



Développement d'un programme pour l'Extraction de Données depuis un Tableau à partir d'Images et de PDF

Soutenu par :

Benguerine Wassim
LAMHATTAT Redouane

Sommaire

1

Introduction

2

Contexte général du Projet

3

Analyse & Conception


4

Réalisation

5

Conclusion

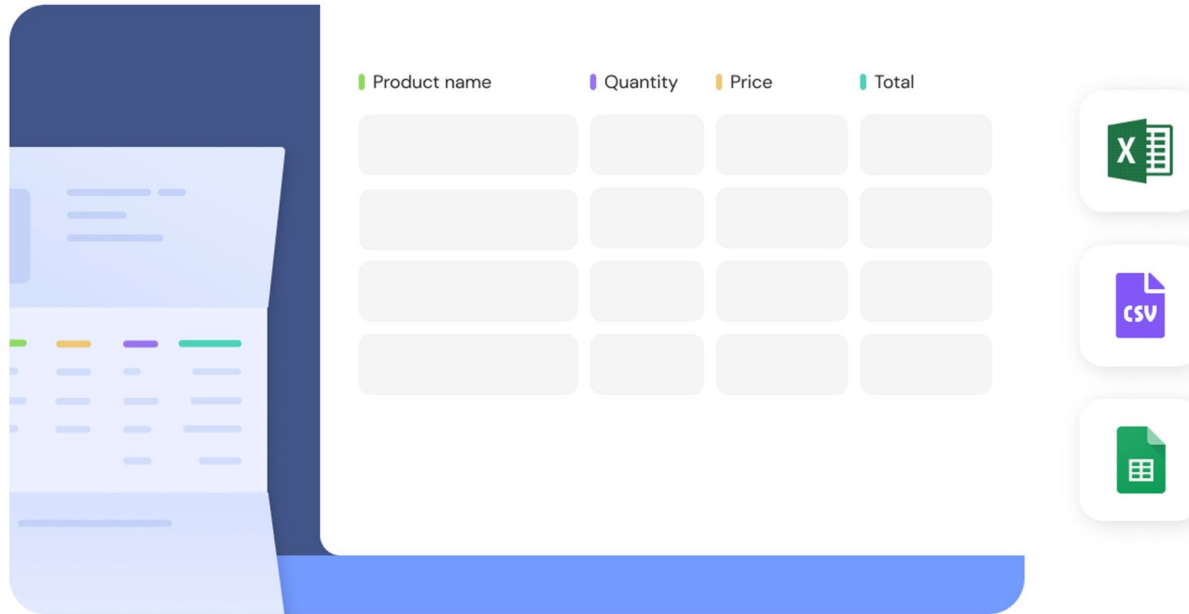
Qu'est-ce qu'une application d'extraction de données ?



Contexte général du Projet

Présentation du projet

Cadrage du projet



Objectifs du notre application



Réaliser une application qui permet de convertir un tableau présent dans une image en un format Excel ou CSV

