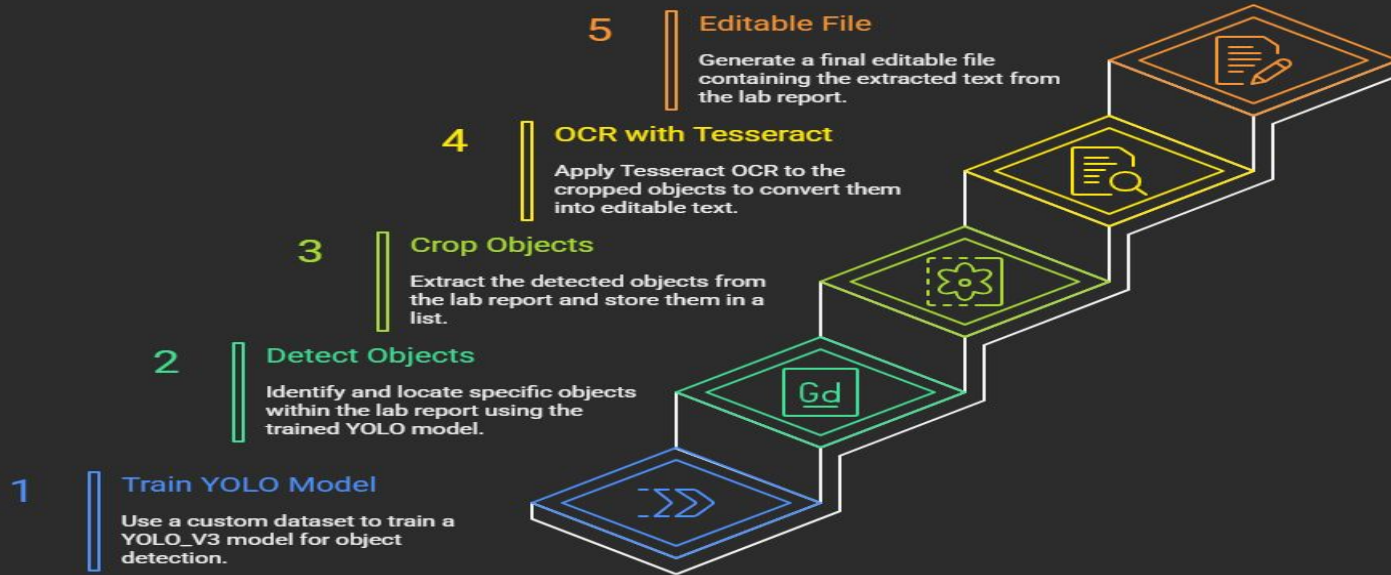# Custom-Object Character Recognition(OCR) on AWS (Google Drive/ Cloud Storage)
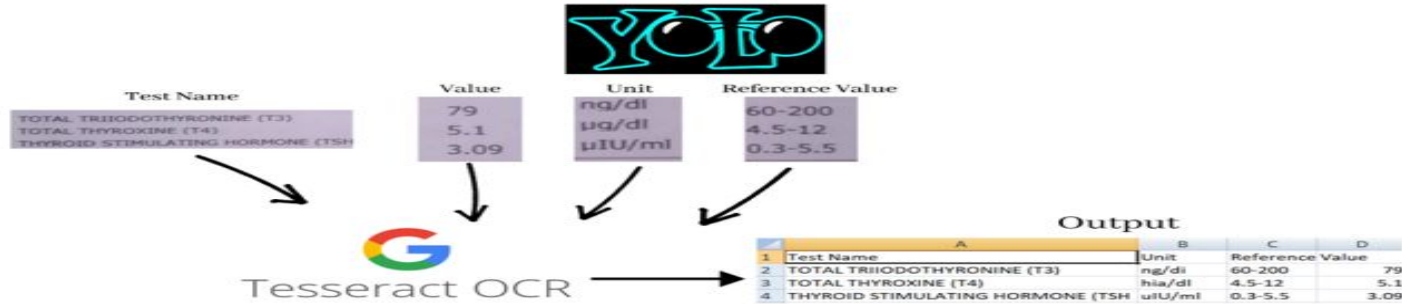
- Next Hikes IT Solutions
- Project_10:OCR on Google drive
- Prepared by: Kalavathi Alegapalli
- Date:10/11/2025

📄 **Structured Medical Data Extraction Using OCR**

- **Efficient Digitization**: Tesseract OCR accurately extracts test names, values, units, and reference ranges from medical reports, enabling automated data entry.
- **Structured Output**: Extracted information is organized into spreadsheet columns for streamlined analysis and record-keeping.
- **Clinical Utility**: This process supports faster diagnostics, trend tracking, and integration into electronic health systems.

