Pdf

Pdf2Img Img

Pdf → Img1 Img2 ... → Detectron model (layout parser) → finds deff elements like img, text, table
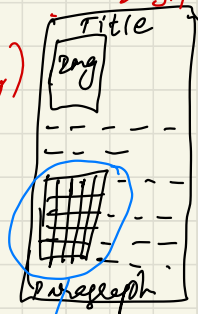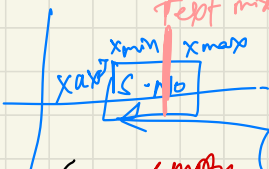
Title
Img
Paragraph

table detected

Use GCN to understand text in image

| 5 | NO | Naba | UV |
| 25 | 76 | 4P | 50 |
| 90 | 80 | 70 | 60 |
| 98 | 4 | 09 | 01 |

collect text

organize as table format

Text middle (Avg)

Xmin   Xmax
Xavg | S-No |

using OCR

| C.N | WO | From | UV |
| 5 | 10 | 15 | 20 |
| 25 | 90 | 60 | 90 |
| 60 | 20 | 65 | 90 |

empty dict to store text blocks

Using range t sep similar text blocks in horizontal way

collect text annotation

"if present in column dict
ok (append to list)
else : new entry

+ Remove columns with shorter length (for now not considering

CSV File

① drawbacks :- complex tables ?

② generalise. more

③ missing data (notable to capture all test from table)