

# End-to-End OCR Project with PaddleOCR

## Abstract

This project documents a complete workflow using **PaddleOCR** for Optical Character Recognition (OCR). It covers the PP-OCR system architecture, dataset preparation and conversion, reproducible training pipelines, and runnable notebooks on Colab and Kaggle. The project demonstrates the full OCR process: text detection, angle classification, and recognition, trained on real-world datasets and evaluated with standard metrics.

---

## 1. Introduction

Optical Character Recognition (OCR) is central to many computer vision tasks such as digitizing documents, recognizing text in natural scenes, and enabling multilingual AI assistants.

**PaddleOCR** is a modular open-source framework that implements **PP-OCR**, a three-stage pipeline:

- **Text Detection** (finding text regions in images)
- **Angle Classification** (normalizing text orientation)
- **Text Recognition** (transcribing detected text into characters)

This project provides an end-to-end study of PaddleOCR, from understanding its evolution (v3  $\rightarrow$  v5) to training and evaluating models on benchmark datasets.

---

## 2. Architecture Review

### 2.1 Core Components

- **Text Detector (DBNet-based):**  
Backbone networks such as MobileNetV3 or ResNet are used. The detector outputs polygonal regions of text.
- **Angle Classifier:**  
A lightweight CNN ensures rotated text is normalized before recognition.
- **Recognizer (CRNN):**  
A Convolutional Recurrent Neural Network with CTC (Connectionist Temporal Classification) loss decodes character sequences.

### 2.2 Evolution of PP-OCR

- **PP-OCRv3 (2022):** Stronger backbones, better augmentations, broader multilingual coverage.
- **PP-OCRv4 (2023):** Smaller and faster lightweight models.



- **PP-OCRv5 (2025):** Higher recognition accuracy, especially for English, Thai, and Greek; improved latency and better multilingual generalization.

*(Insert Diagram: End-to-end PP-OCR pipeline: Detector → Angle Classifier → Recognizer)*

---

### 3. Datasets

#### 3.1 Selected Sources

- **COCO-Text V2.0:** 63,686 images, 239,506 annotated text instances.
- **ICDAR 2015:** Focused on incidental scene text.
- **ICDAR 2019 MLT:** Large multilingual dataset with over 10 scripts.
- **Additional Sources:** LSVT, RCTW-17, MTWI for extended script coverage.

#### 3.2 Formatting for PaddleOCR

- **Detection Labels:** Polygon coordinates + transcription, similar to ICDAR conventions.
  - **Recognition Labels:** Cropped word images with a mapping file.
  - **Tools:** PPOCRLabel was used to convert datasets into PP-OCR formats.
- 

### 4. Training Pipeline

#### 4.1 Setup

- Framework: PaddleOCR on Colab and Kaggle with GPU support.
- Configurations:
  - **Lightweight (MobileNetV3):** Quick experiments, lower resource usage.
  - **Server/Full models:** Higher accuracy, larger computational cost.

#### 4.2 Procedure

- **Detection Training:** PP-OCRv3/v5 detector trained for 10–20 epochs on COCO-Text subset. Metrics logged: Precision, Recall, F-measure.
- **Recognition Training:** CRNN recognizer trained on ICDAR 2019 cropped words. Metrics: Accuracy, Normalized Edit Distance.
- **Optimization:** Mixed precision enabled, Adam optimizer, cosine scheduler, data augmentations (random crop, rotation, blur).

#### 4.3 Evaluation Results (Baseline)

- Detection F-measure: ~0.75 on validation subset.
  - Recognition Accuracy: ~82% on validation set.
-

## 5. Notebooks (Colab and Kaggle)

The project includes runnable notebooks for both Colab and Kaggle. Each notebook contains:

1. Environment setup and PaddleOCR installation.
2. Dataset download and conversion.
3. Training launches with logging (TensorBoard/VisualDL).
4. Visualization of predictions (detected boxes and recognized text).
5. Automatic saving of trained weights and metrics.

---

## 6. Key Insights and Challenges

- **Efficiency:** PP-OCRv5 delivers faster inference without accuracy loss.
- **Multilingual Complexity:** Preparing datasets across multiple scripts is challenging; PPOCRLabel simplifies formatting.
- **Deployment Trade-offs:** Lightweight models offer speed, while larger models provide better accuracy.
- **Resource Constraints:** Large multilingual training requires careful GPU memory management.

---

## 7. Conclusion

This project demonstrates a reproducible, end-to-end OCR workflow with PaddleOCR. From architecture study to dataset formatting, training, and evaluation, the system achieved strong performance and highlighted the trade-offs between speed, accuracy, and multilingual coverage.

---

## References

- PaddleOCR GitHub: <https://github.com/PaddlePaddle/PaddleOCR>
- PaddleOCR Documentation:  
<https://paddlepaddle.github.io/PaddleOCR/main/en/index.html>
- COCO-Text V2.0 Dataset
- ICDAR 2015 and ICDAR 2019 MLT Datasets

---

🔗 This Word-style report mirrors the LaTeX one, just formatted with **headings and lists** instead of LaTeX commands. You can insert:

- **Figures/diagrams** in the "Architecture Review" and "Notebooks" sections.



- **Screenshots** of training logs, sample predictions, or Colab outputs in the "Training Pipeline" section.