# Building a Custom OCR System

# Workflow Overview

**Step 1:** Label images using LabelImg (Test Name, Value, Units, Reference Range).

**Step 2:** Split the dataset into training and testing sets.

**Step 3:** Train YOLO_V5 model in Google Colab for object detection.

**Step 4:** Use YOLO_V5 to detect key regions in lab reports.

**Step 5:** Extract text using Tesseract OCR.

**Step 6:** Build a Streamlit app for user interaction and testing.

# STEP 1 & 2 – IMAGE LABELING & DATA PREPARATION

- **Tool Used:** LabelImg
  → Labeled lab report images in data_images folder with bounding boxes for **Test Name, Value, Units, Reference Range** to train YOLO_V5.
- **XML Parsing & Structuring:**
  → Imported libraries for file handling and XML parsing
  → Extracted image paths and parsed XML for filenames, dimensions, and bounding boxes
  → Combined into a structured **pandas DataFrame**
- **Coordinate Normalization:**
  → Computed normalized center, width, and height for YOLO format
- **Dataset Splitting & Label Encoding:**
  → Split into **80% training / 20% testing** by unique filenames
  → Encoded object labels (Test Name, Value, Units, Reference Range)
- **Directory Setup:**
  → Created train/test folders
  → Moved images and saved labels in .txt format

# 3. Model Training:

Model trained in Google Colab using YOLOv5 with pretrained weights, 640px images, batch size 12, over 200 epochs.
 Achieved high precision, recall, and mAP@0.50–0.95, indicating strong, generalizable object detection performance.

```
Model summary: 157 layers, 7020913 parameters, 0 gradients, 15.8 GFLOPs
                Class     Images   Instances          P          R     mAP50   mAP50-95:
                  all         10          40      0.989          1     0.995       0.733
            Test Name         10          10      0.995          1     0.995       0.851
                Value         10          10      0.996          1     0.995       0.629
                Units         10          10      0.984          1     0.995       0.633
      Reference Range         10          10       0.98          1     0.995       0.821
```

# Step4:Object Detection



- The script loads YOLO from an ONNX file, preprocesses and resizes the image to 640x640, runs inference, applies NMS, and draws labeled bounding boxes.
-  The final annotated image is saved for visualizing detected objects.

# Step5:Text Extraction

# Conclusion

- The pipeline efficiently detects and annotates objects using a YOLO model loaded from an ONNX file.
- Preprocessing ensures input images are resized and padded to meet model requirements (640×640).
- Inference and Non-Maximum Suppression (NMS) refine predictions by filtering low-confidence and overlapping detections.
- Bounding boxes with class labels and confidence scores are drawn for clear visualization.
- The final annotated image is saved, confirming successful end-to-end object detection and rendering.

# References:

1. **Ultralytics YOLOv5 Official Documentation**
   *Source:* Ultralytics. "Train Custom Data with YOLOv5."
   *URL:* https://docs.ultralytics.com/yolov5/tutorials/train_custom_data/
   *Description:* Official guide on preparing datasets for YOLOv5 training, including label formats, folder structure, and common dataset-related errors.

2. **Ultralytics Dataset Format Guide**
   *Source:* Ultralytics. "Dataset Structure and Label Formats."
   *URL:* https://docs.ultralytics.com/datasets/detect/
   *Description:* Explains the correct directory and label structure required for YOLOv5 and YOLOv8 datasets.

3. **YOLOv5 GitHub Repository**
   *Source:* Ultralytics GitHub Repository.
   *URL:* https://github.com/ultralytics/yolov5
   *Description:* The open-source repository for YOLOv5, including training scripts, dataset requirements, and issue discussions related to label errors.

4. **Stack Overflow Discussion: "YOLOv5 no labels found error"**
   *Source:* Stack Overflow.
   *URL:* https://stackoverflow.com/questions/73114650/yolov5-no-labels-found-error
   *Description:* Community discussion confirming that missing or incorrectly structured label files cause this error, with examples of fixes.