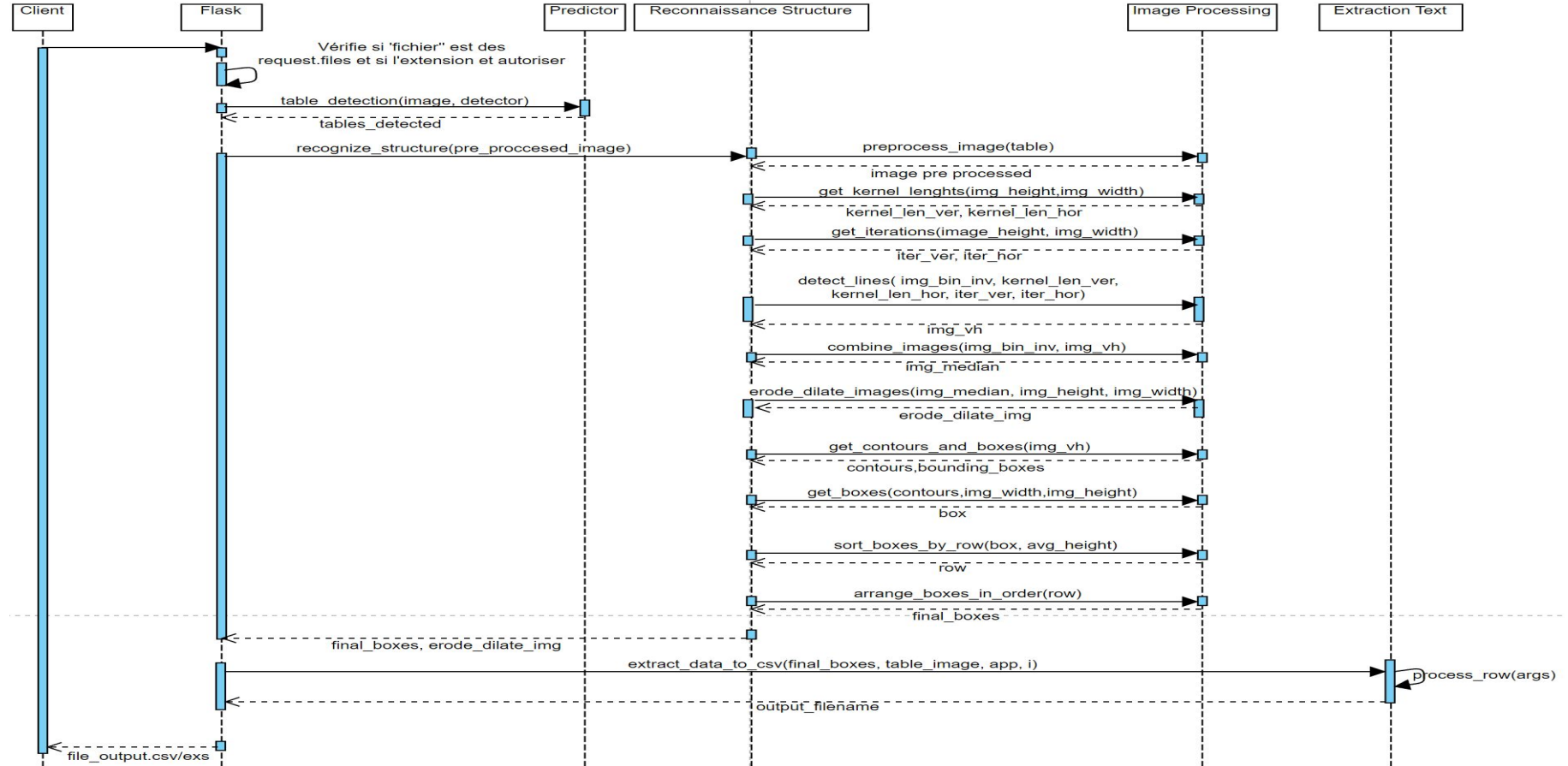
Analyse & Conception

Étude fonctionnelle

Étude technique

Étude fonctionnelle

Étude technique



- Capture des besoins techniques



Detectron2

OpenCV

PaddleOCR



Algorithmes Utilisés

Présentation des Algorithmes

Fonctionnement des Algorithmes



Algorithm: `get_contours_and_boxes`

Input: `img_vh`, a binary image

Output: Contours et boîtes de l'image

Procedure `get_contours_and_boxes(img_vh)`

`img_vh` = `add_border`

`img_vh` contours, = `findContours`

`img_vh` contours, bounding_boxes = `sort_contours` contours

return contours, bounding_boxes

Algorithm: `get_boxes` procedure

Input: contours, image width (`img_w`), image height (`img_h`)

Output: Les boxes satisfaisant les conditions [éviter les boxes trop grandes]

Procedure `get_boxes(contours, img_w, img_h)`

 boxes = []

for each *contour* `c` **in** `contours` **do**

`x`, `y`, `w`, `h` = `cv2.boundingRect(c)`

if `w` < `0.9*img_w` **and** `h` < `0.9*img_h` **then**

 boxes.append({'x': `x`, 'y': `y`, 'width': `w`, 'height': `h`})

end

end

return boxes

➤ Présentation des algorithmes et leur fonctionnement

Algorithm: `arrange_boxes_in_order` procedure

Input: Row of boxes

Output: Boxes organiser

Procedure `arrange_boxes_in_order`(*row*)

`max_columns` \leftarrow 0

`selected_index` \leftarrow 0

for *i* **in** `range(0, length_of(row))` **do**

if `length_of(row[i])` > `max_columns` **then**

`max_columns` \leftarrow `length_of(row[i])`

`selected_index` \leftarrow *i*

end

end

`selected_row` \leftarrow `row[selected_index]`

`centers` \leftarrow `empty_list`

for *j* **in** `range(0, length_of(selected_row))` **do**

`center` = `int(selected_row[j]['x'] + selected_row[j]['width']/2)`

`centers.append(center)`

end

`centers` \leftarrow `np.array(centers)`

`centers.sort()`

`organized_boxes` \leftarrow `empty_list`

for *i* **in** `range(0, length_of(row))` **do**

`temp_list` \leftarrow `create_empty_list_with_length(max_columns)`

for *j* **in** `range(0, length_of(row[i]))` **do**

`distance` \leftarrow `calculate_distance_from_centers(centers, row[i][j])`

`min_distance` \leftarrow `minimum_of(distance)`

`index` \leftarrow `get_index_of(min_distance in distance)`

`append row[i][j] to temp_list[index]`

end

`append temp_list to organized_boxes`

end

return `organized_boxes`

Algorithm: `extract_data_to_xls` procedure

Input: `finalboxes`, `img`, `app`, `table_index`

Output: `output_filename`

Procedure `extract_data_to_xls`(*finalboxes*, *img*, *app*, *table_index*)

