

Extraction d'un tableur depuis une image

Haithem ZEGGARI - Paul Minguet

Encadré par :

- M. PAGANI
- M. FÉRÉE
- M. DEGORRE



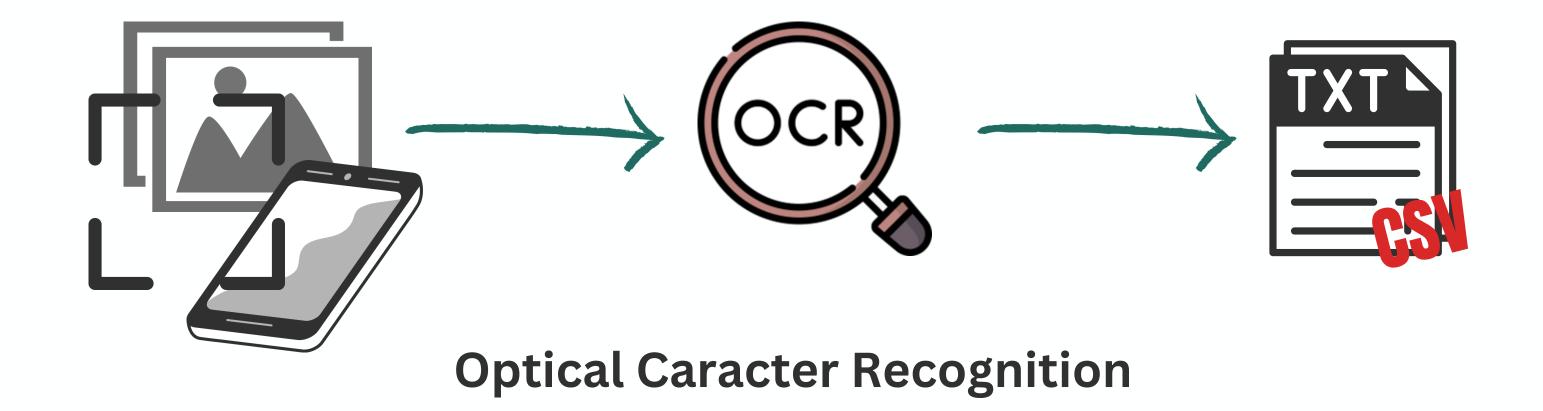


Sommaire

- I. Introduction et Problématique
 - 1.1. Objectifs
 - 1.2. Démonstration
- II. Aspects techniques
 - 2.1. Architecture et Gestion du projet
 - 2.2. Difficultés rencontrées
 - 2.3. Testabilité
 - 2.4. Programmation
- III. Conclusion

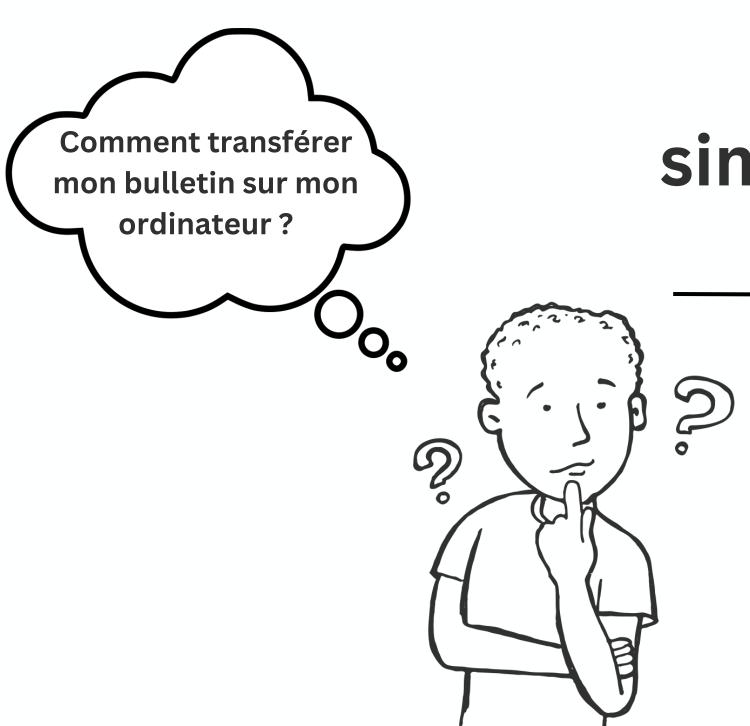


Introduction





Problématique



Comment numériser simplement un tableau sur papier?

