

I M A G E

P D F

C V 2

J S O N

N E R

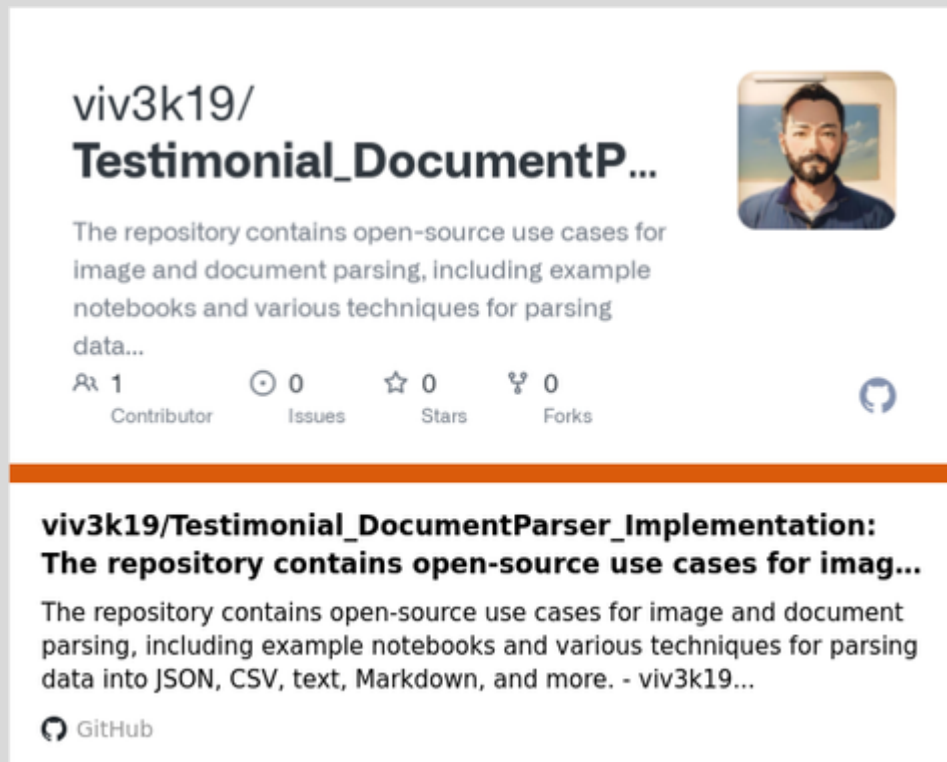
O C R

D o c u m e n t

P a r s e r

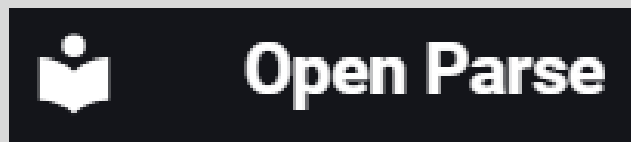
I m p l e m e n t a t i o n

Project Link



click on the image or [here](#)

Research & Experiments



Particular Documentation
for each:
click on the logos

How to Use

- Go to my GitHub repository.
- Select the implementation you want.
- Locate the Colab notebook.
- Download the notebook.
- Upload the notebook to Google Colab.
- Run the Colab notebook.
- The dataset is provided in the repository.
- Analyze the Results.

Edgcases

- Varied Layouts
- Noisy Backgrounds
- Mixed Font Styles and Sizes
- Low-Quality Scans
- Handwritten Text
- Complex Tables
- Language and Character Set Variability
- Embedded Graphics and Charts
- Incorrect Document Orientation
- Inconsistent Terminology
- Variable Page Sizes
- Missing Data
- Dynamic Content
- Noise in Input Data
- Overlapping Text

Testcases

- Layout - blockwise extraction
- Metadata extraction
- RegEx based extraction
- OCR based extraction
- API based extraction
- Pattern based extraction
- NER based extraction
- LLM based extraction
- predefined python based extraction

Note: Identified and addressed within notebooks during research and experimentation.

I M A G E

P D F

C V 2

J S O N

N E R

O C R

T H A N K

Y O U