`extracted_data` \leftarrow `empty_list`

With `ProcessPoolExecutor()` **as** `executor` **do**

for *row* **in** `finalboxes` **do**

`futures` \leftarrow [`executor.submit(process_cells, cell, img)` **for** `cell` **in** `row` **if** `cell`

`row_data` \leftarrow [`f.result()` **for** `f` **in** `futures` **if** `f.result()` **is not** `None`]

`append row_data to extracted_data`

end

`dataframe` \leftarrow `pd.DataFrame(extracted_data)`

`styled_dataframe` \leftarrow `dataframe.style.set_properties(align="left")`

`output_filename` \leftarrow `f"output_table_index.xlsx"`

`output_filepath` \leftarrow `os.path.join(app.config['UPLOAD_FOLDER'], output_filename)`

`styled_dataframe.to_excel(output_filepath)`

return `output_filename`

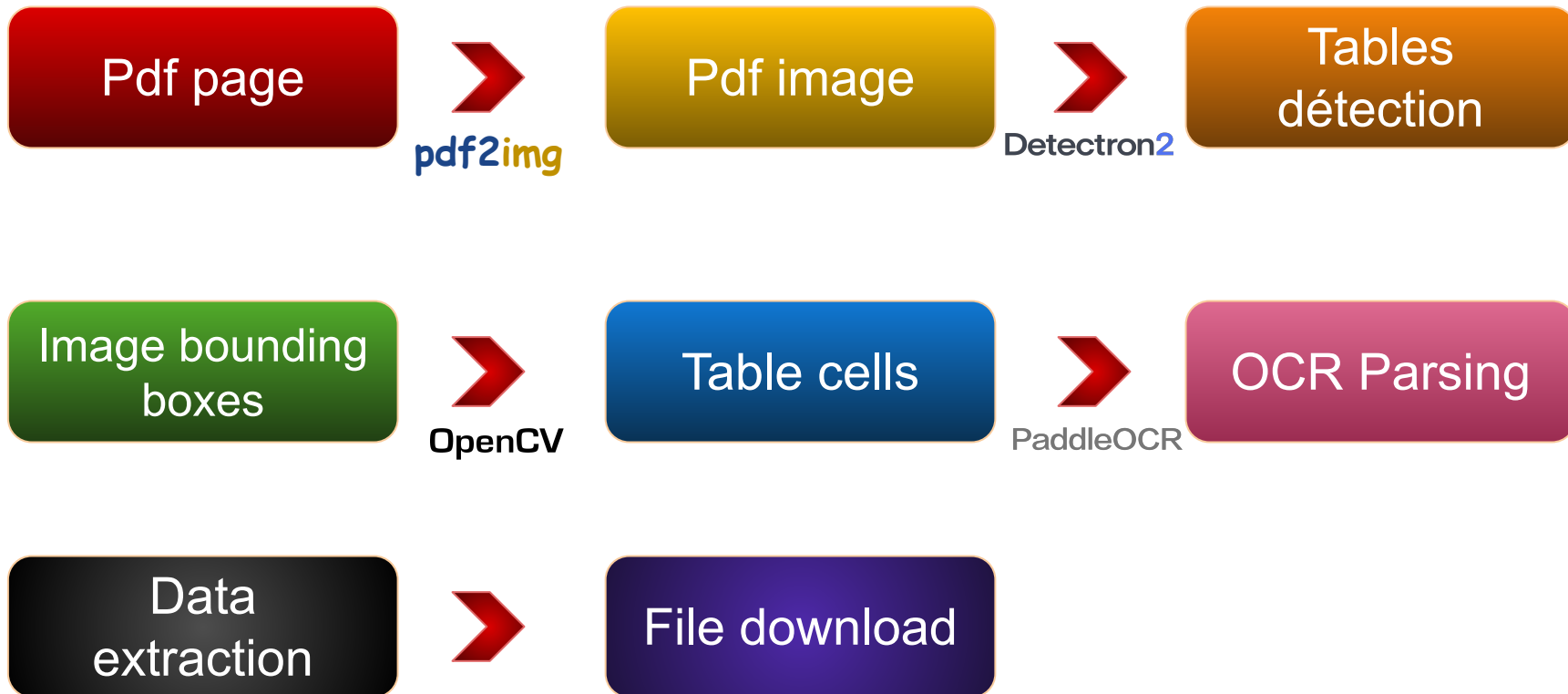


Réalisation

Architecture utilisée

Aperçu du projet

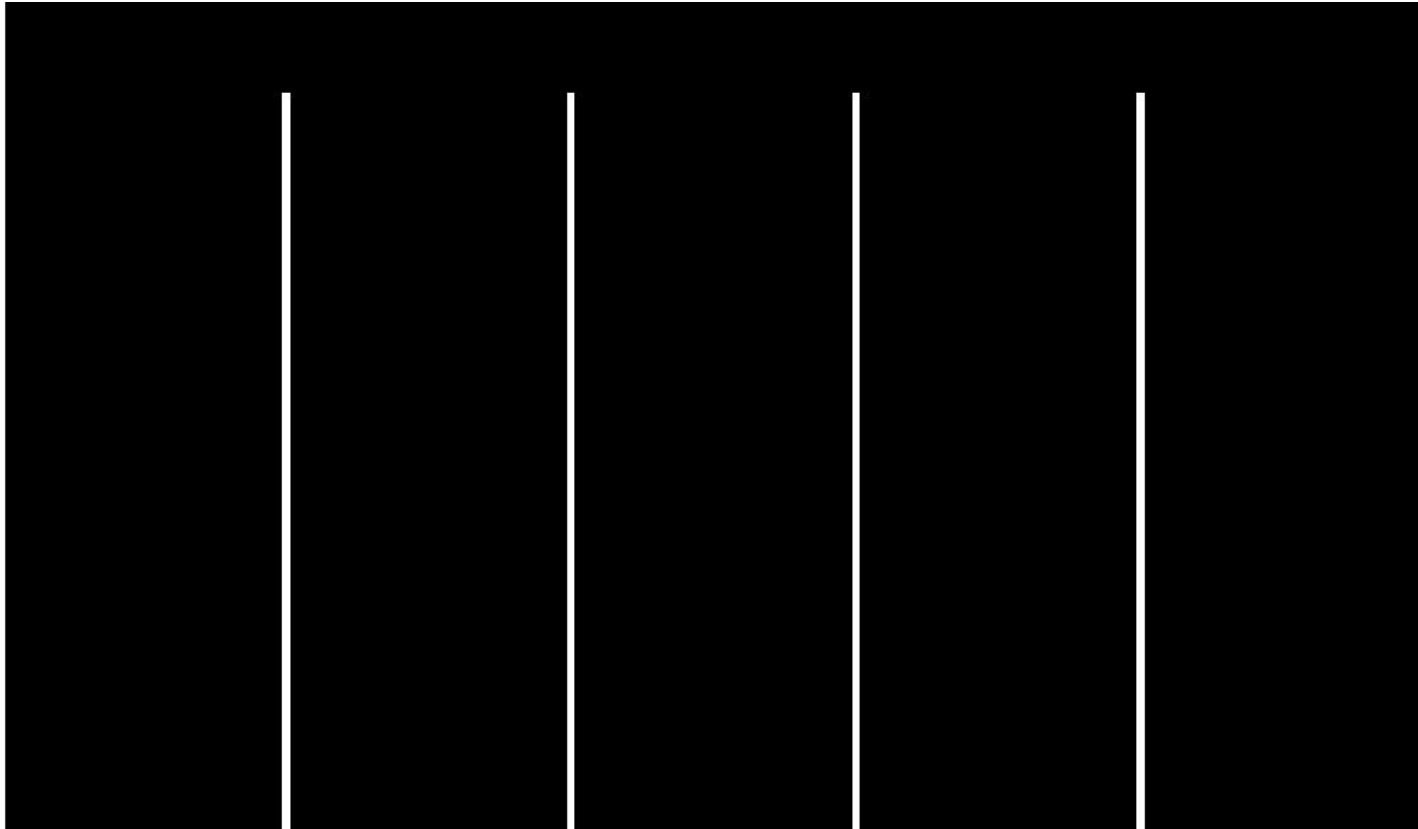
Tests



Order Date	Order ID	Salesperson	Units	Order Amount
1/1/2015	10458	Bossk	13	\$1,400.00
2/1/2015	10459	Dengar	14	\$2,555.95
2/2/2015	10460	Dengar	7	\$344.00
3/1/2015	10461	Bossk	8	\$857.00
3/3/2015	10462	Greedo	9	\$111.55
3/5/2015	10463	Fett	18	\$823.00
4/1/2015	10464	Bossk	17	\$24,455.50
4/28/2015	10465	Fett	13	\$213.00
4/30/2015	10466	Fett	1	\$10.00
5/1/2015	10467	Bossk	11	\$789.70
5/2/2015	10468	Dengar	25	\$21,286.60
5/3/2015	10469	Fett	8	\$1,285.00
6/1/2015	10470	Bossk	11	\$201.00
7/1/2015	10471	Bossk	13	\$859.75