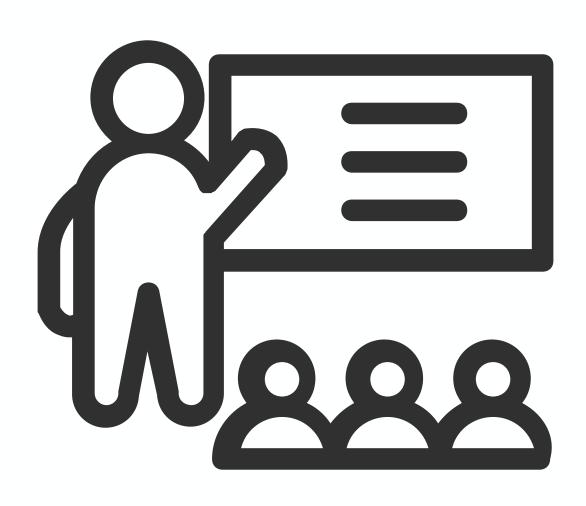


- Création d'un algorithme pour extraire un tableur depuis une image
- Créer un jeu de données conséquent
- Étudier les performances de l'algorithme





Démonstration



Extraction d'un tableur depuis une image



	Nov, 22	Dec, 22	Jan, 23	Feb, 23	Mar, 23	Apr, 23	May, 23
Projet Long							
Étude de l'existant]					
Expérimentation du logiciel d'OCR	L		_]			
Mise en place de l'algorithme et tests			L				
Étude des possibilités d'amélioration							

<u>Diagramme de Gantt du projet</u>



Architecture du projet

- ▼ 🛅 zeggari-minguet-plong-2022

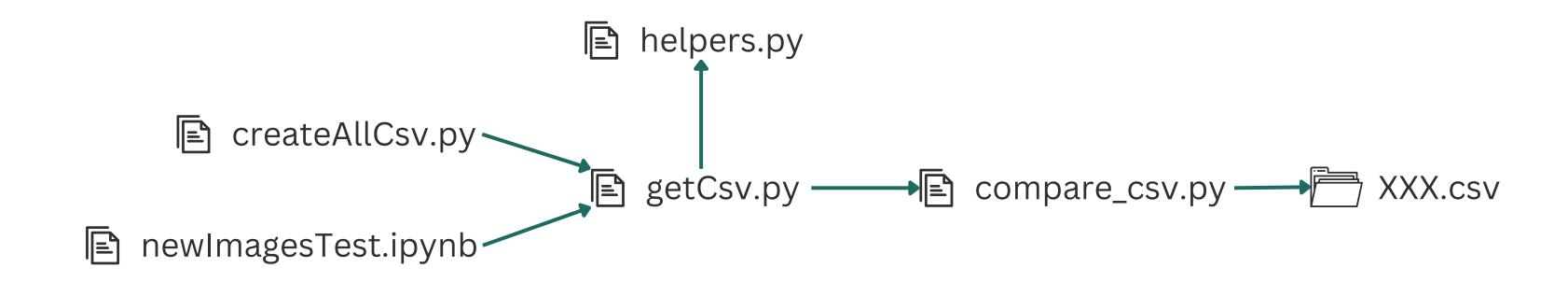
 - **▼** tesseract
 - **▼** image
 - easy-noColors
 - 01.csv
 - (a) 01.jpg
 - **□** 02.csv
 - © 02.jpg

