-based application using Flask. Users can easily upload images, analyze them for objects, classify them, and receive captions, enhancing their understanding and interpretation of the images.

# Alternative Approaches

1. Single Model Approach:

Instead of using separate models for object detection, classification, and captioning, a single model could be trained to perform all tasks simultaneously. This approach reduces complexity and inference time but may require more extensive training and tuning.

2. Different Object Detection Models:

While Faster R-CNN is used in this project, other object detection models like YOLO (You Only Look Once) or SSD (Single Shot MultiBox Detector) could be employed. These models trade-off between accuracy and speed, with YOLO being faster but potentially less accurate than Faster R-CNN.

3. Fine-tuning Pre-trained Models:

Rather than using pre-trained models out-of-the-box, fine-tuning them on specific datasets relevant to the application domain could improve performance. Fine-tuning allows models to adapt to specific object classes or image characteristics present in the target dataset.
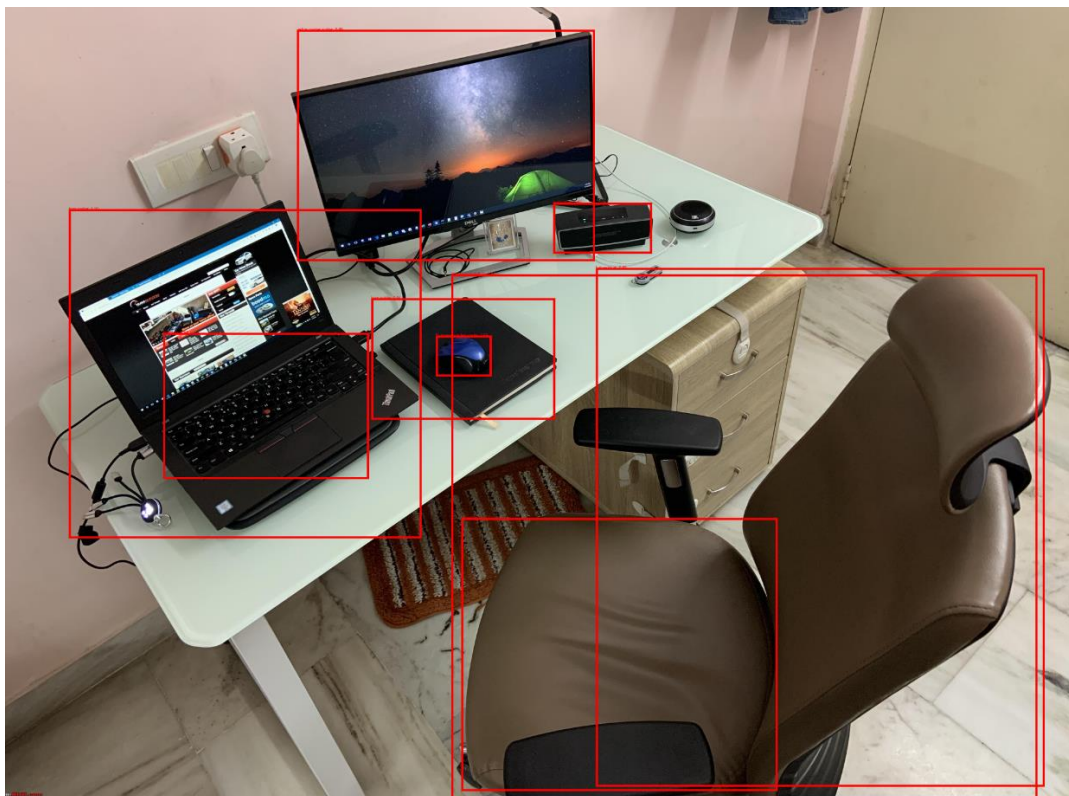
While the chosen approach in this project utilizes a combination of pre-trained models for object detection, classification, and captioning, alternative approaches exist that offer different trade-offs in terms of accuracy, speed, resource requirements, and domain specificity. Understanding these alternatives and their implications is essential for selecting the most suitable approach for a given application.

# Images with output visualizations (boxes, class labels, image captions)

## Input Image:



## Output Image: (boxes)

**Output Image (Class Label & , image captions)**

# Object Detection
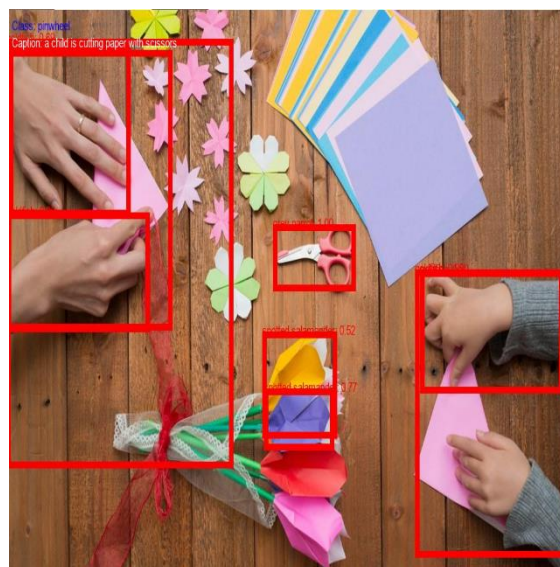
Choose File | image-20.jpg          Detect Objects

Class: desktop computer
Caption: a desk with a laptop computer and a monitor

**(This is above at the image at the top left corner)**

**Input Image: 2**



**Output Image: (boxes)**

# Object Detection

Choose File | image-5.jpg          Detect Objects

Class: pinwheel

Caption: a child is cutting paper with scissors

**(This is above at the image at the top left corner)**