Order Date	Order ID	Salesperson	Units	Order Amount
1/1/2015	10458	Bossk	13	\$1,400.00
2/1/2015	10459	Dengar	14	\$2,555.95
2/2/2015	10460	Dengar	7	\$344.00
3/1/2015	10461	Bossk	8	\$857.00
3/3/2015	10462	Greedo	9	\$111.55
3/5/2015	10463	Fett	18	\$823.00
4/1/2015	10464	Bossk	17	\$24,455.50
4/28/2015	10465	Fett	13	\$213.00
4/30/2015	10466	Fett	1	\$10.00
5/1/2015	10467	Bossk	11	\$789.70
5/2/2015	10468	Dengar	25	\$21,286.60
5/3/2015	10469	Fett	8	\$1,285.00
6/1/2015	10470	Bossk	11	\$201.00
7/1/2015	10471	Bossk	13	\$859.75

A solid black rectangle with no visible content or text.



A 12x5 grid of squares. The top row is highlighted in light gray, while the remaining 11 rows are black. This visualizes the first step in the construction of the Cantor set, where the middle fifth of each interval is removed.

Order Date	Order ID	Salesperson	Units	Order Amount
1/1/2015	10458	Bossk	13	\$1,400.00
2/1/2015	10459	Dengar	14	\$2,555.95
2/2/2015	10460	Dengar	7	\$344.00
3/1/2015	10461	Bossk	8	\$857.00
3/3/2015	10462	Greedo	9	\$111.55
3/5/2015	10463	Fett	18	\$823.00
4/1/2015	10464	Bossk	17	\$24,455.50
4/28/2015	10465	Fett	13	\$213.00
4/30/2015	10466	Fett	1	\$10.00
5/1/2015	10467	Bossk	11	\$789.70
5/2/2015	10468	Dengar	25	\$21,286.60
5/3/2015	10469	Fett	8	\$1,285.00
6/1/2015	10470	Bossk	11	\$201.00
7/1/2015	10471	Bossk	13	\$859.75

Réalisation

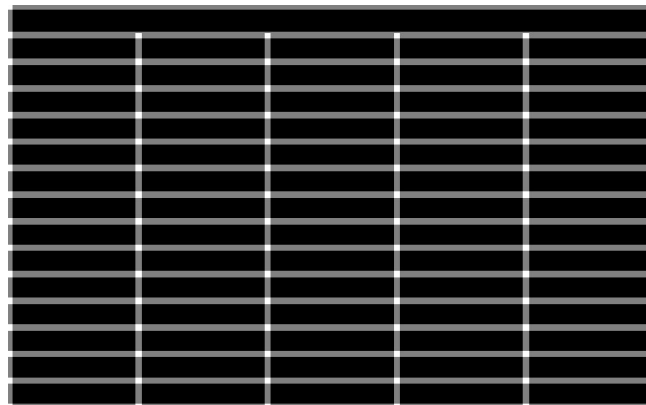
Architecture utilisé

Aperçu du projet



Tests

Order Date	Order ID	Salesperson	Units	Order Amount
1/1/2015	10458	Bossk	13	\$1,400.00
2/1/2015	10459	Dengar	14	\$2,555.95
2/2/2015	10460	Dengar	7	\$344.00
3/1/2015	10461	Bossk	8	\$857.00
3/3/2015	10462	Greedo	9	\$111.55
3/5/2015	10463	Fett	18	\$823.00
4/1/2015	10464	Bossk	17	\$24,455.50
4/28/2015	10465	Fett	13	\$213.00
4/30/2015	10466	Fett	1	\$10.00
5/1/2015	10467	Bossk	11	\$789.70
5/2/2015	10468	Dengar	25	\$21,286.60
5/3/2015	10469	Fett	8	\$1,285.00
6/1/2015	10470	Bossk	11	\$201.00
7/1/2015	10471	Bossk	13	\$859.75



Boxes detected

text extracted

[illegible]

Réalisation

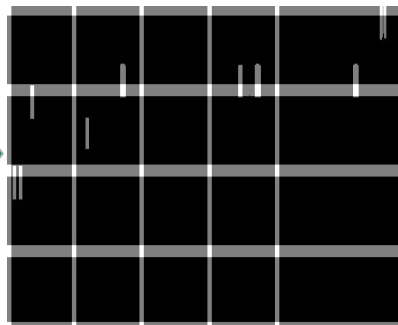
Architecture utilisé

Aperçu du projet



Tests

Size	Lethal to Operational Spacecraft	Number in Orbit	Trackable (i.e. can be cataloged)	Produces Lethal Fragments When Impacting An Operational Spacecraft
Small < 5mm	(Usually) Not	Millions	No	No
Medium 5mm - 10 cm	Usually	~ 500,000 in LEO	No	Maybe
Large > 10 cm	(Almost) Always	~ 21,000	Yes	Yes - 100s to 10,000s