 - hard-holes

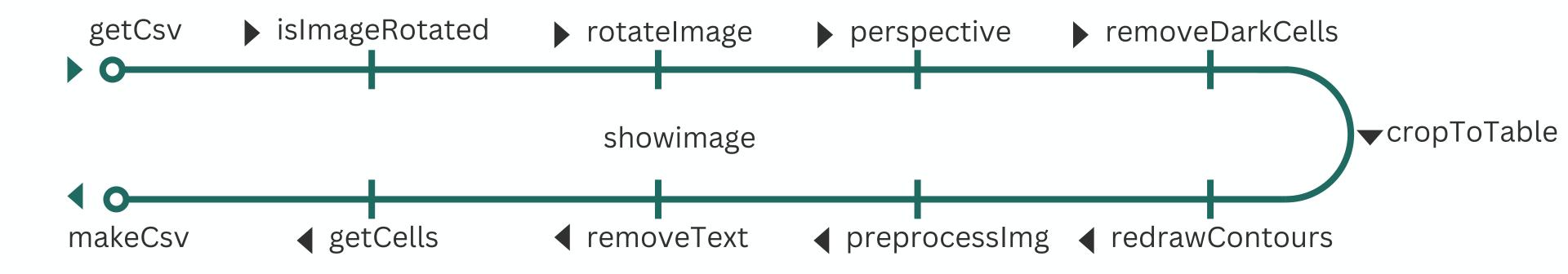
- e observings.md
- compare_csv.py
- createAllCsv.py
- getCsv.py
- helpers.py
- newImagesTest.ipynb
- tests.md
- usage.md
- journal.org
- README.md



Architecture du projet









Difficultés



Compréhension du fonctionnement de Tesseract OCR



Recherche des fonctionnalités d'OpenCV



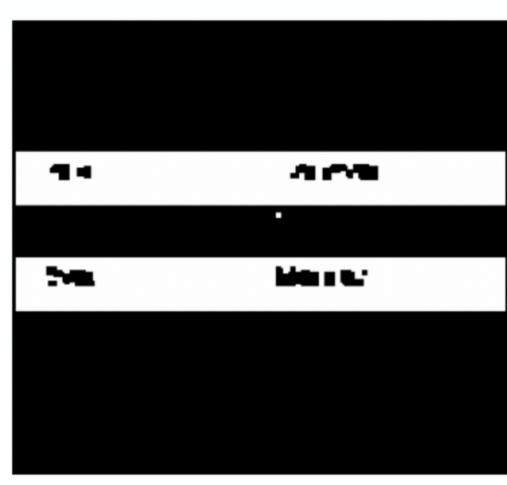
Trouver des solutions pour implémenter certaines fonctionnalités



Récupération des coordonnées des cellules

Nom	Prénom
Daniel	Alves
Fabio	Canavarro
Lio	Messi
Riyad	Mahrez
Cézar	Azbelicueta
Santos	Da Silva
Ali	Ahmed

Tableau à tester



<u>Image"morph"</u>

Nom	Prénom
Daniel	Alves
Fabio	Canavarro
Lio	Messi
Riyad	Mahrez
Cézar	Azbelicueta
Santos	Da Silva
Ali	Ahmed

<u>Cases inversées</u>



Récupération des coordonnées des cellules

Nom	Prénom
Daniel	Alves
Fabio	Canavarro
Lio	Messi
Riyad	Mahrez
Cézar	Azbelicueta
Santos	Da Silva
Ali	Ahmed

Contours recréés

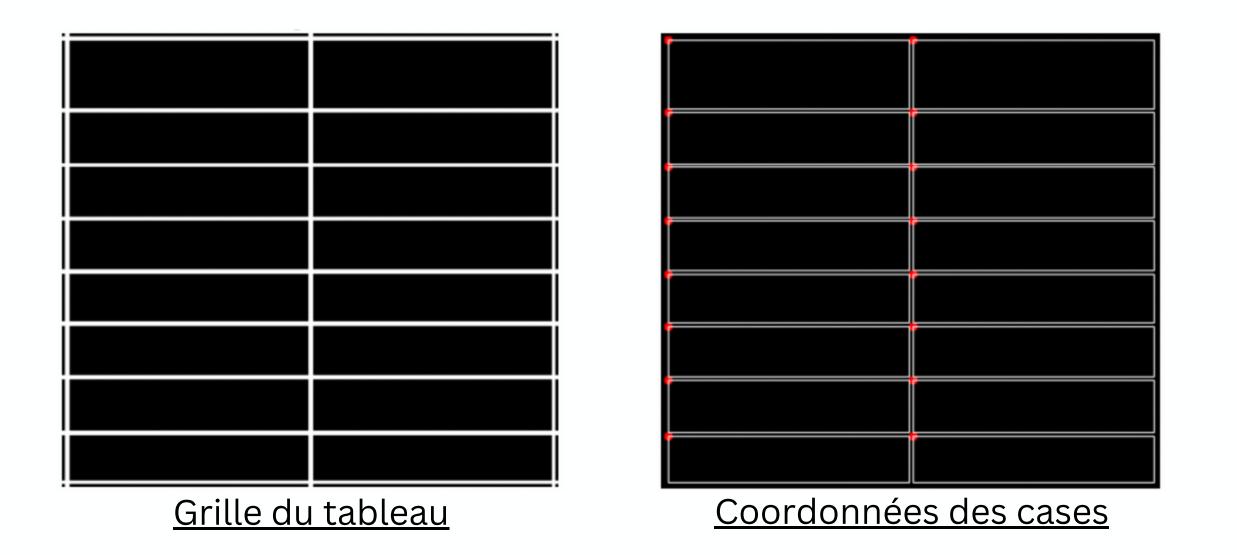
Nom	Prénom
Daniel	Alves
Fabio	Canavarro
Lio	Messi
Riyad	Mahrez
Cézar	Azbelicueta
Santos	Da Silva
Ali	Ahmed