Size	Lethal to Operational Spacecraft	Number in Orbit	Trackable (i.e. can be cataloged)	Produces Lethal Fragments When Impacting An Operational Spacecraft
Small < 5mm	(Usually) Not	Millions	No	No
Medium 5mm - 10 cm	Usually	1 500,000 in LEO	No	Maybe
Large > 10 cm	(Almost) Always	~ 21,000	Yes	Yes - 100s to 10,000s

Boxes detected

text extracted

```
finalboxes:
[[{'x': 10, 'y': 8, 'width': 140, 'height': 59}], [{'x': 158, 'y': 8, 'width': 146, 'height': 59}], [], [{'x': 312, 'y': 8, 'width': 147, 'height': 59}], [{'x': 467, 'y': 8, 'width': 146, 'height': 59}], [{'x': 621, 'y': 8, 'width': 279, 'height': 59}], [{'x': 10, 'y': 75, 'width': 140, 'height': 59}], [{'x': 158, 'y': 75, 'width': 146, 'height': 59}], [{'x': 180, 'y': 94, 'width': 8, 'height': 26}], [], [{'x': 312, 'y': 75, 'width': 147, 'height': 59}], [{'x': 467, 'y': 74, 'width': 146, 'height': 60}], [{'x': 621, 'y': 75, 'width': 279, 'height': 59}], [{'x': 10, 'y': 142, 'width': 140, 'height': 59}], [{'x': 158, 'y': 142, 'width': 146, 'height': 59}], [], [{'x': 312, 'y': 142, 'width': 147, 'height': 59}], [{'x': 467, 'y': 142, 'width': 146, 'height': 59}], [{'x': 621, 'y': 142, 'width': 279, 'height': 59}], [{'x': 10, 'y': 209, 'width': 140, 'height': 56}], [{'x': 158, 'y': 209, 'width': 146, 'height': 56}], [], [{'x': 312, 'y': 209, 'width': 147, 'height': 56}], [{'x': 467, 'y': 209, 'width': 146, 'height': 56}], [{'x': 621, 'y': 209, 'width': 279, 'height': 56}]]

extracted data:
[['Size', 'Lethal to Operational Spacecraft', 'Number in Orbit', 'Trackable (f.e. can be cataloged)', 'Produces Lethal Fragments When Impacting An Operational Spacecraft'], ['Small < 5mm', '(Usually) Not', 'Millions', 'No', 'No'], ['Medium 5mm - 10 cm', 'Usually', '1 500,000 in LEO', 'No', 'Maybe'], ['Large > 10 cm', '(Almost) Always', '~ 21,000', 'Yes', 'Yes - 100s to 10,000s']]
```

Réalisation

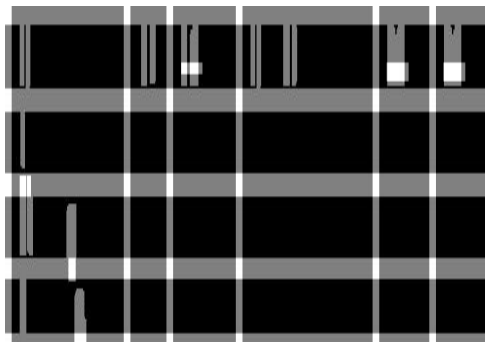
Architecture utilisé

Aperçu du projet



Tests

Statistic	N	Mean	St. Deviation	Min	Max
Complaints	30	64.63	12.73	40	85
Rating	30	65.84	13.35	37	90
Privileges	30	47.95	11.72	30	83



A	B	C	D	E	F	G
	0	1	2	3	4	5
0	Statistic	N	Mean	St. Deviation	Min	Max
1	Complaints	30	64.63	12.73	40	85
2	Rating	30	65.84	13.35	37	90
3	Privileges	30	47.95	11.72	30	83

Boxes detected

text extracted

```
finalboxes:
[[{'x': 10, 'y': 8, 'width': 23, 'height': 29}], [{'x': 37, 'y': 8, 'width': 143, 'height': 29}], [{'x': 188, 'y': 8, 'width': 56, 'height': 29}], [{'x': 273, 'y': 8, 'width': 10, 'height': 29}], [{'x': 252, 'y': 8, 'width': 96, 'height': 29}], [], [{"x": 356, "y": 8, "width": 23, "height": 29}], [{"x": 383, "y": 8, "width": 170, "height": 29}], [{"x": 583, "y": 8, "width": 8, "height": 29}], [{"x": 561, "y": 8, "width": 77, "height": 29}], [{"x": 668, "y": 8, "width": 8, "height": 6}], [{"x": 646, "y": 8, "width": 80, "height": 29}], [{"x": 10, "y": 45, "width": 10, "height": 28}], [{"x": 188, "y": 45, "width": 56, "height": 28}], [{"x": 252, "y": 45, "width": 96, "height": 28}], [{"x": 356, "y": 45, "width": 197, "height": 28}], [{"x": 646, "y": 45, "width": 77, "height": 28}], [{"x": 646, "y": 45, "width": 80, "height": 28}], [{"x": 10, "y": 81, "width": 170, "height": 28}], [{"x": 188, "y": 81, "width": 56, "height": 28}], [{"x": 252, "y": 81, "width": 96, "height": 28}], [{"x": 356, "y": 81, "width": 197, "height": 28}], [{"x": 646, "y": 81, "width": 77, "height": 28}], [{"x": 646, "y": 81, "width": 80, "height": 28}], [{"x": 10, "y": 116, "width": 14, "height": 25}], [{"x": 32, "y": 116, "width": 148, "height": 25}], [{"x": 188, "y": 116, "width": 56, "height": 25}], [{"x": 252, "y": 116, "width": 96, "height": 25}], [{"x": 356, "y": 116, "width": 197, "height": 25}], [{"x": 561, "y": 116, "width": 77, "height": 25}], [{"x": 646, "y": 116, "width": 80, "height": 25}], []]]
extracted data:
[['Statistic', 'N', 'Mean', 'St. Deviation', 'Min', 'Max'], ['Complaints', '30', '64.63', '12.73', '40', '85'], ['Rating', '30', '65.84', '13.35', '37', '90'], ['Privileges', '30', '47.95', '11.72', '30', '83']]
```


Réalisation

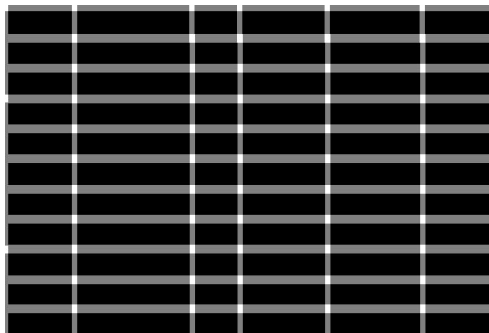
Architecture utilisé

Aperçu du projet



Tests

PPGENDER	PPEDUC	Count	Mean.FwBScore	Median.FwBScore	SD.FwBScore
Male	less than high school	169	47.02367	48	13.32647
Male	high school degree/GED	703	53.85064	53	14.63737
Male	some college/associate	1025	54.65366	55	13.61152
Male	bachelor's degree	790	59.32405	59	12.93089
Male	graduate/professional degree	665	62.36842	61	13.11986
Female	less than high school	260	48.79615	49	12.68633
Female	high school degree/GED	919	53.96409	54	14.61242
Female	some college/associate	908	53.88326	53	13.97973
Female	bachelor's degree	522	58.37931	59	12.97368
Female	graduate/professional degree	433	61.05543	61	12.93061



A	B	C	D	E	F	G
	0	1	2	3	4	5
0	PPGENDER	PPEDUC	Count	Mean.FwBScore	Median.FwBScore	SD.FwBScore
1	Male	less than high school	169	47.02367	48	13.32647
2	Male	high school degree/GED	703	53.85064	53	14.63737
3	Male	some college/associate	1025	54.65366	55	13.61152
4	Male	bachelor's degree	790	59.32405	59	12.93089
5	Male	graduate/professional degree	665	62.36842	61	13.11986
6	Female	less than high school	260	48.79615	49	12.68633
7	Female	high school degree/GED	919	53.96409	54	14.61242
8	Female	some college/associate	908	53.88326	53	13.97973
9	Female	bachelor's degree	522	58.37931	59	12.97368
10	Female	graduate/professional degree	433	61.05543	61	12.93061