<u>Image pretraitée</u>

Tableau sans le texte



Récupération des coordonnées des cellules





Testabilité

- Tests sur le jeu de données grâce à un script
- Tests à la main sur une image





NOM	PRENOM
Hench	John
Atencio	X
Bob	Gurr
Coats	Claude
Baxter	Tony
Blair	Mary
Rohde	Joe
Sherman	Richard
Sherman	Robert
Davis	Marc
Davis	Alice
Sklar	Marty

Pourcentage de similitude : 96.15%



Nom	Prénom	Adresse
Daniel	Alves	25 rue Albert eighnshtein 85632
Fabio	Canavarro	35 Rue Bauchaumont, 75002, Paris, France
Lio	Messi	Gare de l'est
Riyad	Mahrez	17 Bd de Strasbourg, 94200
Cézar	Azbelicueta	-
Santos	Da Silva	-
Ali	Ahmed	-

Pourcentage de similitude : 85.00%



Nom	Prénom	Phone	Email
Daniel	Alves	(258) 879-6320	Dani.alves@gmail.com
Fabio	Canavarro	(102) 787-0023	Fabicanav@gmail.com
Lio	Messi	(222) 589-9623	Lio122-m@gmail.com
Riyad	Mahrez	(213) 751-8794	Mahrezr26@gmail.dz
Cézar	Azbelicueta	(102) 259-6324	C_azb@gmail.com
Santos	Da Silva	(569) 032-7474	santosbr@gmail.com
Ali	Ahmed	(256) 254-0121	aliahled@gmail.com

Pourcentage de similitude : 100%



Nom	Prénom
Daniel	Alves
Fabio	Canavarro
Lio	Messi
Riyad	Mahrez
Cézar	Azbelicueta
Santos	Da Silva
Ali	Ahmed

• "Mahrez" -> ""

Pourcentage de similitude : 93.75%



Nom	Prénom
Daniel	Alves
Fabio	Canavarro
Lio	Messi
Riyad	Mahrez
Cézar	Azbelicueta
Santos	Da Silva
Ali	Ahmed

Pourcentage de similitude : 100%



Nom	Prénom
Daniel	
	Canavarro
Lio	
	Mahrez
Cézar	
	Da Silva
Ali	Ahmed

Pourcentage de similitude : 100%



NOM	PRENOM
Hench	John
	X
Bob	20
Coats	
Baxter	Tony
	18
Rohde	Joe
	Richard
	Robert
Davis	
Davis	Alice
Sklar	Marty

Pourcentage de similitude : 96.15%



NOM	PRENOM	EMAIL
Hench	John	
Atencio	X	
		bgurr@wed.com
Coats	Claude	
Baxter	Tony	tbaxter@wed.com
	Mary	
Rohde	Joe	jrohde@wed.com
Sherman	Richard	dsherman@wed.com
Sherman		rsherman@wed.com
	Marc	mdavis@wed.com
Davis	Alice	
Sklar		lsklar@wed.com

- "X" -> ""
- "@wed" -> "@vwed"
- "Tony" -> "on"
- "Marty" -> "Ma rty"