text extracted

Boxes detected

```
idh': 133, 'height': 26]], [{x': 6, 'y': 283, 'width': 122, 'height': 27}], [{x': 136, 'y': 283, 'width': 214, 'height': 27}], [{x': 358, 'y': 283, 'width': 82, 'height': 27}], [{x': 448, 'y': 283, 'width': 157, 'height': 27}], [{x': 613, 'y': 283, 'width': 171, 'height': 27}], [{x': 792, 'y': 283, 'width': 133, 'height': 27}], [{x': 6, 'y': 318, 'width': 122, 'height': 25}], [{x': 136, 'y': 318, 'width': 214, 'height': 25}], [{x': 358, 'y': 318, 'width': 82, 'height': 25}], [{x': 448, 'y': 318, 'width': 157, 'height': 25}], [{x': 613, 'y': 318, 'width': 171, 'height': 25}], [{x': 792, 'y': 318, 'width': 133, 'height': 25}], [{x': 6, 'y': 351, 'width': 122, 'height': 28}], [{x': 136, 'y': 351, 'width': 214, 'height': 28}], [{x': 358, 'y': 351, 'width': 82, 'height': 28}], [{x': 448, 'y': 351, 'width': 157, 'height': 28}], [{x': 613, 'y': 351, 'width': 171, 'height': 28}], [{x': 792, 'y': 351, 'width': 133, 'height': 28}]]
extracted data:
[['PPGENDER', 'PPEDUC', 'Count', 'Mean.FwBScore', 'Median.FwBScore', 'SD.FwBScore'], ['Male', 'less than high school', '169', '47.02367', '48', '13.32647'], ['Male', 'high school degree/GED', '703', '53.85064', '53', '14.63737'], ['Male', 'some college/associate', '1025', '54.65366', '55', '13.61152'], ['Male', 'bachelor's degree', '790', '59.32405', '59', '12.93089'], ['Male', 'graduate/professional degree', '665', '62.36842', '61', '13.11986'], ['Female', 'less than high school', '260', '48.79615', '49', '12.68633'], ['Female', 'high school degree/GED', '919', '53.96409', '54', '14.61242'], ['Female', 'some college/associate', '908', '53.88326', '53', '13.97973'], ['Female', 'bachelor's degree', '522', '58.37931', '59', '12.97368'], ['Female', 'graduate/professional degree', '433', '61.05543', '61', '12.93061']]
```

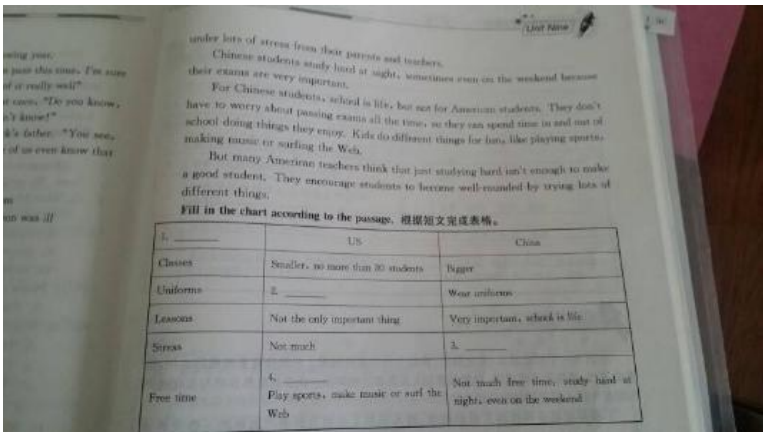

Realisation

Architecture utilisé

Aperçu du projet



Tests



Boxes detected

No text extracted

```
finalboxes:
[[{'x': 0, 'y': 0, 'width': 10, 'height': 97}], [{"x": 17, 'y': 0, 'width': 70, 'height': 15}, {'x': 15, 'y': 22, 'width': 71, 'height': 12}], [], [{"x": 94, 'y': 5, 'width': 114, 'height': 14}, {'x': 94, 'y': 23, 'width': 115, 'height': 14}], [{"x": 216, 'y': 8, 'width': 118, 'height': 11}, {'x': 217, 'y': 24, 'width': 118, 'height': 13}], [{"x": 341, 'y': 0, 'width': 9, 'height': 52}], [{"x": 14, 'y': 41, 'width': 72, 'height': 12}], [{"x": 93, 'y': 41, 'width': 40, 'height': 13}], [{"x": 218, 'y': 42, 'width': 120, 'height': 12}], [{"x": 12, 'y': 60, 'width': 73, 'height': 13}], [{"x": 92, 'y': 60, 'width': 119, 'height': 13}], [{"x": 218, 'y': 59, 'width': 123, 'height': 13}], [{"x": 11, 'y': 80, 'width': 74, 'height': 13}], [{"x": 92, 'y': 80, 'width': 120, 'height': 12}], [{"x": 219, 'y': 78, 'width': 124, 'height': 13}], [{"x": 6, 'y': 100, 'width': 78, 'height': 45}], [{"x": 90, 'y': 99, 'width': 125, 'height': 46}], [{"x": 220, 'y': 96, 'width': 130, 'height': 42}], [], []]]
extracted_data:
[[['', '', ''], ['', '', ''], ['', '', ''], ['', '', ''], ['', '', '2']]]
```



Conclusion & Perspective

Conclusion

Perspective

- Ce que nous avons appris de ce projet

Pour la version 2.0 :

- L'étude des tableaux plus complexes
- Réduire le temps de traitement

Merci pour votre attention