Pourcentage de similitude : 50.00%



NOM	PRENOM	PHONE
Hench	50 50	36
Atencio	X	(839) 589-6673
	72 72 = 11 = 11 = 12	(567) 426-7919
Coats	Claude	(708) 343-5400
Baxter	Tony	(783) 719-1744
Blair	68 - Taran	(666) 449-1166
	Joe	
Sherman		(410) 431-2226
Sherman	Robert	
Davis		(630) 886-1584
	Alice	(654) 348-9009
Sklar	Marty	

- "X" -> ""
- "Tony" -> "on"
- "Marty" -> "Ma rty"

Pourcentage de similitude : 92.30%



Nom	Prénom
Daniel	Alves
Fabio	Canavarro
Lio	Messi
Riyad	Mahrez
Cézar	Azbelicueta
Santos	Da Silva
Ali	Ahmed

• Le texte clair n'est pas pris en compte

Pourcentage de similitude : 56.25%



Nom	Prénom	Phone	Email
Daniel	Alves	(258) 879-6320	Dani.alves@gmail.com
Fabio	Canavarro	(102) 787-0023	Fabicanav@gmail.com
Lio	Messi	(222) 589-9623	Lio122-m@gmail.com
Riyad	Mahrez	(213) 751-8794	Mahrezr26@gmail.dz
Cézar	Azbelicueta	(102) 259-6324	C_azb@gmail.com
Santos	Da Silva	(569) 032-7474	santosbr@gmail.com
Ali	Ahmed	(256) 254-0121	aliahled@gmail.com

• Le texte clair n'est pas pris en compte

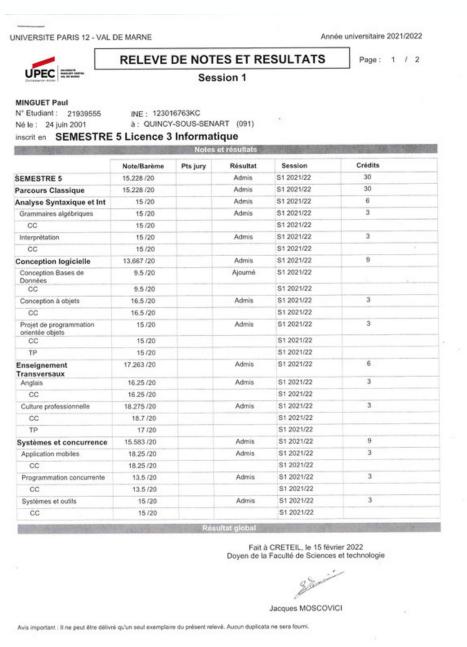
Pourcentage de similitude : 2.50%



NOM	PRENOM
Hench	John
Baxter	Tony
Blair	Mary
Sherman	19
Sherman	10 ²
Davis	Marc
Davis	
Sklar	Marty

• Le texte clair n'est pas pris en compte

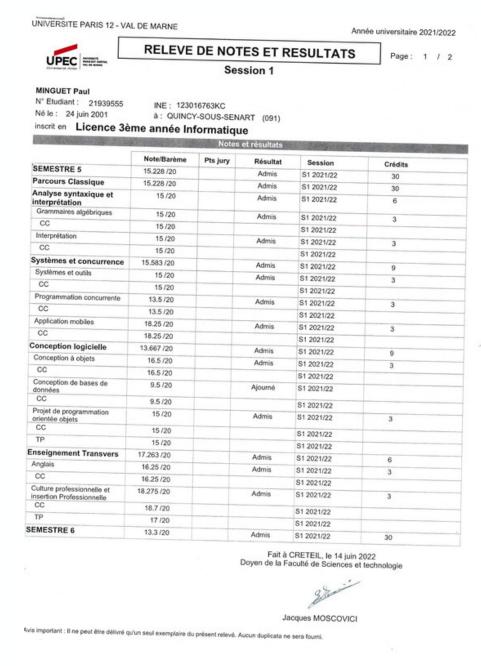
Pourcentage de similitude : 69.24%



Différences:

- "S" -> "\$"
- Bas de l'image rogné

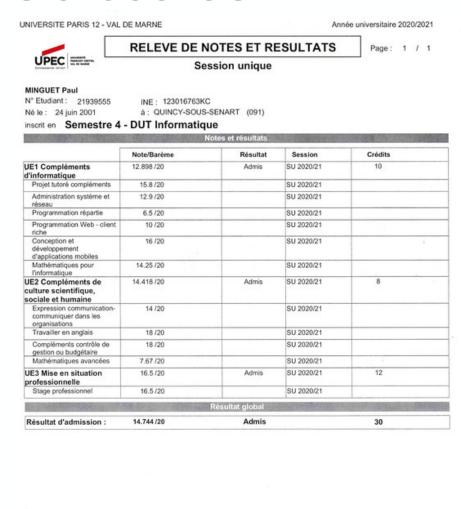
Pourcentage de similitude : 77.22%



Différences:

- "S" -> "\$"
- "CC" -> ""
- "TP" -> ""
- Bas de l'image rogné

Pourcentage de similitude : 75.00%





Différences:

- "tutoré" -> "tutoreé"
- Quelques erreurs de l'OCR

Pourcentage de similitude : 86.67% (81% rotation)





• Quelques erreurs de l'OCR

Pourcentage de similitude : 76.52%





الجمهوريسة الجزائرية الديمقراطية الشعبية République Algérienne Démocratique et Populaire وزارة التعليم العالي و البحث العلمى

المدرسة الوطنية العليا للإعلام الآلى Ecole nationale Supérieure d'Informatique مدبرية الدراسات

Ministère de l'Enseignement Supérieur et de la Recherche Scientifique

RELEVE DE NOTES

Matricule: 16/0116 Année d'étude : 1ère année - Second Cycle Année universitaire : 2018/2019

Prénom : HAITHEM Date de naissance : 25/04/1998

Diplôme préparé : Ingénieur d'état en Informatique

à : CONSTANTINE

RES1	Réseaux 1	4	10.65	Février
5Y51	Systèmes d'exploitation 1	5	5.50	Février
IGL	Introduction au génie logiciel	5	10.00	Février
THP	Théorie des langages de programmation et applications	- 4	11.90	Février
ANUM	Analyse Numérique	4	11.50	Février
ORG	Analyse des organisations	3	14.30	Février
RO	Recherche Opérationnelle: graphes et algorithmes	3	7.20	Février
LANG1	Langue anglaise 1	2	12.17	Février
ARCH	Architecture	4	13.33	Juin
RES2	Réseaux 2	3	12.11	Juin
SYS2	Système d'exploitation 2	4	8.35	Juin
BDD	Bases de données	5	12.00	Juin
MCSI	Méthodologies d'analyse et conception de systèmes d'information	5	12.80	Juin
CPROJ	Conduite de projets	3	10.88	Juin
PROJ	Projet	3	15.51	Juin
SEC	Introduction à la sécurité informatique	1	15.00	Juin
LANG2	Langue anglaise 2	2	14.50	Juin

Moyenne Annuelle: 11.21/20

Fait le: 25/01/2021 Décision du Conseil : Admis

لا بمنح الا نسخة وأحده من هذه الوثيقة .Il n'est délivré qu'un seul exemplaire de ce document

Ecole nationale Supérieure d'informatique بالأعلام بالألم الألم الأطراف المال الإعلام الألم (الأم 1620) BPM68 16270, Oued Smar, Alger, Tél : 023939132 ; Fax : 023939142 ; http://www.esi.dz

Différences:

• Quelques erreurs de l'OCR

Pourcentage de similitude : 92.43%



Caractéristiques de l'image	Pourcentage moyen	Principales erreurs
Simple, texte avec symboles ([,(;:)])	88.08%	
Texte en couleur	96.25%	Certains caractères non / mal reconnus
Trous dans le tableau	92.66%	mon / macreconnas
Trous et couleurs	84.56%	
Couleurs claires	56.04%	Certaines couleurs disparaissent au pré-traitement



Trous et couleurs claires	70.52%	Certaines couleurs disparaissent au pré-traitement
Images scannées	88.04%	Certains caractères non / mal reconnus et certaines images coupées
Tableaux sans bordure	0.00%	
Tableaux tournés	Idems que non tournés	



Conclusion

Ce qu'on a appris :

- Pré-traitement des images
- Techniques de recadrage d'images
- Implémentation d'un logiciel d'OCR dans un script

Version 2.0:

- Amélioration du pré traitement
- Images sans bordures / grille

Pour recommencer:

- Plus de tests sur le pré traitement des images
- Plus de recherche sur la technologie



Merci pour votre attention



Avez-vous des questions?